

MSc Scientific Computing Dissertation

ARM Cluster Linpack Benchmarks

John Duffy

August 2020

1 Introduction

<https://github.com/johnduffymsc/picluster>

1.1 Aims

1.1.1 Investigate Maximum Achievable Linpack Performance

Efficiency... achieved vs theoretical maximum

1.1.2 Investigate Gflops/Watt

Green500 ranking...

1.1.3 Overview of Competitive Available Gflops/£

Buy lots of Pi's, or buy a bigger machine...

Plot Gflops vs £...

1.2 Typography

This is a computer name...

node1

This is a command to type...

\$ grep

This is a command output displayed on your screen...

This is a file listing...

Listing 1: /etc/hosts

```
1 ##
2 # Host Database
3 #
4 # localhost is used to configure the loopback interface
5 # when the system is booting. Do not change this entry.
6 ##
7 127.0.0.1 localhost
8 255.255.255.255 broadcasthost
9 ::1          localhost
10 192.168.0.1 node1
11 192.168.0.2 node2
12 192.168.0.3 node3
13 192.168.0.4 node4
14 192.168.0.5 node5
15 192.168.0.6 node6
16 192.168.0.7 node7
17 192.168.0.8 node8
18 192.168.0.9 node9
```

2 Raspberry Pi 4 Model B

2.1 Description

Photo...

Description...

Highlights...

Limitations...

Reference data sheet in Appendix...

2.2 Theoretical Maximum Performance (Gflop/s)

The Raspberry Pi 4 Model B uses the Broadcom BCM2711 System on a Chip (Soc).

Block diagram from Cortex-A72 Software Optimisation Guide

4 cores

1.5 GHz

128 bit SIMD

4 GB memory (our chosen model)

Caches...

Pipeline...

Simplistically, ...

This ignores instructions pipelining benefits...

3 Pi Cluster

Photo...

Description...

Ubuntu 20.04 LTS 64-bit Preinstalled Server...

Reference Appendix A for detailed build instructions...

Limitations...

Software/update management...

Next PXE/NFS boot...

Cluster management tools

BLAS libraries...

BLAS library management... `update-alternatives --config libblas.so.3-aarch64-linux-gnu`

picluster/tools... appendix ?... use from node1...

4 High-Performance Linpack (HPL) Benchmark

Reference Paper...

[https://www.netlib.org/benchmark/hpl/...](https://www.netlib.org/benchmark/hpl/)

Describe algorithm...

Terminology R_{peak} , R_{max} ..., problem size...

Describe methodology for determining main parameters NB, N, P and Q...

N formula...

Reference <http://hpl-calculator.sourceforge.net>

4.1 Building and Installing HPL

See Appendix...

4.2 HPL.dat

Describe HPL.dat parameters...

Listing 2: Example HPL.dat

```
1 HPLinpack benchmark input file
2 Innovative Computing Laboratory, University of Tennessee
3 HPL.out          output file name (if any)
4 0                device out (6=stdout,7=stderr,file)
5 1                # of problems sizes (N)
6 26208            Ns
7 1                # of NBs
8 32               NBs
9 0                PMAP process mapping (0=Row-,1=Column-major)
10 2               # of process grids (P x Q)
11 1 2             Ps
12 8 4             Qs
13 16.0            threshold
14 3               # of panel fact
```

```

15 0 1 2      PFACTs (0=left, 1=Crout, 2=Right)
16 2          # of recursive stopping criterium
17 2 4        NBMINs (>= 1)
18 1          # of panels in recursion
19 2          NDIVs
20 3          # of recursive panel fact.
21 0 1 2      RFACTs (0=left, 1=Crout, 2=Right)
22 1          # of broadcast
23 0          BCASTs (0=1rg,1=1rM,2=2rg,3=2rM,4=Lng,5=LnM)
24 1          # of lookahead depth
25 0          DEPTHS (>=0)
26 2          SWAP (0=bin-exch,1=long,2=mix)
27 64         swapping threshold
28 0          L1 in (0=transposed,1=no-transposed) form
29 0          U in (0=transposed,1=no-transposed) form
30 1          Equilibration (0=no,1=yes)
31 8          memory alignment in double (> 0)

```

A detailed description of each line of this file is ...

4.3 HPL.out

Describe HPL.out...

It is very easy to use **grep** to find the lines in HPL.out containing the results. And to then conduct a general numeric sort, first by P and then by Gflops, to find Rmax for each P and Q pair, squeezing repeated white space down to a single space for readability.

```
grep WR HPL.out | sort -g -k 4 -k 7 | tr -s ' ' > HPL.out.sorted
```

Listing 3: Example HPL.out.sorted

```

1 WR00C2R2 26208 32 1 8 802.01 1.4965e+01
2 WR00R2C2 26208 32 1 8 799.75 1.5007e+01
3 WR00L2L2 26208 32 1 8 796.04 1.5077e+01
4 WR00C2C2 26208 32 1 8 794.65 1.5103e+01
5 WR00L2C2 26208 32 1 8 793.86 1.5118e+01
6 WR00C2L2 26208 32 1 8 793.67 1.5122e+01
7 WR00R2L2 26208 32 1 8 793.48 1.5126e+01
8 WR00R2R2 26208 32 1 8 790.26 1.5187e+01
9 WR00L2R2 26208 32 1 8 789.16 1.5208e+01
10 WR00R2L4 26208 32 1 8 774.49 1.5497e+01

```

11	WR00C2R4	26208	32	1	8	773.52	1.5516e+01
12	WR00L2L4	26208	32	1	8	770.20	1.5583e+01
13	WR00R2C4	26208	32	1	8	767.92	1.5629e+01
14	WR00L2C4	26208	32	1	8	763.10	1.5728e+01
15	WR00L2R4	26208	32	1	8	762.43	1.5742e+01
16	WR00R2R4	26208	32	1	8	761.92	1.5752e+01
17	WR00C2C4	26208	32	1	8	761.58	1.5759e+01
18	WR00C2L4	26208	32	1	8	757.87	1.5836e+01
19	WR00R2R2	26208	32	2	4	728.78	1.6468e+01
20	WR00R2C2	26208	32	2	4	728.21	1.6481e+01
21	WR00R2L2	26208	32	2	4	726.55	1.6519e+01
22	WR00C2R2	26208	32	2	4	722.38	1.6614e+01
23	WR00L2C2	26208	32	2	4	721.63	1.6632e+01
24	WR00L2L2	26208	32	2	4	721.54	1.6634e+01
25	WR00C2C2	26208	32	2	4	721.25	1.6640e+01
26	WR00C2L2	26208	32	2	4	720.82	1.6650e+01
27	WR00L2R2	26208	32	2	4	720.80	1.6651e+01
28	WR00L2R4	26208	32	2	4	692.09	1.7341e+01
29	WR00R2C4	26208	32	2	4	690.37	1.7385e+01
30	WR00C2L4	26208	32	2	4	686.69	1.7478e+01
31	WR00C2C4	26208	32	2	4	686.23	1.7489e+01
32	WR00C2R4	26208	32	2	4	686.08	1.7493e+01
33	WR00L2L4	26208	32	2	4	686.02	1.7495e+01
34	WR00L2C4	26208	32	2	4	685.88	1.7498e+01
35	WR00R2L4	26208	32	2	4	685.76	1.7502e+01
36	WR00R2R4	26208	32	2	4	684.45	1.7535e+01

4.4 Running xhpl

To run xhpl using the serial version of OpenBLAS...

```
/piccluster/tools/piccluster-set-libblas-openblas-serial
```

```
cd /piccluster/hpl/hpl-2.3/bin/serial
```

```
mpirun -host node1:4 -np 4 xhpl
```

To run xhpl using the OpenMP version of OpenBLAS...

```
/piccluster/tools/piccluster-set-libblas-openblas-openmp
```

```
export OMPNUMTHREADS=4 TODO: SET GLOBALLY AT INSTALLATION
```

```
export BLISNUMTHREADS=4 TODO: SET GLOBALLY AT INSTALLATION
```

TION

```
cd /picluster/hpl/hpl-2.3/bin/serial
```

5 OpenMPI

What is OpenMPI...

6 OpenMP

What is OpenMP...

7 OpenMPI Baseline Benchmarks

Ubuntu 20.04 LTS 64-bit packages, without any tweaks...

1 core... a single ARM Cortex-A72 core...

1 node... a single Raspberry Pi 4 Model B, 4 x ARM Cortex-A72 cores...

Linpack performance scales with problem size... [REFERENCE](#)

80% of memory a good initial guess... [FAQ](#) [REFERENCE](#)...

Methodology...

1 core... to investigate single core performance... caveats... use 1GB of memory...

1 node... to investigate inter-core performance...

2 nodes... to investigate inter-core and inter-node performance...

1..8 nodes ... to investigate over scaling of performance with node count... with optimal N, NB, P and Q parameters determined from 2 node investigation... caveats...

7.1 OpenBLAS

7.2 BLIS

7.3 ARM Performance Libraries

THIS DOES NOT WORK YET, BUT WILL ON THE NEXT RELEASE OF ARM PERFORMANCE LIBRARIES. KEEP FOR FUTURE REFERENCE.

”Arm Performance Libraries provides optimized standard core math libraries for high-performance computing applications on Arm processors. This free version of the libraries provides optimized libraries for Arm® Neoverse™ N1-based Armv8 AArch64 implementations that are compatible with various versions of GCC. You do not require a license for this version of the libraries.”

For clarity the EULA is included as Appendix ?

Downloaded Arm Performance Libraries 20.2.0 with GCC 9.3 for Ubuntu 16.04+...

Note... need to install environment-modules package...

```
1 $ ssh node1
2 $ sudo apt install environment-modules
3 $ mkdir picluster/armpl
4 $ cd picluster/armpl
5 $ tar xvf arm-performance-libraries_20.2_Ubuntu-16.04_gcc-9.3.tar
6 $ rm arm-performance-libraries_20.2_Ubuntu-16.04_gcc-9.3.tar
7 $ sudo ./arm-performance-libraries_20.2_Ubuntu-16.04.sh
```

The default installation directory is /opt/arm...

Compile HPL with armpl...

```
1 $ cd ~/picluster/hpl/hpl-2.3
2 $ cp Make.serial Make.armpl-serial
```

Edit Make.armpl-serial...

```
1 # -----
2 # - Linear Algebra library (BLAS or VSIP) -----
3 # -----
4 # LAinc tells the C compiler where to find the Linear Algebra
  library
5 # header files, LAlib is defined to be the name of
  the library to be
```



```

6 # used. The variable LAdir is only used for defining LAinc and LAlib.
7 #
8 LAdir      = /opt/arm/armpl_20.2_gcc-9.3
9 LAinc      =
10 LAlib      = -L$(LAdir)/lib -larmpl -lgfortran -lamath -lm

```

Compile HPL...

```
1 $ make arch=armpl-serial
```

7.4 1 Core Baseline Benchmark

Problem size restricted to 80% of 1GB...

NB 32 to 256 in increments of 8...

NB	N	NB	N	NB	N	NB	N	NB	N
32	9248	80	9200	128	9216	176	9152	224	9184
40	9240	88	9240	136	9248	184	9200	232	9048
48	9264	96	9216	144	9216	192	9216	240	9120
56	9240	104	9256	152	9120	200	9200	248	9176
64	9216	112	9184	160	9120	208	9152	256	9216
72	9216	120	9240	168	9240	216	9072	-	-

1x1

```
$ mpirun -np 1 xhpl
```

mpirun does bind to core by default for $np \leq 2$

4 x 4.7527e+00 = 19 Gflops

Explain...

Cache misses from peak...

A single core is capable of achieving maximum theoretical performance... CAVEATS
whole L2 cache, whole node 4 GB memory, although problem size limited to
80% of 1 GB...

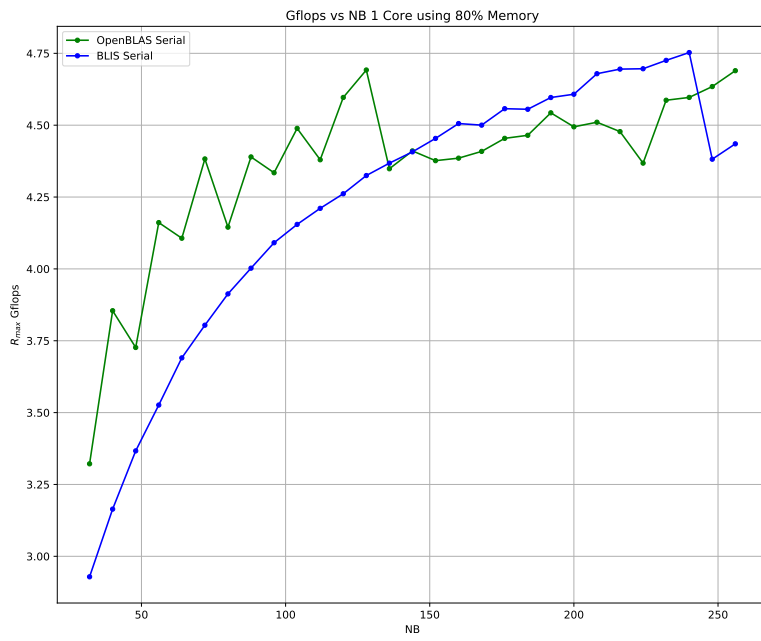


Figure 1: 1 Core R_{max} vs NB using 80% of 1GB memory.

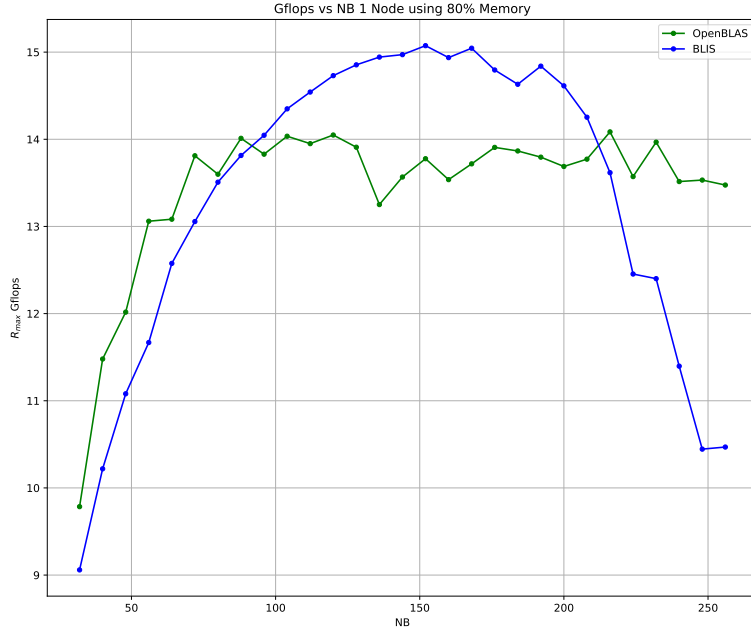


Figure 2: 1 Node R_{max} vs NB using 80% of 4GB memory.

7.5 1 Node Benchmark

1x4

NB	N	NB	N	NB	N	NB	N	NB	N
32	18528	80	18480	128	18432	176	18480	224	18368
40	18520	88	18480	136	18496	184	18400	232	18328
48	18528	96	18528	144	18432	192	18432	240	18480
56	18536	104	18512	152	18392	200	18400	248	18352
64	18496	112	18480	160	18400	208	18512	256	18432
72	18504	120	18480	168	18480	216	18360	-	-

```
1 mpirun -np 4 xhpl
```

mpirun does bind to socket by default for $np \geq 2$

7.6 2 Node Baseline Benchmark

P1 x Q8

P2 x Q4

NB	N	NB	N	NB	N	NB	N	NB	N
32	26208	80	26160	128	26112	176	26048	224	26208
40	26200	88	26136	136	26112	184	26128	232	25984
48	26208	96	26208	144	26208	192	26112	240	26160
56	26208	104	26208	152	26144	200	26200	248	26040
64	26176	112	26208	160	26080	208	26208	256	26112
72	26208	120	26160	168	26208	216	26136	-	-

7.7 8 Node Baseline Benchmark

1x32 2x16 4x8

7.8 Observations

Best NB...

PxQ discussion... 1x8 vs 2x4... ethernet comment...

Iperf...

htop...

top...

perf...

cache misses...

software interrupts...

Suggests... improve network efficiency?

7.9 Baseline Benchmark

As per software installation from Ubuntu Server 20.04 LTS.

OpenBLAS

OpenMP

OpenMPI

HPL-2.3 compiled with Make.rpi4.baseline

NB = 128 (mid-range guess; to be tuned based on L1 cache size)

N for 80% efficiency (from tool)

Recommended: $P \times Q$ as square as possible, with $Q > P$

4 cores per node 1.5 GHz clock speed 4 GB memory per node 2 instructions per cycle (estimated as NEON is 128 bits)

From tool:

Run:

```
mpirun -host node1:4 -np 4 xhpl
mpirun -host node1:4,node2:4 -np 8 xhpl
mpirun -host node1:4,node2:4,node3:4 -np 12 xhpl
etc
```

TABLE RESULTS - Gflops vs node count, time vs node count GRAPH RESULTS - Gflops vs node count, time vs node count

Discussion...

NB size...

P

Q Ratio...

Node number scaling...

7.10 OpenMPI without OpenMP

Describe processor grid layout.

hosts-with-slots file.

7.11 OpenMPI with OpenMP

Describe processor grid layout.

hosts-no-slots file.

8 Performance Optimisation

8.1 Methodology

1. Measure
2. Study results and propose theory
3. Change something based on 2.
4. Measure
5. Repeat steps 1 - 4

9 Build Kernel with Jumbo Frames Support

Standard MTU is 1500 bytes...

Maximum payload size is 1472 bytes...

NB of 184 (x 8 bytes for Double Precision) = 1472 bytes...

NB > 184 => packet fragmentation => reduced network efficiency...

This causes drop of in performance???...

Max MTU on Raspberry Pi 4 Model B is set at build time to 1500...

Not configurable above 1500...

TODO: EXAMPLE OF ERROR MSG...

Need to build the kernel with higher MTU...

Make source packages available...

```

sudo touch /etc/apt/sources.list.d/picluster.list
sudo vim /etc/apt/sources.list.d/picluster.list...
    deb-src http://archive.ubuntu.com/ubuntu focal main
    deb-src http://archive.ubuntu.com/ubuntu focal-updates main
sudo apt update

```

Create a kernel build directory with the correct access permissions to prevent source download warnings.

```

mkdir kernel
sudo chown _apt:root kernel
cd kernel

```

Install the kernel build dependencies...

```

sudo apt-get build-dep linux linux-image-$(uname -r)

```

Download the kernel source...

```

sudo apt-get source linux-image-$(uname -r)

```

Make the required changes to the source... as per REFERENCE

```

cd linux-raspi-5.4.0

sudo vim include/linux/if_vlan.h...
    #define VLAN_ETH_DATA_LEN    9000
    #define VLAN_ETH_FRAME_LEN  9018

sudo vim include/uapi/linux/if_ether.h...
    #define ETH_DATA_LEN        9000
    #define ETH_FRAME_LEN       9014

sudo vim drivers/net/ethernet/broadcom/genet/bcmgenet.c...
    #define RX_BUF_LENGTH      10240

```

Add a Jumbo Frames identifier, "+jf", to the new kernel name...

```

sudo vim debian.raspi/changelog...
    linux (5.4.0-1013.13+jf) focal; urgency=medium

```

Build the kernel...

```
sudo LANG=C fakeroot debian/rules clean
sudo LANG=C fakeroot debian/rules binary
```

Install the new kernel...

```
sudo sudo dpkg -i linux*5.4?????????.deb
```

10 Single Core Optimisation

10.1 Block Size Optimisation

The block size, NB tuning parameter, is used for matrix calculations and also for network transport.

The most efficient block size is related to the L1 cache size. Describe...

11 Single Node Optimisation

12 Cluster Optimisation

12.1 Recompile HPL for Cortex-A72

Block size!

RESULTS

12.2 Recompile OpenBLAS with OpenMP Support

As advised by the DebianScience/LinearAlgebraLibraries, OpenBLAS should be recompiled from source for best performance.??? See conflicting statement below.

The Ubuntu/Debian OpenBLAS package is an ARM64 multi-architecture build of OpenBLAS which includes the ARM Cortex-A72. As stated in the README.Debian, performance improvements will be minimal by compiling from source.

However, Ubuntu/Debian build does not include support for OpenMP. So, to test the combination of OpenMPI/OpenMP it is necessary to recompile OpenBLAS from source.

To download the same version of OpenBLAS as Ubuntu 20.04 Server LTS (v0.3.8):

```
cd ~/phas0077/downloads
wget https://github.com/xianyi/OpenBLAS/archive/v0.3.8.zip -O OpenBLAS-0.3.8.zip
cp OpenBLAS-0.3.8.zip ~/phas0077/projects
cd ~/phas0077/projects
unzip OpenBLAS-0.3.8.zip
rm OpenBLAS-0.3.8.zip
cd OpenBLAS-0.3.8
cp Makefile.rule Makefile.rule.original
vim Makefile.rule
make
make install
make clean
```

PROCESSOR AFFINITY - include later.

Need to edit 'Makefile.rule'.

The file 'Makefile.arm64' already has the following so there is no need to add specific architecture flags to 'Makefile.rule'.

```
ifeq ($(CORE), CORTEXA72)
COMMON_OPT += -march=armv8-a -mtune=cortex-a72
FCOMMON_OPT += -march=armv8-a -mtune=cortex-a72
endif
```

Block size!

TODO: Read DebianScience howto recompile/package

TODO: Check ARM gcc options; -mune, -march for ARM Cortex-A72 USE LD_PRELOAD trick to use recompiled libopenblas.so for testing before packaging as a Debian package. WHAT ABOUT THE OTHER NODES - How do they know about the re-compiled package on node1?

LD_PRELOAD for testing on node1.

ONCE installed as Debian package, prevent it being 'updated/upgraded'.

RESULTS

12.3 Recompile OpenMPI for Cortex-A72

Block size!

TODO: Check ARM gcc options; -mune, -march for ARM Cortex-A72

RESULTS

12.4 Network Tuning

Is it possible to improve performance looking at network parameter? MTU?

TODO: Read OpenMPI docs

RESULTS

12.5 Perf/Dtrace Linpack

See if there are any performance bottlenecks

New package OpenBLASRPi4? for RPi4 specific optimisations?

RESULTS

13 Appendix A - Raspberry Pi Cluster Build

13.1 Introduction

This appendix is intended to be a complete and self contained guide for building a Raspberry Pi Cluster. With the caveat that the cluster has the bare minimum software/functionality necessary to compile and run the High Performance Linpack (HPL) benchmark, namely the build-essential package, two BLAS libraries (OpenBLAS and BLIS), and Open-MPI. A number of performance measurement tools are also installed, such as perf and iperf. The latest version of HPL is downloaded and built from source.

It would be a relatively simple task to add... SLIRM or...

The cluster consists of the following components...

8 x Raspberry Pi 4 Model B 4GB compute nodes, node1 to node8
1 x software development and build node, node9
9 x Official Raspberry Pi 4 Model B power supplies
9 x 32GB Class 10 MicroSD cards
1 x *workstation*, in my case my MacBook Pro,
macbook
1 x 8 port Gigabit Router/Firewall
1 x 16 port Gigabit switch with Jumbo Frame support

Items

Photo

13.2 Preliminary Tasks

1. Update the EE-PROM
2. Get MAC address
3. Generate keys
4. Amend macbook /etc/hosts file...

13.2.1 Update Raspberry Pi EE-PROMs

13.2.2 Get Raspberry Pi MAC Addresses

13.2.3 Generate User Key Pair

On macbook (no passphrase):

```
1 $ ssh-genkey -t rsa -C john
```

This will create two files... in ...

13.2.4 Amend macbook /etc/hosts

On macbook, using your favourite editor, add the following to /etc/hosts:

```
1 192.168.0.1 node1
2 192.168.0.2 node2
3 192.168.0.3 node3
4 192.168.0.4 node4
5 192.168.0.5 node5
6 192.168.0.6 node6
7 192.168.0.7 node7
8 192.168.0.8 node8
9 192.168.0.9 node9
```

This enables...

```
1 ssh john@node1
```

or, the abbreviated...

```
1 ssh node1
```

provided the user name on the macbook is the same as the Linux user created by cloud-init.

13.2.5 Router/Firewall Configuration

Local network behind firewall/switch: 192.168.0.254

WAN address LAN address

Firewall/Switch (Netgear FVS318G)

Describe DHCP reservations mapping IP to MAC addresses.

Describe ssh access

Add relevant PDFs.

13.2.6 Create the Raspberry Pi Ubuntu Server Image

On macbook...

Download Ubuntu 20.04 LTS 64-bit pre-installed server image for the Raspberry Pi 4...

Double click to uncompress the .xz file which leaves the .img file.

Double click to mount the .img in the filesystem...

Amend /Volumes/system-boot/user-data...

```
1 #cloud-config
2
3 # This is the user-data configuration file for cloud-init. By default this s
4 # up an initial user called "ubuntu" with password "ubuntu", which must be
5 # changed at first login. However, many additional actions can be initiated
6 # first boot from this file. The cloud-init documentation has more details:
7 #
8 # https://cloudinit.readthedocs.io/
9 #
10 # Some additional examples are provided in comments below the default
11 # configuration.
12
13 # On first boot, set the (default) ubuntu user's password to "ubuntu" and
14 # expire user passwords
15 chpasswd:
16     expire: false
17     list:
18     - ubuntu:ubuntu
19     - john:john
20
21 # Enable password authentication with the SSH daemon
22 ssh_pwauth: true
23
24 ## On first boot, use ssh-import-id to give the specific users SSH access to
25 ## the default user
```

```

26 #ssh_import_id:
27 #- lp:my_launchpad_username
28 #- gh:my_github_username
29
30 ## Add users and groups to the system, and import keys with the ssh-import-i
31 ## utility
32 #groups:
33 #- robot: [robot]
34 #- robotics: [robot]
35 #- pi
36 #
37 groups:
38 - john: [john]
39
40 #users:
41 #- default
42 #- name: robot
43 # gecos: Mr. Robot
44 # primary_group: robot
45 # groups: users
46 # ssh_import_id: foobar
47 # lock_passwd: false
48 # passwd: $5$hkui88$nvZgIle31cNpryjRf09uArF7DYiBcWEnjq7L1AQNN3
49 users:
50 - default
51 - name: john
52   gecos: John Duffy
53   primary_group: john
54   sudo: ALL=(ALL) NOPASSWD:ALL
55   shell: /bin/bash
56   ssh_authorized_keys:
57     - ssh-rsa AAAAB3NzaC1yc2EAAAADAQABAAQgQDGsnzP+1Q6NgeeKFTd/+Mom+UCYJTL/wzI
58
59 ## Update apt database and upgrade packages on first boot
60 #package_update: true
61 #package_upgrade: true
62 package_update: true
63 package_upgrade: true
64
65 ## Install additional packages on first boot
66 #packages:
67 #- pwgen
68 #- pastebinit
69 #- [libpython2.7, 2.7.3-0ubuntu3.1]
70 packages:
71 - git

```

```

72 - tree
73 - unzip
74 - iperf
75 - net-tools
76 - linux-tools-common
77 - linux-tools-raspi
78 - build-essential
79 - gdb
80 - openmpi-common
81 - openmpi-bin
82 - libblis3-serial
83 - libblis3-openmp
84 - libopenblas0-serial
85 - libopenblas0-openmp
86
87 ## Write arbitrary files to the file-system (including binaries!)
88 #write_files:
89 #- path: /etc/default/keyboard
90 #   content: |
91 #       # KEYBOARD configuration file
92 #       # Consult the keyboard(5) manual page.
93 #       XKBMODEL="pc105"
94 #       XKBLAYOUT="gb"
95 #       XKBVARIANT=""
96 #       XKBOPTIONS="ctrl: nocaps"
97 #   permissions: '0644'
98 #   owner: root:root
99 #- encoding: gzip
100 # path: /usr/bin/hello
101 # content: !!binary |
102 #   H4sIAIDb/U8C/1NW1E/KzNMvzuBKTc7IV8hIzcnJVyJPL8pJ4QIA6N+MVxsAAAA=
103 #   owner: root:root
104 #   permissions: '0755'
105 write_files:
106 - path: /etc/hosts
107   content: |
108     127.0.0.1 localhost
109     192.168.0.1 node1
110     192.168.0.2 node2
111     192.168.0.3 node3
112     192.168.0.4 node4
113     192.168.0.5 node5
114     192.168.0.6 node6
115     192.168.0.7 node7
116     192.168.0.8 node8
117     192.168.0.9 node9

```

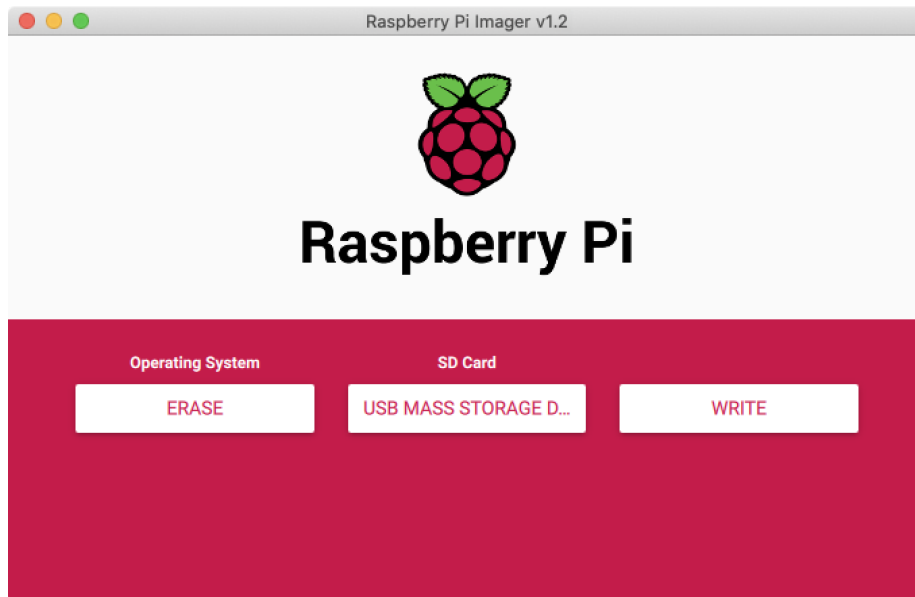


Figure 3: Using Raspberry Pi Imager to erase and format a MicroSD card.

```

118     permissions: '0644'
119     owner: root:root
120
121     ## Run arbitrary commands at rc.local like time
122     #runcmd:
123     #- [ ls, -l, / ]
124     #- [ sh, -xc, "echo $(date) ': hello world!'" ]
125     #- [ wget, "http://ubuntu.com", -O, /run/mydir/index.html ]
126     runcmd:
127     - hostnamectl set-hostname --static node$(hostname -i | cut -d ' ' -f 1 | cu
128     - reboot

```

Eject/unmount .img file

Use Raspberry Pi Imager to erase...

Then use the Raspberry Pi Imager to write preinstalled server image to the MicroSD card...

When complete, remove the MicroSD card from the card reader, place it the Raspberry Pi and plug in the power cable.

The cloud-init configuration process will now start. The Raspberry Pi will ac-



Figure 4: Using Raspberry Pi Imager to write the server image to a MicroSD card.

quire its IP address from the router, setup users, update apt, upgrade the system, download software packages, set the hostname (based on the IP address), and finally the system will reboot.

13.3 Post-Installation Tasks

13.3.1 Enable No Password Access

This is required for Open-MPI...

Our public key was installed on each node by cloud-init. So, we can ssh into each node without a password, and use the abbreviated ssh node1, instead of ssh john@node1 (assuming john is the user name on the workstation).

We need to copy our private key to node1 (only node1)...

```
1 scp ~/.ssh/id_rsa node1:~/.ssh
```

Then to enable access to nodes node2 to node9 without a password from node1, we need to import the ... keys into the node1 knownhosts file...

This is easily done...

From macbook, ssh into node1...

```
ssh node1
```

and then from node1, for each of the nodes node2 to node9:

```
ssh node2
```

This will generate...

```
1 The authenticity of host 'node2 (192.168.0.2)' can't be established.  
2 ECDSA key fingerprint is SHA256:5VgsnN2nPvpfbJmALh3aJd0eT/NvDXqN8TCreQyNaFA.  
3 Are you sure you want to continue connecting (yes/no/[fingerprint])?
```

responding yes, imports the key into the node1 knownhosts file...

```
1 exit
```

Next node...

This is only required to be done on initial contact with nodes node2 to node9
(unless the keys on these nodes change)

13.3.2 Uninstall unattended-upgrades

The package unattended-upgrades is installed automatically...

Can potentially interfere with long running benchmarks...

Remove...

From macbookpro:

```
1 ssh node1 sudo apt remove unattended-upgrades --yes  
2 ssh node2 sudo apt remove unattended-upgrades --yes  
3 ssh node3 sudo apt remove unattended-upgrades --yes  
4 ssh node4 sudo apt remove unattended-upgrades --yes  
5 ssh node5 sudo apt remove unattended-upgrades --yes  
6 ssh node6 sudo apt remove unattended-upgrades --yes  
7 ssh node7 sudo apt remove unattended-upgrades --yes  
8 ssh node8 sudo apt remove unattended-upgrades --yes  
9 ssh node9 sudo apt remove unattended-upgrades --yes
```

Don't forget to update your cluster regularly at convenient times...

See update/upgrade script below...

13.3.3 Add Source Repositories

We are going to be rebuilding some packages from source...

```
1 ssh node1
2 sudo touch /etc/apt/sources.list.d/picluster.list
3 sudo vim /etc/apt/sources.list.d/picluster.list
```

... and add the following source repositories...

Listing 4: /etc/apt/sources.list.d/picluster.list

```
1 deb-src http://archive.ubuntu.com/ubuntu focal main universe
2 deb-src http://archive.ubuntu.com/ubuntu focal-updates main universe
```

... and then update the repository cache?

```
1 sudo apt update
```

13.3.4 Create a Project Repository

Xpand upon...

```
1 ssh node1
2 mkdir picluster
3 cd picluster
4 git init
```

Ensure you do

git add git commit git push

at regular intervals...

13.3.5 Create an System Update/Upgrade Script

Automate...

```
1 #!/usr/bin/bash
2
```

```
3 NODES=9
4
5 for (( i=$NODES; i>0; i-- ))
6 do
7     echo ""
8     echo "UPGRADING node$i..."
9     ssh node$i sudo apt update
10    ssh node$i sudo apt full-upgrade --yes
11    ssh node$i sudo apt autoremove --yes
12    ssh node$i sudo shutdown -r now
13 done
```

13.3.6 Select BLAS Library

We have installed four BLAS libraries...

Confirm all nodes are using the same one initially...

```
ssh node1 sudo update-alternatives --config libblas.so.3-aarch64-linux-gnu
```

TODO screen output...

Confirm option 0, OpenBLAS, is selected. Press return to keep this option and then exit.

14 Appendix B - High-Performance Linpack (HPL) Installation

Download and install the latest version of HPL on node1...

```
1 ssh node1
2 cd picluster
3 mkdir hpl
4 cd hpl
5 wget https://www.netlib.org/benchmark/hpl/hpl-2.3.tar.gz
6 gunzip hpl-2.3.tar.gz
7 tar xvf hpl-2.3.tar
8 rm hpl-2.3.tar
9 cd hpl-2.3
```

Create Make.serial file...

```
1 cd setup
2 bash make_generic
3 cd ..
4 cp setup/Make.UNKNOWN Make.serial
```

Amend Make.serial as per...

Build...

```
1 make arch=serial
```

This creates xhpl and HPL.dat in bin/serial

Copy xhpl to all nodes (only xhpl, and not HPL.dat)...

```
1 ssh node2 mkdir -p picluster/hpl/hpl-2.3/bin/serial
2 ssh node3 mkdir -p picluster/hpl/hpl-2.3/bin/serial
3 ssh node4 mkdir -p picluster/hpl/hpl-2.3/bin/serial
4 ssh node5 mkdir -p picluster/hpl/hpl-2.3/bin/serial
5 ssh node6 mkdir -p picluster/hpl/hpl-2.3/bin/serial
6 ssh node7 mkdir -p picluster/hpl/hpl-2.3/bin/serial
7 ssh node8 mkdir -p picluster/hpl/hpl-2.3/bin/serial
8 ssh node9 mkdir -p picluster/hpl/hpl-2.3/bin/serial
9
10 scp bin/serial/xhpl node2:~picluster/hpl/hpl-2.3/bin/serial
11 scp bin/serial/xhpl node3:~picluster/hpl/hpl-2.3/bin/serial
12 scp bin/serial/xhpl node4:~picluster/hpl/hpl-2.3/bin/serial
13 scp bin/serial/xhpl node5:~picluster/hpl/hpl-2.3/bin/serial
```

```
14 scp bin/serial/xhpl node6:~piccluster/hpl/hpl-2.3/bin/serial
15 scp bin/serial/xhpl node7:~piccluster/hpl/hpl-2.3/bin/serial
16 scp bin/serial/xhpl node8:~piccluster/hpl/hpl-2.3/bin/serial
17 scp bin/serial/xhpl node9:~piccluster/hpl/hpl-2.3/bin/serial
```

15 Appendix ? - High-Performance Linpack (HPL) Installation

Download and install the latest version of HPL on node1...

```
1 ssh node1
2 cd picluster
3 mkdir hpl
4 cd hpl
5 wget https://www.netlib.org/benchmark/hpl/hpl-2.3.tar.gz
6 gunzip hpl-2.3.tar.gz
7 tar xvf hpl-2.3.tar
8 rm hpl-2.3.tar
9 cd hpl-2.3
```

Create Make.serial file...

```
1 cd setup
2 bash make_generic
3 cd ..
4 cp setup/Make.UNKNOWN Make.serial
```

Amend Make.serial as per...

Build...

```
1 make arch=serial
```

This creates xhpl and HPL.dat in bin/serial

Copy xhpl to all nodes (only xhpl, and not HPL.dat)...

```
1 ssh node2 mkdir -p picluster/hpl/hpl-2.3/bin/serial
2 ssh node3 mkdir -p picluster/hpl/hpl-2.3/bin/serial
3 ssh node4 mkdir -p picluster/hpl/hpl-2.3/bin/serial
4 ssh node5 mkdir -p picluster/hpl/hpl-2.3/bin/serial
5 ssh node6 mkdir -p picluster/hpl/hpl-2.3/bin/serial
6 ssh node7 mkdir -p picluster/hpl/hpl-2.3/bin/serial
7 ssh node8 mkdir -p picluster/hpl/hpl-2.3/bin/serial
8 ssh node9 mkdir -p picluster/hpl/hpl-2.3/bin/serial
9
10 scp bin/serial/xhpl node2:~picluster/hpl/hpl-2.3/bin/serial
11 scp bin/serial/xhpl node3:~picluster/hpl/hpl-2.3/bin/serial
12 scp bin/serial/xhpl node4:~picluster/hpl/hpl-2.3/bin/serial
13 scp bin/serial/xhpl node5:~picluster/hpl/hpl-2.3/bin/serial
```

```
14 scp bin/serial/xhpl node6:~piccluster/hpl/hpl-2.3/bin/serial
15 scp bin/serial/xhpl node7:~piccluster/hpl/hpl-2.3/bin/serial
16 scp bin/serial/xhpl node8:~piccluster/hpl/hpl-2.3/bin/serial
17 scp bin/serial/xhpl node9:~piccluster/hpl/hpl-2.3/bin/serial
```


16 Appendix ? - Rebuild Package libblis3-serial:TODO CHECK NAME

```
1 ssh node1
2 mkdir -p picluster/build/blis
3 cd picluster/build/blis
4 apt-get source blis
5 sudo apt-get build-dep blis
6 cd blis-0.6.1
```

Hints from experience... and time savers... for building a development cluster on a local network.

16.1 IP/MAC Addresses

If IP/MAC address assignments get confused, which is easily done during initial build, view IP address assignments on the local network with:

```
arp -a
```

Then delete *incomplete* IP addresses with:

```
sudo arp -d incomplete-ip-address
```

16.2 SSH known_hosts

If *ssh* reports differing keys in 'known-hosts', and warns of a potential 'man-in-the-middle-attack', then just delete 'known-hosts':

```
sudo rm ~/.ssh/known_hosts
```

'known_hosts' will be re-populated as you log into each node.

16.3 tmux

tmux is your friend!

Monitoring long running jobs from a workstation, which goes to sleep after a period of no activity, for example, may interfere with the running of the jobs if a SSH connection is broken.

Use a `tmux` session to start long running jobs, and then detach from the `tmux` session. The job will quite happily run in the background on the cluster. Turn the workstation off and go to bed. In the morning, turn the workstation on and 'attach' to the `tmux` session. All will be well.

16.4 git

`git` is your best friend!

During your cluster build you will accidentally delete files, results etc. After every significant...

17 Appendix ? - cloud-init user-data

Listing 5: picluster/cloudinit/user-data

```
1  #cloud-config
2
3  # This is the user-data configuration file for cloud-init. By default this s
4  # up an initial user called "ubuntu" with password "ubuntu", which must be
5  # changed at first login. However, many additional actions can be initiated
6  # first boot from this file. The cloud-init documentation has more details:
7  #
8  # https://cloudinit.readthedocs.io/
9  #
10 # Some additional examples are provided in comments below the default
11 # configuration.
12
13 # On first boot, set the (default) ubuntu user's password to "ubuntu" and
14 # expire user passwords
15 chpasswd:
16     expire: false
17     list:
18         - ubuntu:ubuntu
19         - john:john
20
21 # Enable password authentication with the SSH daemon
22 ssh_pwauth: true
23
24 ## On first boot, use ssh-import-id to give the specific users SSH access to
25 ## the default user
26 #ssh_import_id:
27 #- lp:my_launchpad_username
28 #- gh:my_github_username
29
30 ## Add users and groups to the system, and import keys with the ssh-import-i
31 ## utility
32 #groups:
33 #- robot: [robot]
34 #- robotics: [robot]
35 #- pi
36 #
37 groups:
38     - john: [john]
39
40 #users:
41 #- default
42 #- name: robot
```

```

43 # gecos: Mr. Robot
44 # primary_group: robot
45 # groups: users
46 # ssh_import_id: foobar
47 # lock_passwd: false
48 # passwd: $5$hkui88$nvZgIle31cNpryjRf09uArF7DYiBcWEnjq7L1AQNN3
49 users:
50 - default
51 - name: john
52   gecos: John Duffy
53   primary_group: john
54   sudo: ALL=(ALL) NOPASSWD:ALL
55   shell: /bin/bash
56   ssh_authorized_keys:
57     - ssh-rsa AAAAB3NzaC1yc2EAAAADAQABAAQgQDGsnzP+1Q6NgeeKFTd/+Mom+UCYJTL/wzI
58
59 ## Update apt database and upgrade packages on first boot
60 #package_update: true
61 #package_upgrade: true
62 package_update: true
63 package_upgrade: true
64
65 ## Install additional packages on first boot
66 #packages:
67 #- pwgen
68 #- pastebinit
69 #- [libpython2.7, 2.7.3-0ubuntu3.1]
70 packages:
71 - git
72 - tree
73 - unzip
74 - iperf
75 - net-tools
76 - linux-tools-common
77 - linux-tools-raspi
78 - build-essential
79 - gdb
80 - openmpi-common
81 - openmpi-bin
82 - libblis3-serial
83 - libblis3-openmp
84 - libopenblas0-serial
85 - libopenblas0-openmp
86
87 ## Write arbitrary files to the file-system (including binaries!)
88 #write_files:

```

```

89 #- path: /etc/default/keyboard
90 # content: |
91 #     # KEYBOARD configuration file
92 #     # Consult the keyboard(5) manual page.
93 #     XKBMODEL="pc105"
94 #     XKBLAYOUT="gb"
95 #     XKBVARIANT=""
96 #     XKBOPTIONS="ctrl: nocaps"
97 # permissions: '0644'
98 # owner: root:root
99 #- encoding: gzip
100 # path: /usr/bin/hello
101 # content: !!binary |
102 #     H4sIAIDb/U8C/1NW1E/KzNMvzuBKTc7IV8hIzcnJVyjPL8pJ4QIA6N+MVxsAAAA=
103 # owner: root:root
104 # permissions: '0755'
105 write_files:
106 - path: /etc/hosts
107   content: |
108     127.0.0.1 localhost
109     192.168.0.1 node1
110     192.168.0.2 node2
111     192.168.0.3 node3
112     192.168.0.4 node4
113     192.168.0.5 node5
114     192.168.0.6 node6
115     192.168.0.7 node7
116     192.168.0.8 node8
117     192.168.0.9 node9
118   permissions: '0644'
119   owner: root:root
120
121 ## Run arbitrary commands at rc.local like time
122 #runcmd:
123 #- [ ls, -l, / ]
124 #- [ sh, -xc, "echo $(date) ': hello world!'" ]
125 #- [ wget, "http://ubuntu.com", -O, /run/mydir/index.html ]
126 runcmd:
127 - hostnamectl set-hostname --static node$(hostname -i | cut -d ' ' -f 1 | cut)
128 - reboot

```

18 Appendix ? - Pi Cluster Tools

Listing 6: picluster/tools/picluster-update

```
1 #!/usr/bin/bash
2
3 # A simple bash script to upgrade the cluster.
4
5 NODES=9
6
7 for (( i=$NODES; i>0; i-- ))
8 do
9     echo ""
10    echo "UPGRADING node$i..."
11    ssh node$i sudo apt update
12    ssh node$i sudo apt full-upgrade --yes
13    ssh node$i sudo apt autoremove --yes
14    ssh node$i sudo shutdown -r now
15 done
```

Listing 7: picluster/tools/picluster-reboot

```
1 #!/usr/bin/bash
2
3 # A simple bash script to reboot the cluster.
4
5 NODES=9
6
7 for (( i=$NODES; i>0; i-- ))
8 do
9     echo "Rebooting node$i..."
10    ssh node$i sudo shutdown -r now
11 done
```

Listing 8: picluster/tools/picluster-shutdown

```
1 #!/usr/bin/bash
2
3 # A simple bash script to shutdown the cluster.
4
5 NODES=9
6
7 for (( i=$NODES; i>0; i-- ))
8 do
9     echo "Shutting down node$i..."
10    ssh node$i sudo shutdown -h now
11 done
```

Listing 9: picluster/tools/picluster-libblas-query

```

1  #!/usr/bin/bash
2
3  # A simple bash script to query the current alternative for libblas.
4
5  NODES=9
6
7  for (( i=1; i<=$NODES; i++ ))
8  do
9      printf "node$i: "
10     ssh node$i update-alternatives --query libblas.so.3-aarch64-linux-gnu \
11         | grep Value: \
12         | gawk '{print $2}'
13 done

```

Listing 10: picluster/tools/picluster-libblas-set-openblas-serial

```

1  #!/usr/bin/bash
2
3  # A simple bash script to query the current alternative for libblas.
4
5  NODES=9
6
7  for (( i=1; i<=$NODES; i++ ))
8  do
9      printf "node$i: "
10     ssh node$i update-alternatives --query libblas.so.3-aarch64-linux-gnu \
11         | grep Value: \
12         | gawk '{print $2}'
13 done

```

Listing 11: picluster/tools/picluster-libblas-set-openblas-openmp

```

1  #!/usr/bin/bash
2
3  # A simple bash script to query the current alternative for libblas.
4
5  NODES=9
6
7  for (( i=1; i<=$NODES; i++ ))
8  do
9      printf "node$i: "
10     ssh node$i update-alternatives --query libblas.so.3-aarch64-linux-gnu \
11         | grep Value: \
12         | gawk '{print $2}'
13 done

```

Listing 12: picluster/tools/picluster-libblas-set-blis-serial

```
1 #!/usr/bin/bash
2
3 # A simple bash script to query the current alternative for libblas.
4
5 NODES=9
6
7 for (( i=1; i<=$NODES; i++ ))
8 do
9     printf "node$i: "
10    ssh node$i update-alternatives --query libblas.so.3-aarch64-linux-gnu \
11    | grep Value: \
12    | gawk '{print $2}'
13 done
```

Listing 13: picluster/tools/picluster-libblas-set-blis-openmp

```
1 #!/usr/bin/bash
2
3 # A simple bash script to query the current alternative for libblas.
4
5 NODES=9
6
7 for (( i=1; i<=$NODES; i++ ))
8 do
9     printf "node$i: "
10    ssh node$i update-alternatives --query libblas.so.3-aarch64-linux-gnu \
11    | grep Value: \
12    | gawk '{print $2}'
13 done
```