Arm Performance Reports User Guide

Version 20.0.3



Non-Confidential Proprietary Notice

This document is protected by copyright and other related rights and the practice or implementation of the information contained in this document may be protected by one or more patents or pending patent applications. No part of this document may be reproduced in any form by any means without the express prior written permission of Arm. **No license, express or implied, by estoppel or otherwise to any intellectual property rights is granted by this document unless specifically stated**.

Your access to the information in this document is conditional upon your acceptance that you will not use or permit others to use the information for the purposes of determining whether implementations infringe any third party patents.

THIS DOCUMENT IS PROVIDED "AS IS". ARM PROVIDES NO REPRESENTATIONS AND NO WARRANTIES, EXPRESS, IMPLIED OR STATUTORY, INCLUDING, WITHOUT LIMITATION, THE IMPLIED WARRANTIES OF MERCHANTABILITY, SATISFACTORY QUALITY, NONINFRINGEMENT OR FITNESS FOR A PARTICULAR PURPOSE WITH RESPECT TO THE DOCUMENT. For the avoidance of doubt, Arm makes no representation with respect to, and has undertaken no analysis to identify or understand the scope and content of, third party patents, copyrights, trade secrets, or other rights.

This document may include technical inaccuracies or typographical errors.

TO THE EXTENT NOT PROHIBITED BY LAW, IN NO EVENT WILL ARM BE LIABLE FOR ANY DAMAGES, INCLUDING WITHOUT LIMITATION ANY DIRECT, INDIRECT, SPECIAL, INCIDENTAL, PUNITIVE, OR CONSEQUENTIAL DAMAGES, HOWEVER CAUSED AND REGARDLESS OF THE THEORY OF LIABILITY, ARISING OUT OF ANY USE OF THIS DOCUMENT, EVEN IF ARM HAS BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

This document consists solely of commercial items. You shall be responsible for ensuring that any use, duplication or disclosure of this document complies fully with any relevant export laws and regulations to assure that this document or any portion thereof is not exported, directly or indirectly, in violation of such export laws. Use of the word "partner" in reference to Arm's customers is not intended to create or refer to any partnership relationship with any other company. Arm may make changes to this document at any time and without notice.

If any of the provisions contained in these terms conflict with any of the provisions of any click through or signed written agreement covering this document with Arm, then the click through or signed written agreement prevails over and supersedes the conflicting provisions of these terms. This document may be translated into other languages for convenience, and you agree that if there is any conflict between the English version of this document and any translation, the terms of the English version of the Agreement shall prevail.

The Arm corporate logo and words marked with ®or ™are registered trademarks or trademarks of Arm Limited (or its subsidiaries) in the US and/or elsewhere. All rights reserved. Other brands and names mentioned in this document may be the trademarks of their respective owners. Please follow Arm's trademark usage guidelines at http://www.arm.com/company/policies/trademarks.

Copyright ©2018, 2019 Arm Limited (or its affiliates). All rights reserved. Arm Limited. Company 02557590 registered in England. 110 Fulbourn Road, Cambridge, England CB1 9NJ. LES-PRE-20349

Confidentiality Status

This document is Non-Confidential. The right to use, copy and disclose this document may be subject to license restrictions in accordance with the terms of the agreement entered into by Arm and the party that Arm delivered this document to.

Unrestricted Access is an Arm internal classification.

Product Status

The information in this document is Final, that is for a developed product.

Web Address

http://www.arm.com

Contents

Co	Contents 3							
1	Intr	oduction	7					
	1.1	Online resources	7					
2	Inst	allation	8					
	2.1	Linux installation	8					
		2.1.1 Graphical install	8					
		2.1.2 Text-mode install	9					
	2.2	License files	10					
	2.3	Workstation and evaluation licenses	10					
	2.4							
	2.5	Architecture licensing	11					
		2.5.1 Using multiple architecture licenses	11					
	2.6	Environment variables	12					
		2.6.1 Report customization	12					
		2.6.2 Warning suppression	12					
		2.6.3 I/O behavior	12					
		2.6.4 Licensing	13					
		2.6.5 Timeouts	13					
		2.6.6 Sampler	13					
		2.6.7 Simple troubleshooting	16					
3	D	ning with an avample program	17					
3		ning with an example program Overview of the example source code	17					
	3.1		17					
	3.2	Compiling	17					
	2.2	3.2.1 Cray X-series						
	3.3 3.4	Running	18 19					
4		ning with real programs	20					
	4.1	Preparing a program for profiling	20					
		4.1.1 Linking	20					
		4.1.2 Dynamic linking on Cray X-Series systems	21					
		4.1.3 Static linking	22					
		4.1.4 Static linking on Cray X-Series systems	24					
		4.1.5 Dynamic and static linking on Cray X-Series systems using the modules envi-	25					
		ronment	25					
	4.2		25					
	4.2	Express Launch mode	26					
	4.2	4.2.1 Compatible MPIs	26					
	4.3	Compatibility Launch mode	26					
	4.4	Generating a performance report	27					
	4.5	Specifying output locations	28					
	4.6	Support for DCIM systems	28					
		4.6.1 Customizing your DCIM script	28					
		4.6.2 Customising the gmetric location	29					
	4.7	Enable and disable metrics	29					
5	Sum	marizing an existing MAP file	30					

6	Inter	preting	g performance reports	31
	6.1	HTML	performance reports	31
	6.2	Report	summary	33
		6.2.1	Compute	33
		6.2.2	MPI	33
		6.2.3	Input/Output	33
	6.3	CPU b	reakdown	33
		6.3.1	Single core code	34
		6.3.2	OpenMP code	34
		6.3.3	Scalar numeric ops	34
		6.3.4	Vector numeric ops	34
		6.3.5	Memory accesses	34
		6.3.6	Waiting for accelerators	34
	6.4		netrics breakdown	35
	0.1	6.4.1	Cycles per instruction	35
		6.4.2	Stalled cycles	35
		6.4.3	L2 cache misses	35
		6.4.4		35
			L3 cache miss per instruction	
		6.4.5	FLOPS scalar lower bound	35
		6.4.6	FLOPS vector lower bound	35
	c =	6.4.7	Memory accesses	36
	6.5	_	IP breakdown	36
		6.5.1	Computation	36
		6.5.2	Synchronization	36
		6.5.3	Physical core utilization	36
		6.5.4	System load	36
	6.6		s breakdown	37
		6.6.1	Computation	37
		6.6.2	Synchronization	37
		6.6.3	Physical core utilization	37
		6.6.4	System load	37
	6. 7	MPI br	eakdown	37
		6.7.1	Time in collective calls	38
		6.7.2	Time in point-to-point calls	38
		6.7.3	Effective process collective rate	38
		6.7.4	Effective process point-to-point rate	38
	6.8	I/O bre	akdown	38
		6.8.1	Time in reads	39
		6.8.2	Time in writes	39
		6.8.3	Effective process read rate	39
		6.8.4	Effective process write rate	39
		6.8.5	Lustre metrics	39
	6.9		y breakdown	40
	0.0	6.9.1	Mean process memory usage	40
		6.9.2	Peak process memory usage	40
		6.9.3	Peak node memory usage	40
	6.10		rator breakdown	41
	0.10		GPU utilization	41
			Global memory accesses	41
			Mean GPU memory usage	41
		0.10.4	Peak GPU memory usage	41

	6.11	Energy breakdown	42
		6.11.1 CPU	42
		6.11.2 Accelerator	42
		6.11.3 System	42
		6.11.4 Mean node power	42
		6.11.5 Peak node power	42
		6.11.6 Requirements	43
	6.12	Textual performance reports	43
		CSV performance reports	43
		Worked examples	44
		6.14.1 Code characterization and run size comparison	44
		6.14.2 Deeper CPU metric analysis	44
		6.14.3 I/O performance bottlenecks	44
		of the temperature sections of the temperature of t	
7	Con	figurable Perf metrics	45
	7.1	Permissions	45
	7.2	Probing target hosts	45
	7.3	Specifying Perf metrics via the command line	46
	7.4	Specifying Perf metrics via a file	46
	7.5	Viewing events	47
	7.6	Advanced configuration	47
8		figuration	49
	8.1	Compute node access	49
A	Gett	ing support	50
B	Suni	ported platforms	51
В		ported platforms Performance Reports	51
В		Ported platforms Performance Reports	51 51
	B.1	·	
	B.1	Performance Reports	51
	B.1 Kno MPI	Performance Reports	51 53 54
C	B.1 Kno MPI	Performance Reports	51 53 54
C	B.1 Kno MPI	Performance Reports	51 53 54
C	B.1 Kno MPI	Performance Reports	51 53 54 54
C	B.1 Kno MPI D.1 D.2	Performance Reports wn issues distribution notes Bull MPI	51 53 54 54 54
C	B.1 Kno MPI D.1 D.2 D.3	Performance Reports wn issues distribution notes Bull MPI	51 53 54 54 54 54
C	B.1 Kno MPI D.1 D.2 D.3 D.4	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2	51 53 54 54 54 54 55
C	B.1 Kno MPI D.1 D.2 D.3 D.4 D.5	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2 MPICH 3	51 53 54 54 54 54 55 55
C	B.1 Kno MPI D.1 D.2 D.3 D.4 D.5	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2 MPICH 3 Open MPI	51 53 54 54 54 55 55 55
C	B.1 Kno MPI D.1 D.2 D.3 D.4 D.5 D.6	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2 MPICH 3 Open MPI D.6.1 Open MPI 3.x on IBM Power with the GNU compiler	513 533 544 544 555 555 555
C	B.1 Kno MPI D.1 D.2 D.3 D.4 D.5 D.6	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2 MPICH 3 Open MPI D.6.1 Open MPI 3.x on IBM Power with the GNU compiler Platform MPI	513 534 544 545 555 555 555 566
C	B.1 Kno MPI D.1 D.2 D.3 D.4 D.5 D.6	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2 MPICH 3 Open MPI D.6.1 Open MPI 3.x on IBM Power with the GNU compiler Platform MPI SGI MPT / SGI Altix	513 544 544 555 555 556 566 566
C	B.1 Kno MPI D.1 D.2 D.3 D.4 D.5 D.6 D.7 D.8 D.9	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2 MPICH 3 Open MPI D.6.1 Open MPI 3.x on IBM Power with the GNU compiler Platform MPI SGI MPT / SGI Altix SLURM npiler notes	51 53 54 54 54 55 55 55 56 56 56
C D	B.1 Kno MPI D.1 D.2 D.3 D.4 D.5 D.6 D.7 D.8 D.9	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2 MPICH 3 Open MPI D.6.1 Open MPI 3.x on IBM Power with the GNU compiler Platform MPI SGI MPT / SGI Altix SLURM	51 53 54 54 54 55 55 55 56 56 56 57
C D	B.1 Kno MPI D.1 D.2 D.3 D.4 D.5 D.6 D.7 D.8 D.9 Com	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2 MPICH 3 Open MPI D.6.1 Open MPI 3.x on IBM Power with the GNU compiler Platform MPI SGI MPT / SGI Altix SLURM npiler notes	51 53 54 54 54 55 55 55 56 56 56
C D	B.1 Kno MPI D.1 D.2 D.3 D.4 D.5 D.6 D.7 D.8 D.9 Com E.1	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2 MPICH 3 Open MPI D.6.1 Open MPI 3.x on IBM Power with the GNU compiler Platform MPI SGI MPT / SGI Altix SLURM appler notes AMD OpenCL compiler	51 53 54 54 54 55 55 55 56 56 56 57
C D	B.1 Kno MPI D.1 D.2 D.3 D.4 D.5 D.6 D.7 D.8 D.9 Com E.1 E.2	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2 MPICH 3 Open MPI D.6.1 Open MPI 3.x on IBM Power with the GNU compiler Platform MPI SGI MPT / SGI Altix SLURM npiler notes AMD OpenCL compiler Berkeley UPC compiler	51 53 54 54 54 55 55 56 56 56 57 57
C D	B.1 Kno MPI D.1 D.2 D.3 D.4 D.5 D.6 D.7 D.8 D.9 Cont E.1 E.2 E.3	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2 MPICH 3 Open MPI D.6.1 Open MPI 3.x on IBM Power with the GNU compiler Platform MPI SGI MPT / SGI Altix SLURM upiler notes AMD OpenCL compiler Berkeley UPC compiler Cray compiler environment	51 53 54 54 54 54 55 55 55 56 56 56 57 57
C D	B.1 Kno MPI D.1 D.2 D.3 D.4 D.5 D.6 D.7 D.8 D.9 Cont E.1 E.2 E.3	Performance Reports wn issues distribution notes Bull MPI Cray MPT Intel MPI MPICH 2 MPICH 3 Open MPI D.6.1 Open MPI 3.x on IBM Power with the GNU compiler Platform MPI SGI MPT / SGI Altix SLURM upiler notes AMD OpenCL compiler Berkeley UPC compiler Cray compiler environment GNU	51 53 54 54 54 55 55 55 56 56 57 57 57

F	Platform notes			58
	F.1	Intel X	eon	58
		F.1.1	Enabling RAPL energy and power counters when profiling	58
	F.2	NVIDI	A CUDA	58
	F.3	Arm .		58
		F.3.1	Arm®v8 (AArch64) known issues	58
	F.4	POWE	R8 and POWER9 (POWER 64-bit) known issues	58
G	Gen	eral tro	ubleshooting	59
	G .1	Starting	g a program	59
		G.1.1	Problems starting scalar programs	59
		G.1.2	Problems starting multi-process programs	59
		G.1.3	No shared home directory	60
		nance Reports specific issues	60	
		G.2.1	MPI wrapper libraries	60
		G.2.2	Thread support limitations	60
		G.2.3	No thread activity while blocking on an MPI call	61
		G.2.4	I am not getting enough samples	61
		G.2.5	Performance Reports is reporting time spent in a function definition	61
		G.2.6	Performance Reports is not correctly identifying vectorized instructions	62
		G.2.7	Performance Reports takes a long time to gather and analyze my OpenBLAS-	
			linked application	62
		G.2.8	Performance Reports over-reports MPI, I/O, accelerator or synchronization time	62
	G. 3	Obtaini	ing support	63
G.4 Arm IPMI Energy Agent		Arm IP	MI Energy Agent	64
		G/11	Requirements	64

Introduction

Arm Performance Reports is a low-overhead tool that produces one-page text and HTML reports summarizing and characterizing both scalar and MPI application performance.

Arm Performance Reports provides the most effective way to characterize and understand the performance of HPC application runs.

One single page HTML report elegantly answers a range of vital questions for any HPC site:

- Is this application optimized for the system it is running on?
- Does it benefit from running at this scale?
- Are there I/O or networking bottlenecks affecting performance?
- Which hardware, software or configuration changes can be made to improve performance further?

It is based on MAP's low-overhead adaptive sampling technology that keeps data volumes collected and application overhead low:

- Runs transparently on optimized production-ready codes by adding a single command to your scripts.
- Just 5% application slowdown even with thousands of MPI processes.

Chapters 3 to 6 of this manual describe Arm Performance Reports in more detail.

Online resources

You can find links to tutorials, training material, webinars and white papers in our online knowledge center:

Knowledge Center Arm help and tutorials

Known issues and the latest version of this user guide may be found on the support web pages:

Support Arm Developer website

Installation

A release of Arm Performance Reports can be downloaded from the Arm Developer website.

Both a graphical and text-based installer are provided. See the following sections for details.

Linux installation

Graphical install

Untar the package and run the installer executable, using:

```
tar xf arm-reports-20.0.3-<distro>-<arch>.tar
cd arm-reports-20.0.3-<distro>-<arch>
./installer
```

replacing <distro> and <arch> with the OS distribution and architecture of your tar package, respectively. For example, the tarball package for Redhat 7.4 OS and Armv8-A (AArch64) architecture is: arm-reports-20.0.3 -Redhat-7.4-aarch64.tar

The installer consists of a number of pages where you can choose install options. Use the *Next* and *Back* buttons to move between pages or *Cancel* to cancel the installation.

The *Install Type* page allows you to choose which user(s) to install Arm Performance Reports for.

If you are an administrator (root) you can install Arm Performance Reports for *All Users* in a common directory, such as /opt or /usr/local, otherwise only the *Just For Me* option is enabled.

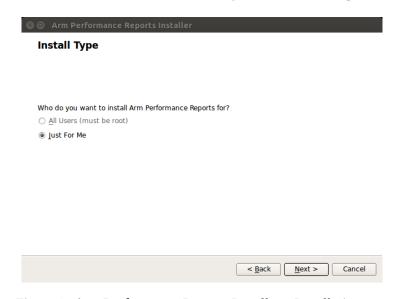


Figure 1: Arm Performance Reports Installer—Installation type

Once you have selected the installation type, you are prompted to specify the directory you would like to install Arm Performance Reports in. For a cluster installation, choose a directory that is shared between the cluster login or frontend node and the compute nodes. Alternatively, install it on or copy it to the same location on each node.

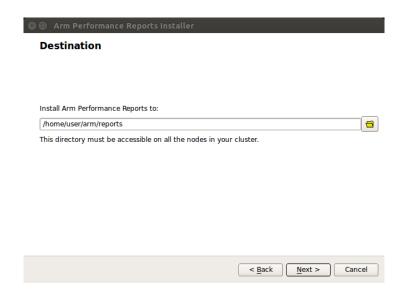


Figure 2: Arm Performance Reports Installer—Installation directory

You are shown the progress of the installation on the *Install* page.

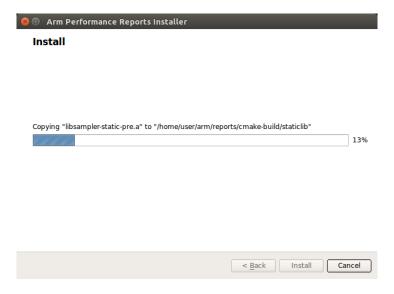


Figure 3: Install in progress

Arm Performance Reports does not have a GUI and does not add any desktop icons.

It is important to follow the instructions in the README file that is contained in the tar file. In particular, you need a valid license file. Use the following link to obtain an evaluation license Get software.

Due to the large number of different site configurations and MPI distributions that are supported by Arm Performance Reports, it is inevitable that you may need to take further steps to get everything fully integrated into your environment. For example, it may be necessary to ensure that environment variables are propagated to remote nodes, and ensure that the tool libraries and executables are available on the remote nodes.

Text-mode install

The text-mode install script textinstall. sh is useful if you are installing remotely.

To install using the text-mode install script, untar the package and run the textinstall.sh script, using:

```
tar xf arm-reports-20.0.3-<distro>-<arch>.tar
cd arm-reports-20.0.3-<distro>-<arch>
./text-install.sh
```

replacing <distro> and <arch> with the OS distribution and architecture of your tar package, respectively. For example, the tarball package for Redhat 7.4 OS and Armv8-A (AArch64) architecture is: arm-reports-20.0.3 -Redhat-7.4-aarch64.tar

Next, you are prompted with the license agreement. To read the license, press *Return*. Following the license prompt, you are requested to enter the directory where you want to install Arm Performance Reports. This directory must be accessible on all the nodes in your cluster. Enter a directory for the installation.

Alternatively, to run the text-mode install script textinstall.sh, accept the license, and point to an installation directory in one step, pass the arguments --accept-licence and <installation-directory> when executing textinstall.sh. For example:

```
./textinstall.sh --accept-licence <installation_directory>
```

replacing the <installation-directory> with a directory of your choice.

License files

Arm Performance Reports requires a license file for its operation.

Time-limited evaluation licenses are available from the Arm Developer website.

Workstation and evaluation licenses

Workstation and Evaluation license files for Arm Performance Reports do not require Arm Licence Server and should be copied directly to {installation-directory}/licences, for example, /home/user/arm/reports/licences/Licence.ddt. Do not edit the files as this prevents them from working.

You may specify an alternative location of the license directory using an environment variable: ALLINEA_LICENCE_DIR. For example:

```
export ALLINEA_LICENCE_DIR=${HOME}/SomeOtherLicenceDir
```

ALLINEA_LICENSE_DIR is an alias for ALLINEA_LICENCE_DIR.

Supercomputing and other floating licenses

For users with Supercomputing and other floating licenses, the Arm Licence Server must be running on the designated license server machine prior to running Arm Performance Reports.

The Arm Licence Server and instructions for its installation and usage may be downloaded from the Arm Developer website.

The license server download is on the Arm Forge download page.

A floating license consists of two files: the server license, a file name Licence.xxxx, and a client license file Licence.

The client file should be copied to {installation-directory}/licences, for example, /home/user/arm/reports/licences/Licence.

You need to edit the hostname line to contain the host name or IP address of the machine running the Licence Server.

See the Licence Server user guide for instructions on how to install the server license.

Architecture licensing

Licenses issued after the release of Arm Performance Reports 6.1 specify the compute node architectures that they may be used with. Licenses issued prior to this release will enable the x86_64 architecture by default.

Existing users for other architectures will be supplied with new

Using multiple architecture licenses

If you are using multiple license files to specify multiple architectures, Arm recommends that you follow these steps.

Procedure

- 1. Ensure that the default licenses directory is empty.
- 2. Create a directory for each architecture.
- 3. When you want to target a specific architecture, set ALLINEA_LICENSE_DIR to the relevant directory.

Alternatively, set ALLINEA_LICENSE_FILE to specify the license file.

Example

On a site that targets two architectures, x86_64 and AArch64, create a directory for each architecture, and name them licenses_x86_64 and licenses_aarch64. Then, to target the architectures, set the license directories as follows:

To target AArch64:

export ALLINEA_LICENSE_DIR=/path/to/licenses/licenses_aarch64

To target x86_64:

export ALLINEA_LICENSE_DIR=/path/to/licenses/licenses_x86_64

Environment variables

Report customization

Environment variables to customize your reports:

ALLINEA_NOTES

Any text in this environment variable will be included in all reports produced.

ALLINEA MAP_TO_DCIM

Allows you to specify a .map file when using the --dcim-output argument.

ALLINEA DCIM SCRIPT

Path to the script to use to communicate with DCIM. Default is

\\${ALLINEA_TOOLS_PATH}/performance-reports/ganglia-connector/pr-dcim }.

ALLINEA_GMETRIC

Path to the gmetric instance to use. This is specific to the pr-dcim script. Default is which gmetric.

Warning suppression

Environment variables for warning suppression (for use when autodetection is resulting in erroneous messages):

ALLINEA NO APPLICATION PROBE

Do not attempt to auto-detect MPI or CUDA executables.

ALLINEA_DETECT_APRUN_VERSION

Automatically detect Cray MPT by passing --version to the aprun wrapper and parsing the output.

I/O behavior

Environment variables for handling default I/O behavior:

ALLINEA_NEVER_FORWARD_STDIN

Never forward the stdin of the perf-report command stdin to the program being analyzed, even if not using the GUI. Normally Arm Performance Reports only forwards stdin when running without the GUI.

ALLINEA_ENABLE_ALL_REPORTS_GENERATION

Enables the option in Arm Performance Reports to generate all types of results at once, using the .all extension.

Licensing

Environment variables to handle licensing:

ALLINEA_LICENCE_FILE

Location of the license file. Default is \${ALLINEA_TOOLS_PATH}/Licence

ALLINEA_FORCE_LICENCE_FILE

Location of the license file. This ensures the license file being pointed to is used.

ALLINEA LICENCE DIR

Location of the licenses directory. Default is \${ALLINEA_TOOLS_PATH}/licences.

ALLINEA_MAC_INTERFACE

Specify the host name of the network interface the license is tied to.

Timeouts

Environment variables for handling timeouts:

ALLINEA NO_TIMEOUT

Do not time out if nodes do not connect after a specified length of time. This may be necessary if the MPI subsystem takes unusually long to start processes.

ALLINEA PROCESS_TIMEOUT

Length of time (in ms) to wait for a process to connect to the front end.

ALLINEA_MPI_FINALIZE_TIMEOUT_MS

Length of time (in ms) to wait for MPI_Finalize to end and the program to exit. Default is 300000 (5 minutes). 0 waits forever.

Sampler

Environment variables for handling sampler-related setup, runtime behavior, and backend processing:

ALLINEA_SAMPLER_INTERVAL

Arm Performance Reports takes a sample in each 20 milliseconds period, giving it a default sampling rate of 50Hz. This will be automatically decreased as the run proceeds to ensure a constant number of samples are taken. See ALLINEA SAMPLER NUM SAMPLES.

If your program runs for a very short period of time, you may benefit by decreasing the initial sampling interval. For example, ALLINEA_SAMPLER_INTERVAL=1 sets an initial sampling rate of 1000Hz, or once per millisecond. Higher sampling rates are not supported.

Increasing the sampling frequency from the default is not recommended if there are lots of threads or very deep stacks in the target program because this may not leave sufficient time to complete one sample before the next sample is started.

Note: Custom values for ALLINEA_SAMPLER_INTERVAL may be overwritten by values set from the combination of ALLINEA_SAMPLER_INTERVAL_PER_THREAD and the expected number of threads (from OMP_NUM_THREADS). For more information, see ALLINEA_SAMPLER_INTERVAL_PER_THREAD.

ALLINEA_SAMPLER_INTERVAL_PER_THREAD

To keep overhead low, Arm Performance Reports imposes a minimum sampling interval based on the number of threads. By default, this is 2 milliseconds per thread, thus for eleven or more threads Arm Performance Reports will increase the initial sampling interval to more than 20 milliseconds.

To adjust this behavior set ALLINEA_SAMPLER_INTERVAL_PER_THREAD to the minimum per thread sample time, in milliseconds.

Lowering this value from the default is not recommended if there are lots of threads as this may not leave sufficient time to complete one sample before the next sample is started.

Notes:

- Whether OpenMP is enabled or disabled in Arm Performance Reports, the final script or scheduler values set for OMP_NUM_THREADS will be used to calculate the sampling interval per thread (ALLINEA_SAMPLER_INTERVAL_PER_THREAD). When configuring your job for submission, check whether your final submission script, scheduler or the Arm Performance Reports GUI has a default value for OMP_NUM_THREADS.
- Custom values for ALLINEA_SAMPLER_INTERVAL will be overwritten by values set from the combination of ALLINEA_SAMPLER_INTERVAL_PER_THREAD and the expected number of threads from OMP_NUM_THREADS.

ALLINEA_MPI_WRAPPER

Arm Performance Reports ships with a number of precompiled wrappers, when your MPI is supported Arm Performance Reports will automatically select and use the appropriate wrapper.

To manually compile a wrapper specifically for your system, set ALLINEA_WRAPPER_COMPILE=1 and MPICC and run < path to Arm Performance Reports installation > /map/wrapper/build_wrapper.

This generates the wrapper library ~/.allinea/wrapper/libmap-sampler-pmpi-<hostname>.so with symlinks to the following files:

- ~/.allinea/wrapper/libmap-sampler-pmpi-<hostname>.so.1
- ~/.allinea/wrapper/libmap-sampler-pmpi-<hostname>.so.1.0
- ~/.allinea/wrapper/libmap-sampler-pmpi-<hostname>.so.1.0.0.

ALLINEA WRAPPER COMPILE

To direct Arm Performance Reports to fall back to creating and compiling a just-in-time wrapper, set ALLINEA_WRAPPER_COMPILE=1.

In order to be able to generate a just-in-time wrapper an appropriate compiler must be available on the machine where Arm Performance Reports is running, or on the remote host when using remote connect.

Arm Performance Reports will attempt to auto detect your MPI compiler, however, setting the MPICC environment variable to the path to the correct compiler is recommended.

ALLINEA_MPIRUN

The path of mpirun, mpiexec or equivalent.

If set, ALLINEA_MPIRUN has higher priority than that set in the GUI and the mpirun found in PATH.

ALLINEA_SAMPLER_NUM_SAMPLES

Arm Performance Reports collects 1000 samples per process by default. To avoid generating too much data on long runs, the sampling rate is automatically decreased as the run progresses to ensure only 1000 evenly spaced samples are stored.

You may adjust this by setting ALLINEA_SAMPLER_NUM_SAMPLES=<positiveinteger>.

Note: It is strongly recommended that you leave this value at the default setting. Higher values are not generally beneficial and add extra memory overheads while running your code. With 512 processes, the default setting already collects half a million samples over the job, the effective sampling rate can be very high indeed.

ALLINEA_KEEP_OUTPUT_LINES

Specifies the number of lines of program output to record in .map files. Setting to 0 will remove the line limit restriction, although this is not recommended as it may result in very large .map files if the profiled program produces lots of output.

ALLINEA_KEEP_OUTPUT_LINE_LENGTH

The maximum line length for program output that will be recorded in .map files. Lines containing more characters than this limit will be truncated. Setting to 0 will remove the line length restriction. This is not recommended because it may result in very large .map files if the profiled program produces lots of output per line.

ALLINEA_PRESERVE_WRAPPER

To gather data from MPI calls Arm Performance Reports generates a wrapper to the chosen MPI implementation.

By default, the generated code and shared objects are deleted when Arm Performance Reports no longer needs them.

To prevent Arm Performance Reports from deleting these files set ALLINEA_PRESERVE_WRAPPER=

Note: If you are using remote launch then this variable must be exported in the remote script.

ALLINEA_SAMPLER_NO_TIME_MPI_CALLS

To prevent Arm Performance Reports from timing the time spent in MPI calls, set ALLINEA_SAMPLER_ NO_TIME_MPI_CALLS.

ALLINEA_SAMPLER_TRY_USE_SMAPS

To allow Arm Performance Reports to use /proc/[pid]/smaps to gather memory usage data, set this ALLINEA_SAMPLER_TRY_USE_SMAPS. This is not recommended since it slows down sampling significantly.

MPICC

To create the MPI wrapper Arm Performance Reports will try to use MPICC, then if that fails search for a suitable MPI compiler command in PATH. If the MPI compiler used to compile the target binary is not in PATH (or if there are multiple MPI compilers in PATH) then MPICC should be set.

Simple troubleshooting

Environment variables for simple troubleshooting:

ALLINEA_DEBUG_HEURISTICS

To print the weights and heuristics used to autodetect which MPI is loaded, set to 1.

Running with an example program

This section takes you through compiling and running one of the the example programs.

Overview of the example source code

Compiling

Arm provides a simple 1-D wave equation solver that is useful as a profiling example program. Both C and Fortran variants are provided:

- examples/wave.c
- examples/wave.f90

Both are built using the same makefile, wave.makefile. To navigate and run wave.makefile, use:

```
cd {installation-directory}/examples/
make -f wave.makefile
```

There is also a mixed-mode MPI+OpenMP variant in examples/wave_openmp.c, which is built with the openmp.makefile makefile.

Note: The makefiles for all supplied examples are located in the {installation-directory}/ examples directory.

Depending on the default compiler on your system you may see some errors when running the makefile, for example:

```
pgf90-Error-Unknown switch: -fno-inline
```

By default, this example makefile is set up for the GNU compilers. To setup the makefile for a different compiler, open the examples/wave.makefile file, uncomment the appriopriate compilation command for the compiler you want to use, and comment those of the GNU compiler.

Notes:

- The compilation commands for other popular compilers are already present within the makefile, separated by compiler.
- Although the example makefiles include the -g flag, Arm Performance Reports does *not* require this and you should not use them in your own makefiles.

In most cases Arm Performance Reports can run on an unmodified binary with no recompilation or linking required.

Cray X-series

On Cray X-series systems the example program must either be dynamically linked (using -dynamic) or explicitly linked with the Arm profiling libraries.

Example how to dynamically link:

```
cc -dynamic -g -03 wave.c -o wave -lm -lrt
ftn -dynamic -G2 -O3 wave.f90 -o wave -lm -lrt
```

Example how to explicitly link with the Arm profiling libraries: First create the libraries using the command make-profiler-libraries --platform=cray --lib-type=static:

```
Created the libraries in /home/user/examples:
   libmap-sampler.a
   libmap-sampler-pmpi.a
To instrument a program, add these compiler options:
   compilation for use with MAP - not required for Performance
      Reports:
      -g (or -G2 for native Cray fortran) (and -O3 etc.)
   linking (both MAP and Performance Reports):
      -Wl,@/home/user/examplesm/allinea-profiler.ld ...
         EXISTING MPI LIBRARIES
   If your link line specifies EXISTING MPI LIBRARIES (e.g. -lmpi)
      , then
   these must appear *after* the Arm sampler and MPI wrapper
      libraries in
   the link line. There's a comprehensive description of the link
       ordering
   requirements in the `Preparing a Program for Profiling' section
       of either
  userguide-forge.pdf or userguide-reports.pdf, located in
   /opt/arm/forge/doc/.
```

Then follow the instructions in the output to link the example program with the Arm profiling libraries:

```
cc -g -03 wave.c -o wave -g -Wl,@allinea-profiler.ld -lm -lrt
ftn -G2 -O3 wave.f90 -o wave -G2 -Wl,@allinea-profiler.ld -lm -lrt
```

Running

As this example uses MPI you need to run on a compute node on your cluster. Your site's help pages and support staff can tell you exactly how to do this on your machine. The simplest way when running small programs is often to request an interactive session, as follows:

If you see output similar to this then the example program is compiled and working correctly.

Generating a performance report

Make sure the Arm Performance Reports module for your system has been loaded:

```
$ perf-report --version
Arm Performance Reports
Copyright (c) 2002-2019 Arm Limited (or its affiliates). All
    rights reserved.
```

If this command cannot be found consult the site documentation to find the name of the correct module.

Once the module is loaded, you can add the perf-report command in front of your existing mpiexec command-line:

```
perf-report mpiexec -n 4 examples/wave_c
```

If your program is submitted through a batch queuing system, then modify your submission script to load the Arm module and add the 'perf-report' line in front of the mpiexec command you want to generate a report for.

The program runs as usual, although startup and shutdown may take a few minutes longer while Arm Performance Reports generates and links the appropriate wrapper libraries before running and collects the data at the end of the run. The runtime of your code (between MPI_Init and MPI_Finalize should not be affected by more than a few percent at most.

After the run finishes, a performance report is saved to the current working directory, using a name based on the application executable:

Note that both .txt and .html versions are automatically generated.

Running with real programs

This section will take you through compiling and running your own programs.

Arm Performance Reports is designed to run on unmodified production executables, so in general no preparation step is necessary. However, there is one important exception: statically linked applications require additional libraries at the linking step.

Preparing a program for profiling

In most cases you do not need to recompile your program to use it with Performance Reports, although in some cases it may need to be relinked, as explained in section 4.1.1 Linking.

Note: Ensure that the program does not set or unset the SIGPROF signal handler because this interferes with the Performance Reports profiling function and can cause it to fail.

Note: It is not recommended to use Performance Reports to profile programs that contain instructions to perform MPI profiling using MPI wrappers and the MPI standard profiling interface, PMPI. This is because Performance Reports's own MPI wrappers may conflict with those contained in the program, producing incorrect metrics.

CUDA programs

When compiling CUDA kernels do not generate debug information for device code (the -G or --device-debug flag) as this can significantly impair runtime performance. Use -lineinfo instead, for example:

nvcc device.cu -c -o device.o -g -lineinfo -03

Linking

To collect data from your program, Performance Reports uses two small profiler libraries, map-sampler and map-sampler-pmpi. These must be linked with your program. On most systems Performance Reports can do this automatically without any action by you. This is done via the system's LD_PRELOAD mechanism, which allows an extra library into your program when starting it.

Note: Although these libraries contain the word 'map' they are used for both Arm Performance Reports and Arm MAP.

This automatic linking when starting your program only works if your program is dynamically-linked. Programs may be dynamically-linked or statically-linked, and for MPI programs this is normally determined by your MPI library. Most MPI libraries are configured with --enable-dynamic by default, and mpicc/mpif90 produce dynamically-linked executables that Performance Reports can automatically collect data from.

The map-sampler-pmpi library is a temporary file that is precompiled and copied or compiled at runtime in the directory ~/.allinea/wrapper.

If your home directory will not be accessible by all nodes in your cluster you can change where the map-sampler-pmpi library will be created by altering the shared directory as described in G.1.3 No shared home directory.

The temporary library will be created in the .allinea/wrapper subdirectory to this shared directory.

For Cray X-Series Systems the shared directory is not applicable, instead map-sampler-pmpi is copied into a hidden .allinea sub-directory of the current working directory.

If Performance Reports warns you that it could not pre-load the sampler libraries, this often means that your MPI library was not configured with --enable-dynamic, or that the LD_PRELOAD mechanism is not supported on your platform. You now have three options:

- 1. Try compiling and linking your code dynamically. On most platforms this allows Performance Reports to use the LD_PRELOAD mechanism to automatically insert its libraries into your application at runtime.
- 2. Link MAP's map-sampler and map-sampler-pmpi libraries with your program at link time manually.
 - See 4.1.2 Dynamic linking on Cray X-Series systems, or 4.1.3 Static linking and 4.1.4 Static linking on Cray X-Series systems.
- 3. Finally, it may be that your system supports dynamic linking but you have a statically-linked MPI. You can try to recompile the MPI implementation with --enable-dynamic, or find a dynamically-linked version on your system and recompile your program using that version. This will produce a dynamically-linked program that Performance Reports can automatically collect data from.

Dynamic linking on Cray X-Series systems

If the LD_PRELOAD mechanism is not supported on your Cray X-Series system, you can try to dynamically link your program explicitly with the Performance Reports sampling libraries.

Compiling the Arm MPI Wrapper Library

First you must compile the Arm MPI wrapper library for your system using the make-profiler-libraries --platform=cray --lib-type=shared command.

Note: Performance Reports also uses this library.

```
user@login:~/myprogram$ make-profiler-libraries --platform=cray
   --lib-type=shared
Created the libraries in /home/user/myprogram:
libmap-sampler.so
                        (and .so.1, .so.1.0, .so.1.0.0)
libmap-sampler-pmpi.so (and .so.1, .so.1.0, .so.1.0.0)
To instrument a program, add these compiler options:
compilation for use with MAP - not required for Performance
   Reports:
    -g (or '-G2' for native Cray Fortran) (and -O3 etc.)
linking (both MAP and Performance Reports):
    -dynamic -L/home/user/myprogram -lmap-sampler-pmpi -lmap-
       sampler -W1, --eh-frame-hdr
\b{Note:} These libraries must be on the same NFS/Lustre/GPFS
   filesystem as your
program.
```

Before running your program (interactively or from a queue), set

```
LD_LIBRARY_PATH:
export LD_LIBRARY_PATH=/home/user/myprogram:$LD_LIBRARY_PATH
map ...
or add -W1,-rpath=/home/user/myprogram when linking your program.
```

Linking with the Arm MPI Wrapper Library

```
mpicc -G2 -o hello hello.c -dynamic -L/home/user/myprogram \
    -lmap-sampler-pmpi -lmap-sampler -W1,--eh-frame-hdr
```

PGI Compiler

When linking OpenMP programs you must pass the -Bdynamic command line argument to the compiler when linking dynamically.

When linking C++ programs you must pass the -pgc++libs command line argument to the compiler when linking.

Static linking

If you compile your program statically, that is your MPI uses a static library or you pass the <code>-static</code> option to the compiler, then you must explicitly link your program with the Arm sampler and MPI wrapper libraries.

Compiling the Arm MPI Wrapper Library

First you must compile the Arm MPI wrapper library for your system using the make-profiler-libraries --lib-type=static command.

Note: Performance Reports also uses this library.

/opt/arm/forge/doc/.

```
user@login:~/myprogram$ make-profiler-libraries --lib-type=static
Created the libraries in /home/user/myprogram:
   libmap-sampler.a
   libmap-sampler-pmpi.a
To instrument a program, add these compiler options:
   compilation for use with MAP - not required for Performance
      Reports:
      -q (and -03 etc.)
   linking (both MAP and Performance Reports):
      -Wl,@/home/user/myprogram/allinea-profiler.ld ...
         EXISTING_MPI_LIBRARIES
   If your link line specifies EXISTING_MPI_LIBRARIES (e.g. -lmpi)
   these must appear *after* the Arm sampler and MPI wrapper
      libraries in
   the link line. There's a comprehensive description of the link
       ordering
   requirements in the 'Preparing a Program for Profiling' section
       of either
   userquide-forge.pdf or userquide-reports.pdf, located in
```

Linking with the Arm MPI Wrapper Library

The -Wl,@/home/user/myprogram/allinea-profiler.ld syntax tells the compiler to look in /home/user/myprogram/allinea-profiler.ld for instructions on how to link with the Arm sampler. Usually this is sufficient, but not in all cases. The rest of this section explains how to manually add the Arm sampler to your link line.

PGI Compiler

When linking C++ programs you must pass the -pgc++libs command line argument to the compiler when linking.

The PGI C runtime static library contains an undefined reference to __kmpc_fork_call, which will cause compilation to fail when linking allinea-profiler.ld. Add --undefined __wrap_-_kmpc_fork_call to your link line before linking to the Arm sampler to resolve this.

The PGI compiler must be 14.9 or later. Using earlier versions of the PGI compiler will fail with an error such as "Error: symbol 'MPI_F_MPI_IN_PLACE' can not be both weak and common" due to a bug in the PGI compiler's weak object support.

If you do not have access to PGI compiler 14.9 or later try compiling and the linking Arm MPI wrapper as a shared library as described in 4.1.2 Dynamic linking on Cray X-Series systems Ommit the option --platform=cray if you are not on a Cray.

Cray

When linking C++ programs you may encounter a conflict between the Cray C++ runtime and the GNU C++ runtime used by the Performance Reports libraries with an error similar to the one below:

```
/opt/cray/cce/8.2.5/CC/x86-64/lib/x86-64/libcray-c++-rts.a(rtti.o)
    : In function `__cxa_bad_typeid':
/ptmp/ulib/buildslaves/cfe-82-edition-build/tbs/cfe/lib_src/rtti.c
    :1062: multiple definition of `__cxa_bad_typeid'
/opt/gcc/4.4.4/snos/lib64/libstdc++.a(eh_aux_runtime.o):/tmp/peint
    /gcc/repackage/4.4.4c/BUILD/snos_objdir/x86_64-suse-linux/
    libstdc++-v3/libsupc++/../../../xt-gcc-4.4.4/libstdc++-v3/
    libsupc++/eh_aux_runtime.cc:46: first defined here
```

You can resolve this conflict by removing -lstdc++ and -lgcc_eh from allinea-profiler.ld.

-lpthread

When linking -Wl, @allinea-profiler.ld must go before the -lpthread command line argument if present.

Manual Linking

When linking your program you must add the path to the profiler libraries (-L/path/to/profiler-libraries), and the libraries themselves (-lmap-sampler-pmpi, -lmap-sampler).

The MPI wrapper library (-lmap-sampler-pmpi) must go:

- 1. After your program's object (.0) files.
- 2. *After* your program's own static libraries, for example -lmylibrary.
- 3. *After* the path to the profiler libraries (-L/path/to/profiler-libraries).
- 4. *Before* the MPI's Fortran wrapper library, if any. For example -lmpichf.
- 5. *Before* the MPI's implementation library usually -lmpi.

6. *Before* the Arm sampler library -lmap-sampler.

The sampler library -lmap-sampler must go:

- 1. After the Arm MPI wrapper library.
- 2. *After* your program's object (.0) files.
- 3. *After* your program's own static libraries, for example -lmylibrary.
- 4. After -Wl, --undefined, allinea_init_sampler_now.
- 5. *After* the path to the profiler libraries L/path/to/profiler-libraries.
- 6. Before -lstdc++, -lgcc_eh, -lrt, -lpthread, -ldl, -lm and -lc.

For example:

```
mpicc hello.c -o hello -g -L/users/ddt/arm \
   -lmap-sampler-pmpi \
   -Wl, --undefined, allinea_init_sampler_now \
   -lmap-sampler -lstdc++ -lgcc_eh -lrt \
   -Wl,--whole-archive -lpthread \
   -Wl,--no-whole-archive \
   -Wl,--eh-frame-hdr \
   -1d1 \
   -1m
mpif90 hello.f90 -o hello -g -L/users/ddt/arm \
   -lmap-sampler-pmpi \
   -Wl, --undefined, allinea_init_sampler_now \
   -lmap-sampler -lstdc++ -lgcc_eh -lrt \
   -Wl,--whole-archive -lpthread \
   -Wl,--no-whole-archive \
   -Wl,--eh-frame-hdr \
   -1d1 \
   -1m
```

Static linking on Cray X-Series systems

Compiling the MPI Wrapper Library

On Cray X-Series systems use make-profiler-libraries --platform=cray --lib-type=static instead:

```
Created the libraries in /home/user/myprogram:
    libmap-sampler.a
    libmap-sampler-pmpi.a

To instrument a program, add these compiler options:
    compilation for use with MAP - not required for Performance
    Reports:
        -g (or -G2 for native Cray Fortran) (and -O3 etc.)
    linking (both MAP and Performance Reports):
        -Wl,@/home/user/myprogram/allinea-profiler.ld ...
        EXISTING MPI LIBRARIES
```

```
If your link line specifies EXISTING_MPI_LIBRARIES (e.g. -lmpi) , then
```

these must appear *after* the Arm sampler and MPI wrapper libraries in

the link line. There's a comprehensive description of the link ordering

requirements in the 'Preparing a Program for Profiling' section of either

userguide-forge.pdf or userguide-reports.pdf, located in /opt/arm/forge/doc/.

Linking with the MPI Wrapper Library

```
cc hello.c -o hello -g -Wl,@allinea-profiler.ld
```

ftn hello.f90 -o hello -g -Wl,@allinea-profiler.ld

Dynamic and static linking on Cray X-Series systems using the modules environment

If your system has the Arm module files installed, you can load them and build your application as usual. See section 4.1.6.

- 1. module load reports or ensure that make-profiler-libraries is in your PATH.
- 2. module load map-link-static or module load map-link-dynamic.
- 3. Recompile your program.

map-link modules installation on Cray X-Series

To facilitate dynamic and static linking of user programs with the Arm MPI Wrapper and Sampler libraries Cray X-Series System Administrators can integrate the map-link-dynamic and map-link-static modules into their module system. Templates for these modules are supplied as part of the Arm Performance Reports package.

Copy files share/modules/cray/map-link-* into a dedicated directory on the system.

For each of the two module files copied:

- 1. Find the line starting with **conflict** and correct the prefix to refer to the location the module files were installed, for example, arm/map-link-static. The correct prefix depends on the subdirectory (if any) under the module search path the map-link-* modules were installed.
- 2. Find the line starting with **set MAP_LIBRARIES_DIRECTORY "NONE"** and replace "NONE" with a user writable directory accessible from the login and compute nodes.

After installed you can verify whether or not the prefix has been set correctly with 'module avail', the prefix shown by this command for the map-link-* modules should match the prefix set in the 'conflict' line of the module sources.

Express Launch mode

Arm Performance Reports can be launched by typing its command name in front of an existing mpiexec command:

```
$ perf-report mpiexec -n 256 examples/wave_c 30
```

This startup method is called *Express Launch* and is the simplest way to get started. If your MPI is not yet supported in this mode, you will see an error message like this:

\$ 'MPICH 1 standard' programs cannot be started using Express Launch syntax (launching with an mpirun command).

```
Try this instead:
perf-report --np=256 ./wave_c 20
```

Type perf-report --help for more information.

This is referred to as *Compatibility Mode*, in which the mpiexec command is not included and the arguments to mpiexec are passed via a --mpiargs="args here" parameter.

One advantage of Express Launch mode is that it is easy to modify existing queue submission scripts to run your program under one of the Arm Performance Reports products.

Normal redirection syntax may be used to redirect standard input and standard output.

Compatible MPIs

The following lists the MPI implementations supported by Express Launch:

- Bullx MPI
- Cray X-Series (MPI/SHMEM/CAF)
- Intel MPI
- MPICH 2
- MPICH 3
- Open MPI (MPI/SHMEM)
- Oracle MPT
- Open MPI (Cray XT/XE/XK)
- · Spectrum MPI
- Spectrum MPI (PMIx)
- Cray XT/XE/XK (UPC)

Compatibility Launch mode

Compatibility Mode must be used if Arm Performance Reports does not support Express Launch mode for your MPI, or, for some MPIs, if it is not able to access the compute nodes directly (for example, using ssh).

To use Compatibility Mode replace the mpiexec command with the perf-report command. For example:

```
mpiexec --np=256 ./wave_c 20
```

This would become:

```
perf-report --np=256 ./wave_c 20
```

Only a small number of mpiexec arguments are supported by perf-report (for example, -n and -np). Other arguments must be passed using the --mpiargs="args here" parameter.

For example:

```
mpiexec --np=256 --nooversubscribe ./wave_c 20
```

Becomes:

```
perf-report --mpiargs="--nooversubscribe" --np=256 ./wave_c 20
```

Normal redirection syntax may be used to redirect standard input and standard output.

Generating a performance report

Make sure the Arm Performance Reports module for your system has been loaded:

```
$ perf-report --version
Arm Performance Reports
Copyright (c) 2002-2019 Arm Limited (or its affiliates). All
    rights reserved.
```

If this command cannot be found consult the site documentation to find the name of the correct module.

Once the module is loaded, you can simply add the perf-report command in front of your existing mpiexec command-line:

```
perf-report mpiexec -n 4 examples/wave_c
```

If your program is submitted through a batch queuing system, then modify your submission script to load the Arm module and add the 'perf-report' line in front of the mpiexec command you want to generate a report for.

The program runs as usual, although startup and shutdown may take a few minutes longer while Arm Performance Reports generates and links the appropriate wrapper libraries before running and collects the data at the end of the run. The runtime of your code (between MPI_Init and MPI_Finalize should not be affected by more than a few percent at most.

After the run finishes, a performance report is saved to the current working directory, using a name based on the application executable:

Note that both .txt and .html versions are automatically generated.

You can include a short description of the run or other notes on configuration and compilation settings by setting the environment variable ALLINEA_NOTES before running perf-report:

\$ ALLINEA_NOTES="Run with inp421.dat and mc=1" perf-report mpiexec -n 512 ./parEval.bin --use-mc=1 inp421.dat

The string in the ALLINEA_NOTES environment variable is included in all report files produced.

Specifying output locations

By default, performance reports are placed in the current working directory using an auto-generated name based on the application executable name, for example:

```
wave_f_16p_2013-11-18_23-30.html
wave_f_2p_8t_2013-11-18_23-30.html
```

This is formed by the name, the size of the job, the date, and the time. If using OpenMP, the value of $OMP_NUM_THREADS$ is also included in the name after the size of the job. The name will be made unique if necessary by adding a $_1/_2/...$ suffix.

You can specify a different location for output files using the --output argument:

- -- output=my-report.txt will create a plain text report in the file my-report.txt in the current directory.
- --output=/home/mark/public/my-report.html will create an HTML report in the file/home/mark/public/my-report.html.
- -- output=my-report will create a plain text report in the file my-report.txt and an HTML report in the file my-report.html, both in the current directory.
- --output=/tmp will create reports with names based on the application executable name in /tmp/, for example, /tmp/wave_f_16p_2013-11-18_23\-30.txt and /tmp/wave_f_16p_2013-11-18_23\-30.html.

Support for DCIM systems

Performance Reports includes support for Data Center Infrastructure Management (DCIM) systems.

You can output all the metrics generated by the Performance Reports to a script using the --dcim-output argument. By default, the pr-dcim script is called and the collected metrics are sent to Ganglia (a System Monitoring tool).

The pr-dcim script looks for a gmetric implementation as part of the Ganglia software, and call it as many times as there are metrics.

Customizing your DCIM script

The default pr-dcim script is located in {installation-directory}/performance-reports/ganglia-connector/pr-dcim.

However, you can use your own custom script by specifying the ALLINEA_DCIM_SCRIPT environment variable.

This option is recommended if you are using a System Monitoring tool other than Ganglia.

Such a script is expecting arguments as follows, each argument can be specified once per metric:

- -V{METRIC}={VALUE} (mandatory) specifies that the metric METRIC has the value VALUE.
- -U{METRIC}={UNITS} (optional) specifies that the metric METRIC is expressed in UNITS.
- -T{METRIC}={TITLE} (optional) specifies that the metric METRIC has title TITLE.
- -t{METRIC}={TYPE} (optional) specifies that the metric METRIC has TYPE data type.

Customising the gmetric location

You can specify the path to your gmetric implementation by using the ALLINEA_GMETRIC environment variable.

Your gmetric version must accept the following command line arguments:

- -n {NAME} (mandatory) specifies the name of the metric (starts with com.allinea).
- -t {TYPE} (mandatory) specifies the type of the metric (for example, double or int32).
- - V {VALUE} (mandatory) specifies the value of the metric.
- -g {GROUP} (optional) specifies which groups the metric belongs to (for example allinea).
- -u $\{UNIT\}$ (optional) specifies the unit of the metric. For example, %, Watts, Seconds, and so on.
- -T {TITLE} (optional) specifies the title of the metric.

Enable and disable metrics

- --enable-metrics=METRICS
- --disable-metrics=METRICS

Allows you to specify comma-separated lists which explicitly enable or disable metrics for which data is to be collected. If the metrics specified cannot be found, an error message is displayed and Performance Reports exits. Metrics which are always enabled or disabled cannot be explicitly disabled or enabled. A metrics source library which has all its metrics disabled, either in the XML definition or via --disable-metrics, will not be loaded. Metrics which can be explicitly enabled or disabled can be listed using the --list-metrics option.

Summarizing an existing MAP file

Arm Performance Reports can be used to summarize an application profile generated by Arm MAP. To produce a performance report from an existing MAP output file called profile.map, simply run:

\$ perf-report profile.map

Command-line options which would alter the execution of a program being profiled, such as specifying the number of MPI ranks, have no effect. Options affecting how Performance Reports produces its report, such as --output, work as expected.

For best results the Performance Reports and MAP versions should match, for example, Performance Reports 20.0.3 with MAP 20.0.3. Performance Reports can use MAP files from versions of MAP as old as 5.0.

Interpreting performance reports

This section explains how to interpret the reports produced by Arm Performance Reports.

Reports are generated in both HTML and textual formats for each run of your application by default. The information presented in both of these formats is the same.

If you want to combine Arm Performance Reports with other tools, consider using the CSV output format.

See CSV performance reports for more details.

HTML performance reports

Viewing HTML files is best done on your local machine. Many sites have places you can put HTML files to be viewed from within the intranet. These directories are a good place to automatically send your performance reports to. Alternatively, you can use SCp or SShfs to make the reports available on your computer:

```
$ scp login1:arm/reports/examples/wave_c_4p*.html .
$ firefox wave_c_4p*.html
```

The following report was generated by running the wave_openmp.c example program with 8 MPI processes and 2 OpenMP threads per process on a typical HPC cluster:

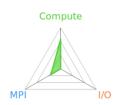
arm PERFORMANCE REPORTS fommand: mpirun -np 8 examples/wave_openmp 60 esources: 1 node (8 physical, 8 logical cores per node)

Memory: 15 GiB per node

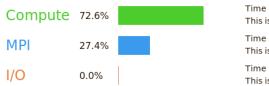
Tasks: 8 processes, OMP_NUM_THREADS was 2

mars

Start time: Tue Nov 7 2017 15:35:50 (UTC)
Total time: 61 seconds (about 1 minutes)
Full path: /scratch/user/reports/examples



Summary: wave openmp is Compute-bound in this configuration



Time spent running application code. High values are usually good. This is **high**; check the CPU performance section for advice

Time spent in MPI calls. High values are usually bad.

This is $\boldsymbol{low};$ this code may benefit from a higher process count

Time spent in filesystem I/O. High values are usually bad.

This is **negligible**; there's no need to investigate I/O performance

This application run was Compute-bound. A breakdown of this time and advice for investigating further is in the CPU section below.

As little time is spent in MPI calls, this code may also benefit from running at larger scales.

CPU

A breakdown of the 72.6% CPU time:

Single-core code 8.2% |
OpenMP regions 91.8% |
Scalar numeric ops 5.1% |
Vector numeric ops 0.0% |
Memory accesses 56.9% |

The per-core performance is memory-bound. Use a profiler to identify time-consuming loops and check their cache performance.

No time is spent in vectorized instructions. Check the compiler's vectorization advice to see why key loops could not be vectorized.

MPI

A breakdown of the 27.4% MPI time:

Time in collective calls

Time in point-to-point calls

Effective process collective rate

Effective process point-to-point rate

305 kB/s

Most of the time is spent in point-to-point calls with a very low transfer rate. This suggests load imbalance is causing synchronization overhead; use an MPI profiler to investigate.

I/O

A breakdown of the 0.0% I/O time:

Time in reads 0.0% |
Time in writes 0.0% |
Effective process read rate 0.00 bytes/s |
Effective process write rate 0.00 bytes/s |

No time is spent in I/O operations. There's nothing to optimize

OpenMP

A breakdown of the 91.8% time in OpenMP regions:

Computation 9.9% Synchronization 90.1% Physical core utilization 100.0% System load 167.0%

Significant time is spent synchronizing threads in parallel regions. Check the affected regions with a profiler.

The system load is high. Ensure background system processes are not running.

Memory

Per-process memory usage may also affect scaling:

Mean process memory usage 38.6 MiB

Peak process memory usage 53.7 MiB

Peak node memory usage 17.0%

The peak node memory usage is very low. Running with fewer MPI processes and more data on each process may be more efficient.

Energy

A breakdown of how energy was used:

CPU not supported % |
System not supported % |
Mean node power not supported W |
Peak node power 0.00 W |

Energy metrics are not available on this system.

CPU metrics are not supported (no intel_rapl module)

Figure 4: *A performance report for the wave_openmp.c example*

When you run a report on this example program, the results might be different to the report shown here depending on the performance and network architecture of the machine on which you run it, but the basic structure of these reports is always the same. This common structure makes comparisons between reports simple, direct and intuitive. Each section of the report is described in the following sections.

Report summary

This characterizes how the wall clock time of the application was spent, and is organized by Compute, MPI, and I/O.

In this example file, you can see that Arm Performance Reports has identified that the program is compute-bound, which means that most of its time is spent inside application code rather than communicating or using the filesystem.

The pieces of advice that the program offers, such as "this code may benefit from running at larger scales", are good starting points for future investigations. They are designed to be informative to scientific users with no previous MPI tuning experience.

The triangular radar chart, in the top-right corner of the report, reflects the values of these three key measurements: compute, MPI and I/O. It is helpful to recognize and compare these triangular shapes when switching between multiple reports.

Compute

Time spent computing. This is the percentage of wall clock time spent in application and in library code, excluding time spent in MPI calls and I/O calls.

MPI

Time spent communicating. This is the percentage of wall clock time spent in MPI calls such as MPI_-Send, MPI_Reduce and MPI_Barrier.

Input/Output

Time spent reading from and writing to the filesystem. This is the percentage of wall clock time spent in system library calls such as read, write and close.

Note: All time spent in MPI-IO calls is included here, even though some communication between processes might also be performed by the MPI library. MPI_File_close is treated as time spent writing, which is often, but not always, correct.

CPU breakdown

Note: All of the metrics described in this section are only available on x86_64 systems.

This section organizes the time spent in application and library code further by analyzing the kinds of instructions that this time was spent on.

Note: All percentages here are relative to the compute time, not to the entire application run. Time spent in MPI and I/O calls is not represented inside this section.

Single core code

The percentage of wall clock in which the application executed using only one core per process, rather than multithreaded or OpenMP code. If you have a multithreaded or OpenMP application, a high value here indicates that your application is bound by Amdahl's law and that scaling to larger numbers of threads will not meaningfully improve performance.

OpenMP code

The percentage of wall clock time spent in OpenMP regions. The higher this is, the better. This metric is only shown if the program spent a measurable amount of time inside at least one OpenMP region.

Scalar numeric ops

The percentage of time spent executing arithmetic operations such as add, mul, div. This does not include time spent using the more efficient vectorized versions of these operations.

Vector numeric ops

The percentage of time spent executing vectorized arithmetic operations such as Intel's SSE2 / AVX extensions.

Generally it is good if a scientific code spends most of its time in these operations, because that is the only way to achieve anything close to the peak performance of modern processors.

If this value is low, you can check the vectorization report of the compiler to understand why the most time consuming loops are not using these operations. Compilers need a lot of help to efficiently vectorize nontrivial loops and the investment in time is often rewarded with 2x–4x performance improvements.

Memory accesses

The percentage of time spent in memory access operations, such as mov, load, store. A portion of the time spent in instructions that use indirect addressing is also included here. A high figure here shows the application is memory-bound and is not able to take full advantage of the CPU resources. Often it is possible to reduce this figure by analyzing loops for poor cache performance and problematic memory access patterns, improving performance significantly.

A high percentage of time spent in memory accesses in an OpenMP program is often a scalability problem. If each core spends most of its time waiting for memory, or the L3 cache, then adding further cores rarely improves matters. Equally, false sharing, in which cores block attempts to access the same cache lines, and the over-use of the atomic pragma, show up as increased time spent in memory accesses.

Waiting for accelerators

The percentage of time that the CPU is waiting for the accelerator.

CPU metrics breakdown

This section presents key CPU performance measurements gathered using the Linux perf event subsystem.

Note: Metrics described in this section are only available on Armv8 and IBM Power systems. These metrics are not available on virtual machines. Linux perf events performance events counters must be accessible on all systems on which the target program runs. See section F.3.1 or F.4 for details.

Cycles per instruction

The average amount of CPU cycles lapsed for each retired instruction. This metric is affected by CPU frequency scaling and various issues, particularly hardware interrupt counts.

Stalled cycles

Note: This metric is available on Armv8 and IBM Power 9 systems only.

The percentage of CPU cycles that lapsed, on which operation instructions are not issued.

L2 cache misses

Note: This metric is available on Armv8 systems only.

The percentage of L2 data cache accesses which resulted in a miss.

L3 cache miss per instruction

Note: This metric is available on IBM Power 9 systems only.

The ratio of L3 data cache demand loads to instructions completed.

FLOPS scalar lower bound

This is a lower bound because its value is calculated from FLOPS vector lower bound, which does not account for the length of vector operations.

Note: This metric is available on IBM Power 8 systems only.

The rate at which floating-point scalar operations finished.

FLOPS vector lower bound

This is a lower bound because the counted value does not account for the length of vector operations.

Note: This metric is available on IBM Power 8 systems only.

The rate at which vector floating-point instructions completed.

Memory accesses

Note: This metric is available on IBM Power 8 systems only.

The rate at which the processor's data cache reloaded from a memory location, including L4 from local, remote, or distant due to a demand load.

OpenMP breakdown

This section breaks down the time spent in OpenMP regions into computation and synchronization and includes additional metrics that help to diagnose OpenMP performance problems. It is only shown if a measurable amount of time was spent inside OpenMP regions.

Computation

The percentage of time threads in OpenMP regions that is spent computing rather than waiting or sleeping. Keeping this high is one important way to ensure that OpenMP codes scale well. If this is high, look at the CPU breakdown to see whether that time is being used optimally on floating-point operations for example, or whether the cores are mostly waiting for memory accesses.

Synchronization

The percentage of time threads in OpenMP regions spent waiting or sleeping. By default, each OpenMP region ends with an implicit barrier. If the workload is imbalanced and some threads finish sooner and wait, this value will increase. Also, there is some overhead associated with entering and leaving OpenMP regions, and a high synchronization time, might show that the threading is too fine-grained. In general, OpenMP performance is better when outer loops are parallelized, rather than inner loops.

Physical core utilization

Modern CPUs often have multiple *logical* cores for each *physical* cores. This is often referred to as hyperthreading. These logical cores can share logic and arithmetic units. Some programs perform better when using additional logical cores, but most HPC codes do not.

If the value here is greater than 100, OMP_NUM_THREADS is set to a larger number of threads than physical cores that are available and performance can be impacted, usually appearing as a larger percentage of time in OpenMP synchronization or memory accesses.

System load

The number of active (running or runnable) threads as a percentage of the number of physical CPU cores that are present in the compute node. This value can exceed 100% if you are using hyper-threading, the cores are *oversubscribed*, or other system processes and daemons start running and take CPU resources away from your program. A value consistently less than 100% might indicate your program is not taking full advantage of the CPU resources available on a compute node.

Threads breakdown

This section organizes the time spent by worker threads (non-main threads) into computation and synchronization, and includes additional metrics that help to diagnose multicore performance problems. This section is replaced by the OpenMP Breakdown if a measurable amount of application time was spent in OpenMP regions.

Computation

The percentage of time that worker threads spend computing rather than waiting in locks and synchronization primitives. If this is high, look at the CPU breakdown to see whether that time is used optimally on floating-point operations for example, or whether the cores are mostly waiting for memory accesses.

Synchronization

The percentage of time worker threads spend waiting in locks and synchronization primitives. This only includes time in which those threads were active on a core and does not include time spent sleeping while other useful work is being done. A large value here indicates a performance and scalability problem that can be detected with a multicore profiler such as Arm MAP.

Physical core utilization

Modern CPUs often have multiple *logical* cores for each *physical* core. This is often referred to as hyperthreading. These logical cores can share logic and arithmetic units. Some programs perform better when using additional logical cores, but most HPC codes do not.

The value here shows the percentage utilization of physical cores. A value over 100% indicates that more threads are executing than there are physical cores, indicating that hyper-threading is in use.

Only threads actively and simultaneously consuming CPU time are included in this metric. A program can have many helper threads that do little except sleep, and are not shown.

System load

The number of active (running or runnable) threads as a percentage of the number of physical CPU cores present in the compute node. This value may exceed 100% if you are using hyper-threading, if the cores are *oversubscribed*, or if other system processes and daemons start running and take CPU resources away from your program. A value consistently less than 100% may indicate your program is not taking full advantage of the CPU resources available on a compute node.

MPI breakdown

This section organizes the time spent in MPI calls reported in the summary. It is only interesting if the program spends a significant amount of its time in MPI calls.

All the rates quoted here are inbound + outbound rates. This means that the rate of communication is being measured from the process to the MPI API, not of the underlying hardware directly.

This application-perspective is found throughout Arm Performance Reports, and in this case allows the results to capture effects such as faster intra-node performance, zero-copy transfers, and other effects.

Note: For programs that make MPI calls from multiple threads (MPI is in MPI_THREAD_SERIALIZED or MPI_THREAD_MULTIPLE mode), Arm Performance Reports only displays metrics for MPI calls made on the main thread.

Time in collective calls

The percentage of time spent in collective MPI operations such as MPI_Scatter, MPI_Reduce and MPI Barrier.

Time in point-to-point calls

The percentage of time spent in point-to-point MPI operations such as MPI_Send and MPI_Recv.

Effective process collective rate

The average transfer per-process rate during collective operations, from the perspective of the application code and not the transfer layer. For example, an MPI_Alltoall that takes 1 second to send 10 Mb to 50 processes and receive 10 Mb from 50 processes has an effective transfer rate of 10x50x2 = 1000 Mb/s.

Collective rates can often be higher than the peak point-to-point rate if the network topology matches the application's communication patterns well.

Effective process point-to-point rate

The average per-process transfer rate during point-to-point operations, from the perspective of the application code and not the transfer layer. Asynchronous calls that allow the application to overlap communication and computation, such as MPI_ISend can achieve much higher effective transfer rates than synchronous calls.

Overlapping communication and computation is often a good strategy to improve application performance and scalability.

I/O breakdown

This section organizes the amount of time spent in library and system calls relating to I/O, such as read, write and close. I/O that is generated by MPI network traffic is not included. In most cases, this should be a direct measure of the amount of time spent reading and writing to the filesystem, whether local or networked.

Some systems, such as the Cray X-series, do not have I/O accounting enabled for all filesystems. On these systems only Lustre I/O is reported in this section.

Even if your application does not perform I/O, a non-zero amount of I/O will be reported because of internal I/O performed by Arm Performance Reports.

Time in reads

The percentage of time spent on average in read operations from the perspective of the application, not the filesystem. Time spent in the stat system call is also included here.

Time in writes

The percentage of time spent on average in write and sync operations from the perspective of the application, not the filesystem.

Opening and closing files is also included here, because measurements have shown that the latest networked filesystems can spend significant amounts of time opening files with create or write permissions.

Effective process read rate

The average transfer rate during read operations from the perspective of the application. A cached read has a much higher read rate than one that has to hit a physical disk. This is particularly important to optimize for because current clusters often have complex storage hierarchies with multiple levels of caching.

Effective process write rate

The average transfer rate during write and sync operations from the application's perspective. A buffered write will have a much higher write rate than one that has to hit a physical disk. However, unless there is significant time between writing and closing the file, the penalty will be paid during the synchronous close operation instead. All these complexities are captured in this measurement.

Lustre metrics

Lustre metrics are enabled if your compute nodes have one or more Lustre filesystems mounted. Lustre metrics are obtained from a Lustre client process that runs on each node. Therefore, the data gives the information gathered on a per-node basis. The data is also cumulative over all of the processes run on a node, not only the application being profiled. Consequently, there might be some data reported to be read and written, even if the application itself does not perform file I/O through Lustre.

However, an assumption is made that the majority of data that is read and written through the Lustre client will be from an I/O intensive application, not from background processes. This assumption has been observed to be reasonable. For generated application profiles with more than a few megabytes of data that is read or written, almost all of the data reported in Arm Performance Reports is attributed to the application being profiled.

The data that is gathered from the Lustre client process is the read and write rate of data to Lustre, and a count of some metadata operations. Lustre does not just store pure data, but associates this data with metadata, which describes where data is stored on the parallel file system and how to access it. This metadata is stored separately from data, and needs to be accessed whenever new files are opened, closed, or files are resized. Metadata operations consume time and add to the latency in accessing the data. Therefore, frequent metadata operations can slow down the performance of I/O to Lustre.

Arm Performance Reports reports on the total number of metadata operations, and also the total number of file opens that are encountered by a Lustre client. With the information provided in Arm Performance Reports you can observe the rate at which data is read and written to Lustre through the Lustre client,

and also identify whether a slow read or write rate can be correlated to a high rate of expensive metadata operations.

Notes:

- For jobs run on multiple nodes, the reported values are the mean across the nodes.
- If you have more than one Lustre filesystem mounted on the compute nodes, the values are summed across all Lustre filesystems.
- Metadata metrics are only available if you have the Advanced Metrics Pack add-on for Arm Performance Reports.

Lustre read transfer: The number of bytes read per second from Lustre.

Lustre write transfer: The number of bytes written per second to Lustre.

Lustre file opens: The number of file open operations per second on a Lustre filesystem.

Lustre metadata operations: The number of metadata operations per second on a Lustre filesystem. Metadata operations include file open, close and create as well as operations such as readdir, rename, and unlink.

Note: Depending on the circumstances and implementation, 'file open' might count as multiple operations, for example, when it creates a new file or truncates an existing one.

Memory breakdown

Unlike the other sections, the memory section does not refer to one particular portion of the job. Instead, it summarizes memory usage across all processes and nodes over the entire duration. All of these metrics refer to RSS, meaning physical RAM usage, and not virtual memory usage. Most HPC jobs attempt to stay within the physical RAM of their node for performance reasons.

Mean process memory usage

The average amount of memory used per-process across the entire length of the job.

Peak process memory usage

The peak memory usage that is seen by one process at any moment during the job. If this varies a lot from the mean process memory usage, it might be a sign of either imbalanced workloads between processes or a memory leak within a process.

Note: This is not a true high-watermark, but rather the peak memory seen during statistical sampling. For most scientific codes this is not a meaningful difference because rapid allocation and deallocation of large amounts of memory is generally avoided for performance reasons.

Peak node memory usage

The peak percentage of memory that is seen being used on any single node during the entire run. If this is close to 100%, swapping might be occurring, or the job might be likely to hit hard system-imposed limits. If this is low, it might be more efficient in CPU hours to run with a smaller number of nodes and a larger workload per node.

Accelerator breakdown

Accelerators

A breakdown of how accelerators were used:

GPU utilization 47.8%

Global memory accesses 1.6%

Mean GPU memory usage 0.8%

Peak GPU memory usage 0.8%

GPU utilization is low; identify CPU bottlenecks with a profiler and offload them to the accelerator.

The peak GPU memory usage is low. It may be more efficient to offload a larger portion of the dataset to each device.

Figure 5: *Accelerator metrics report*

This section shows the utilization of NVIDIA CUDA accelerators by the job.

GPU utilization

The average percentage of the GPU cards working when at least one CUDA kernel is running.

Global memory accesses

The average percentage of time that the GPU cards were reading or writing to global (device) memory.

Mean GPU memory usage

The average amount of memory in use on the GPU cards.

Peak GPU memory usage

The maximum amount of memory in use on the GPU cards.

Energy breakdown

Energy A breakdown of how the 1.77 Wh was used: CPU 55.7% System 44.3% Mean node power 106 W Peak node power 107 W Significant time is spent on memory accesses. Reducing the CPU clock frequency could reduce the total energy usage.

Figure 6: Energy metrics report

This section shows the energy used by the job, organized by component, such as CPU and accelerators.

CPU

The percentage of the total energy used by the CPUs.

CPU power measurement requires an Intel CPU with RAPL support, for example Sandy Bridge or newer, and the intel_rapl powercap kernel module to be loaded.

Accelerator

The percentage of energy used by the accelerators. This metric is only shown when a CUDA card is present.

System

The percentage of energy used by other components not shown above. If CPU and accelerator metrics are not available, the system energy will be 100%.

Mean node power

The average of the mean power consumption of all the nodes in Watts.

Peak node power

The node with the highest peak of power consumption in Watts.

Requirements

CPU power measurement requires an Intel CPU with RAPL support, for example Sandy Bridge or newer, and the intel_rapl powercap kernel module to be loaded.

Node power monitoring is implemented through one of two methods: the Arm IPMI energy agent which can read IPMI power sensors, or the Cray HSS energy counters.

For more information on how to install the Arm IPMI energy agent please see G.4 Arm IPMI Energy Agent. The Cray HSS energy counters are known to be available on Cray XK6 and XC30 machines.

Accelerator power measurement requires a NVIDIA GPU that supports power monitoring. This can be checked on the command-line with nvidia-smi -q -d power. If the reported power values are reported as "N/A", power monitoring is not supported.

Textual performance reports

The same information is presented as described in 6.1 HTML performance reports, but in a format better suited to automatic data extraction and reading from a terminal:

Command: mpiexec -n 16 examples/wave_c 60

Resources: 1 node (12 physical, 24 logical cores per node, 2

GPUs per node available)

Memory: 15 GB per node, 11 GB per GPU

Tasks: 16 processes

Machine: node042

Started on: Tue Feb 25 12:14:06 2014 Total time: 60 seconds (1 minute)

Full path: /global/users/mark/arm/reports/examples

Notes:

Summary: wave_c is compute-bound in this configuration

Compute: 82.4% |======|

MPI: 17.6% |=| 1/0: 0.0% |

This application run was compute-bound. A breakdown of this time and advice for investigating further is found in the compute section below.

Because minimal time is spent in MPI calls, this code might also benefit from running at larger scales.

. . .

A combination of grep and sed can be useful for extracting and comparing values between multiple runs, or for automatically placing this data into a centralized database.

CSV performance reports

A CSV (comma-separated values) output file can be generated using the --output argument and specifying a filename with the .CSV extension:

```
perf-report --output=myFile.csv ...
```

The CSV file will contain lines in a NAME, VALUE format for each of the reported fields. This is convenient for passing to an automated analysis tool, such as a plotting program. It can also be imported into a spreadsheet for analyzing values among executions.

Worked examples

The best way to understand how to use and interpret performance reports is by example. You can download several sets of real-world reports with analysis and commentary from the Arm Developer website.

At the time of writing there are three collections available, which are described in the following sections.

Code characterization and run size comparison

A set of runs from well-known HPC codes at different scales showing different problems:

Characterization of HPC codes and problems

Deeper CPU metric analysis

A look at the impact of hyper-threading on the performance of a code as seen through the CPU instructions breakdown:

Exploring hyperthreading

I/O performance bottlenecks

The open source MAD-bench I/O benchmark is run in several different configurations, including on a laptop, and the performance implications are analyzed:

Understanding I/O behavior

Configurable Perf metrics

The Perf metrics use the Linux kernel perf_event_open() system call to provide additional CPU related metrics available for Performance Reports. They can be used on any system supported by the Linux perf command (also called perf_event). These cannot be tracked on typical virtual machines.

Perf metrics count the rate of one or more performance events occurring in a program. There are some software events provided by the Linux kernel but most are hardware events tracked by the Performance Monitoring Unit (PMU) of the CPU. *Generalized* hardware events are event name aliases that the Linux kernel identifies.

The quantity and combinations (in some cases) of events that can be simultaneously tracked is limited by the hardware. This feature does not support multiplexing performance events.

If the set of events you requested can not be tracked at the same time, Performance Reports ends the profiling session immediately with an error message. Try requesting fewer events, or a different combination. See the PMU reference manual for your architecture for more information on incompatible events.

Permissions

On some systems, using the Perf hardware counters can be restricted by the value of /proc/sys/k-ernel/perf_event_paranoid.

perf_event_paranoid	Description
3	Disable use of Perf events
2	Allow only user-space measurements
1	Allow kernel and user-space measurements
0	Allow access to CPU-specific data but not raw trace-point sam-
	ples.
-1	No restrictions

The value of /proc/sys/kernel/perf_event_paranoid must be 2 or lower to collect Perf metrics. To set this until the next reboot, run the following commands:

sudo sysctl -w kernel.perf_event_paranoid=2

To permanently set the paranoid level, add the following line to: /etc/sysctl.conf.

kernel.perf_event_paranoid=2

Probing target hosts

You must probe an example of a typical host machine before using these metrics. As well as other properties, this collects the CPU ID used to identify the set of potential hardware events for the host, and tests which generalized events are supported.

Ensure that /proc/sys/kernel/perf_event_paranoid is set to 2 or lower (Permissions) before performing the probe.

Note: It is not necessary to probe every potential host, a single compute node in a homogeneous cluster is sufficient.

If your home directory is writable you can generate a probe file and install it in your config directory by running the following on the intended host:

```
/path/to/forge/bin/forge-probe --install=user
```

If the Forge installation directory is writable, you can generate and install the probe file for the current host with:

```
/path/to/forge/bin/forge-probe --install=global
```

To generate the probe file, but install it manually, execute:

```
/path/to/forge/bin/forge-probe
```

The probe is named <hostname>_probe.json and is generated in your current working directory. You must manually copy it to the location specified in the forge-probe output. This is typically only necessary when the compute node that you are probing does not have write access to your home file system.

Check that the expected probe files are correctly installed with:

```
/path/to/forge/bin/map --target-host=list
```

This shows something like:

```
0x0000000420f5160 (thunderx2) e.g. node07.myarmhost.com
GenuineIntel-6-4E (skylake) e.g. node01.myintelhost.com
```

If you have exactly one probe file installed, this is automatically assumed to be the target host. If there are multiple installed probe files, you must specify the intended target whenever you use the configurable Perf metrics feature. When using the command line, use the --target-host argument. You can specify the intended target CPU ID (such as, 0x00000000420f5160), family name (such as, thunderx2), or a unique substring of the hostname (myarmhost).

Specifying Perf metrics via the command line

You can list available events for a given probed host using:

```
/path/to/forge/bin/perf-report map --target-host=myarmhost \
    --perf-metrics=avail
```

Note: Use list instead of avail to see the events listed on separate lines.

Specify the events you want using a semicolon separated list:

```
/path/to/forge/bin/perf-report --profile --target-host=myarmhost \
    --perf-metrics="cpu-cycles; bus-cycles; instructions" mpirun ...
```

Specifying Perf metrics via a file

The --perf-metrics argument can also take the name of a plain text file:

```
/path/to/forge/bin/perf-report --profile --target-host=myhost \
    --perf-metrics=./myevents.txt mpirun ...
```

myevents.txt lists the events to track on separate lines, such as:

```
cpu-cycles
bus-cycles
instructions
```

--perf-metrics=template outputs a more complex template that lists all possible events with accompanying descriptions. Redirect this output to a file and uncomment the events to track, for example:

```
/path/to/forge/bin/perf-report --target-host=myhost \
    --perf-metrics=template > myevents.txt

vim myevents.txt

/path/to/forge/bin/perf-report --profile \
    --perf-metrics=myevents.txt mpirun ...
```

Viewing events

You can view Perf event counts in the **CPU Metrics** section. All these metrics are reported as events per second with a suitable SI prefix (such as, K, M, G) that is automatically determined.

The default values that are reported are the mean of means:

- 1. The mean value is taken across all processes for each sample (averaging across processes).
- 2. The mean value is taken of those per-sample results (averaging across time).

Advanced configuration

You can override the default settings used by Performance Reports when making perf_event_open calls. Specify one or more flags in a preamble section in square brackets at the start of the perf metrics definition string (either on the command line or at the top of a template file).

```
/path/to/forge/bin/perf-report --profile --target-host=myarmhost \
    --perf-metrics="[optional,noinherit]; instructions; cpu-cycles"
```

Possible options are:

- **[optional]:** Do not abort the program if the requested metrics cannot be collected. Set this if you wish to continue profiling even if the no Perf metric results is returned.
- [noinherit]: Disable multithreading support (new threads will not inherit the event counter configuration). If you specified events, they are only collected on the main thread (in the case of MPI programs, the thread that called MPI_thread_init).
- **[nopinned]:** Disable pinning events on the PMU. If you have specified this, event counting might be multiplexed. Arm does not recommend doing this as it interacts poorly with the Forge sampling strategy.
- [noexclude=kernel]: Do not exclude kernel events that happen in kernel space. This might require a more permissive perf_event_paranoid level.
- **[noexclude=hv]:** Do not exclude events that happen in the hypervisor. This is mainly for PMUs that have built-in support for handling this (such as IBM Power). Most machines require extra support for handling hypervisor measurements.

• **[noexclude=idle]:** Do not exclude counting software events when the CPU is running the idle task. This is only relevant for software events.

Configuration

Arm Performance Reports generally requires no configuration before use.

If you only intend to use Arm Performance Reports and have checked that it works on your system without extra setup then you can safely ignore the rest of this section.

Compute node access

When Arm Performance Reports needs to access another machine as part of starting one of *MPICH 1*–3, *Intel MPI*, and *SGI MPT*, it attempts to use the secure shell, ssh, by default.

However, this may not always be appropriate, SSh may be disabled or be running on a different port to the normal port 22. In this case, you can create a file called remote-exec which is placed in your ~/.allinea directory and Arm Performance Reports will use this instead.

Arm Performance Reports checks for the script at ~/.allinea/remote-exec, and it will be executed as follows:

```
remote-exec HOSTNAME APPNAME [ARG1] [ARG2] ...
```

The script should start APPNAME on HOSTNAME with the arguments ARG1 ARG2 without further input (no password prompts). Standard output from APPNAME will appear on the standard output of remote-exec.

SSH based remote-exec

A remote-exec script using ssh running on a non-standard port could look as follows:

```
#!/bin/sh
ssh -P {port-number} $*
```

In order for this to work without prompting for a password, you should generate a public and private SSH key, and ensure that the public key has been added to the ~/.ssh/authorized_keys file on machines you wish to use.

See the ssh-keygen manual page for more information.

Testing

Once you have set up your remote-exec script, it is recommended that you test it from the command line. For example:

```
~/.allinea/remote-exec TESTHOST uname -n
```

This returns the output of uname -n on TESTHOST, without prompting for a password.

If you are having trouble setting up remote-exec, please contact Arm support at Arm support for assistance.

Windows

The functionality described above is also provided by the Windows remote client. There are two differences:

- The script is named remote-exec.cmd rather than remote-exec.
- The default implementation uses the plink.exe executable supplied with Arm Performance Reports.

Getting support

While this document attempts to cover as many parts of the installation, features and uses of our tool as possible, there will be scenarios or configurations that are not covered, or are only briefly mentioned, or you may on occasion experience a problem using the product.

You can contact Arm support at Arm support.

Please provide as much detail as you can about the scenario in hand, such as:

• Version number of Arm Performance Reports (for example, perf-report --version) and your operating system and the distribution, for example Red Hat Enterprise Linux 6.4.

This information is all available by using the --version option:

```
bash$ perf-report --version
```

Arm Performance Reports

Copyright (c) 2002-2019 Arm Limited (or its affiliates). All rights reserved.

Version: 18.0.2

Build ID: 556f23c4895e

Build Platform: Ubuntu 16.04 x86_64 Build Date: Jan 25 2018 22:42:00

Frontend OS: Ubuntu 16.04.2 LTS

Nodes' OS: unknown

Last connected ddt-debugger: unknown

- The compiler used and its version number.
- The MPI library and version if appropriate.
- A description of the issue: what you expected to happen and what actually happened.
- An exact copy of any warning or error messages that you may have encountered.

Supported platforms

This section describes the architectures, operating systems, MPI distributions, compilers and accelerators that Performance Reports supports.

Performance Reports

Architecture	Operating sys-	Toolkits	Compilers	Accelerators
	tems			
Intel and	Red Hat Enterprise	Open MPI 1.8, 1.10, 2.0,	GNU C/C++/Fortran	Nvidia
AMD	Linux/CentOS 7	2.1, 3.0, and 3.1	Compiler 4.3.x to 8.1.x	CUDA
(x86_64)	SuSE Linux Enter-	MPICH 3.1 to 3.2	LLVM Clang 3.3 to 6.0	Toolkit 8.0
	prise 11.3, 12, and	MVAPICH2 2.0 to 2.3	Intel Parallel Studio XE	to 9.2
	15	Intel MPI 5.1.x to 2019.5	2013.x to 2019.5	
	Ubuntu 14.04,	Cray MPT 6.3.1 to 7.7.1	PGI Compiler 14.0 to	
	16.04, and 18.04	SGI MPT 2.10 to 2.15	18.7	
	Open SuSE 12, 13,	HPE MPI 1.1	Cray Compiling Environ-	
	42.3, and 15.0	IBM Platform MPI 9.1.x	ment 8.3.x to 8.7.x	
		Bullx MPI 1.2.7.x to	Nvidia CUDA Compiler	
		1.2.9.x	8.0 to 9.2	
Armv8	Red Hat Enterprise	Open MPI 1.8, 1.10, 2.0,	GNU C/C++/Fortran	_
(AArch64)	Linux/CentOS 7	2.1, 3.0, and 3.1	Compiler 4.3.x to 8.1.x	
	SuSE Linux Enter-	MPICH 3.1 to 3.2	LLVM Clang 3.3 to 6.0	
	prise 12 and 15	MVAPICH2 2.0 to 2.3	Arm C/C++/Fortran	
	Ubuntu 16.04 and	Cray MPT 7.7.1	Compiler 18.0 to 19.2	
	18.04	HPE MPI 1.1	Cray Compiling Environ-	
		Bullx MPI 1.2.7.x to	ment 8.7.x	
		1.2.9.x		
Intel Xeon	Red Hat Enterprise	Open MPI 1.8, 1.10, 2.0,	GNU C/C++/Fortran	_
Phi (knl)	Linux/CentOS 7	2.1, 3.0, and 3.1	Compiler 4.3.x to 8.1.x	
	SuSE Linux Enter-	Intel MPI 5.1.x to 2019.5	LLVM Clang 3.3 to 6.0	
	prise 12	MPICH 3.1 to 3.2	Intel Parallel Studio XE	
		MVAPICH2 2.0 to 2.3	2013.x to 2019.5	
		Cray MPT 5.0.x to 7.7.x	PGI Compiler 14.0 to	
		HPE MPI 1.1	18.5	
		Bullx MPI 1.2.7.x to	Cray Compiling Environ-	
		1.2.9.x	ment 8.3.x to 8.7.x	
IBM Power	Red Hat Enterprise	Open MPI 1.8, 1.10, 2.0,	GNU C/C++/Fortran	_
(ppc64le)	Linux/CentOS	2.1, 3.0, and 3.1	Compiler 4.3.x to 8.1.x	
	7.2+	Spectrum MPI 10.2	IBM XL C/C++ Compiler	
			13.1.x	
			IBM XL Fortran Com-	
			piler 15.1.x	
			IBM XL Compiler 16.1.x	
			PGI Compiler 18.1, 18.5,	
			and 18.7	

The Arm profiling libraries must be explicitly linked with statically linked programs which mostly applies

to the Cray X-Series.

Batch schedulers: SLURM 2.6.3+ and 14.03+ (srun only).

Known issues

The most significant known issues for the latest release are summarized here:

- I/O metrics are not fully available on some systems, including Cray systems.
- CPU instruction metrics are only available on x86_64 systems.
- Thread activity is not sampled whilst a process is inside an MPI call with a duration spanning multiple samples.
- Xeon Phi systems using many MPI ranks per card should set ALLINEA_REDUCE_MEMORY_-USAGE=1.
- Version 14.9 or later of the PGI compilers is required to compile the Arm Performance Reports MPI wrappers as a static library.
- MPICH 3.0.3 and 3.0.4 do not work with MAP, DDT or Performance Reports due to a defect in MPICH 3.0.3/4. MPICH 3.1 addressed this and is fully supported.
- Open MPI 2.1.3 works with Arm Performance Reports. Previous versions of Open MPI 2.1.x do not work due to a bug in the Open MPI debug interface.
- On Cray X-series systems only *native* SLURM is supported, *hybrid* mode is not supported.
- Performance Reports may fail to finalize a profiling session if the cores are oversubscribed on AArch64 architectures. For example when attempting to profile a 64 process MPI program on a machine with only 8 cores. This will appear as a hang after finishing a profile.

See also additional known issues here:

Category	Known Issues
MPI Distribution	D MPI distribution notes
Compiler	E Compiler notes
Platform	F Platform notes
General	G General troubleshooting

MPI distribution notes

This appendix has brief notes on many of the MPI distributions supported by Arm Performance Reports.

Advice on settings and problems particular to a distribution are given here.

Bull MPI

Bull X-MPI is supported.

Cray MPT

Arm Performance Reports users may wish to read 4.1.3 Static linking on Cray X-Series Systems.

Arm Performance Reports has been tested with Cray XK7 and XC30 systems.

Arm Performance Reports requires Arm's sampling libraries to be linked with the application before running on this platform.

See 4.1.1 Linking for a set-by-step guide.

Arm supplies module files in {Reports-installation-directory}/share/modules/cray.

See 4.1.5 Dynamic and static linking on Cray X-Series systems using the modules environment to simplify linking with the sampling libraries.

Known Issues:

• By default scripts wrapping Cray MPT will not be detected, but you can force the detection by setting the ALLINEA_DETECT_APRUN_VERSION environment variable to "yes" before starting Performance Reports.

Intel MPI

Arm Performance Reports has been tested with Intel MPI 4.1.x, 5.0.x and onwards.

Known Issue: If you use Spectrum LSF as workload manager in combination with Intel MPI and you get for example one of the following errors:

- < target program > exited before it finished starting up. One or more processes were killed or died without warning
- < target program > encountered an error before it initialised the MPI environment. Thread 0 terminated with signal SIGKILL

or the job is killed otherwise during launching then you may need to set/export I_MPI_LSF_USE_ COLLECTIVE_LAUNCH=1 before executing Arm Performance Reports. See Using IntelMPI under LSF quick guide and Resolve the problem of the Intel MPI job ...hang in the cluster for more details.

MPICH 2

If you see the error undefined reference to MPI_Status_c2f during initialization or if manually building the sampling libraries as described in 4.1.1 Linking, then you need to rebuild MPICH 2 with Fortran support.

MPICH 3

MPICH 3.0.3 and 3.0.4 do not work with Arm Performance Reports due to an MPICH bug. MPICH 3.1 addresses this and is supported.

Open MPI

Arm Performance Reports products have been tested with Open MPI 1.6.x, 1.8.x, 1.10.x, 2.0.x, and 3.0.x.

The following versions of Open MPI do not work with Arm Performance Reports because of bugs in the Open MPI debug interface:

- Open MPI 2.1.0 to 2.1.2.
- Open MPI 3.0.0 when compiled with the Arm Compiler for Linux on Arm®v8 (AArch64) systems.
- Open MPI 3.0.x when compiled with some versions of the GNU compiler on Arm®v8 (AArch64) systems.
- Open MPI <= 3.x.4 and <= 4.0.1 when compiled with some versions of IBM XLC/XLF or PGI compilers.
- Open MPI 3.1.0 and 3.1.1.
- Open MPI 3.x with any version of PMIx < 2.
- Open MPI 4.0.1 with PMIx 3.1.2.

To resolve any of the above issues, instead select *Open MPI (Compatibility)* for the MPI Implementation.

Open MPI 3.x on IBM Power with the GNU compiler

To use Open MPI versions 3.0.0 to 3.0.4 (inclusive) and Open MPI versions 3.1.0 to 3.1.3 (inclusive) with the GNU compiler on IBM Power systems, you might need to configure the Open MPI build with CFLAGS=-fasynchronous-unwind-tables. Configuring the Open MPI build with CFLAGS=-fasynchronous-unwind-tables fixes a startup bug where Arm Performance Reports is unable to step out of MPI_Init into your main function. The startup bug occurs because of missing debug information and optimization in the Open MPI library. If you already configure with -g, you do not need to add this extra flag. An example configure command is:

```
./configure --prefix=/software/openmpi-3.1.2
    CFLAGS=-fasynchronous-unwind-tables
```

If you do not have the option to recompile your MPI, an alternative workaround is to select *Open MPI* (*Compatibility*) for the MPI Implementation. This issue is fixed in later versions.

Platform MPI

Platform MPI 9.x is supported, but only with the mpirun command. Currently mpiexec is not supported.

SGI MPT / SGI Altix

SGI MPT 2.10+ is supported.

Some SGI systems can not compile programs on the batch nodes, for example because the gcc package is not installed.

If this applies to your system you must explicitly compile the Arm MPI wrapper library using the make-profiler-libraries command and then explicitly link your programs against the Arm profiler and sampler libraries.

The mpio.h header file shipped with SGI MPT 2.10 contains a mismatch between the declaration of MPI_File_set_view and some other similar functions and their PMPI equivalents, for example PMPI_File_set_view. This prevents Arm Performance Reports from generating the MPI wrapper library. Please contact SGI for a fix.

If you are using SGI MPT with SLURM and would normally use mpiexec_mpt to launch your program you will need to use srun --mpi=pmi2 directly.

Preloading the Arm profiler and MPI wrapper libraries is not supported in Express Launch mode. Arm recommends you explicitly link your programs against these libraries to work around this problem. If this is not possible you can manually compile the MPI wrapper, and explicitly set LD_PRELOAD in the launch line.

SLURM

The use of the --export argument to srun (SLURM 14.11 or newer) is not supported. In this case you can avoid using --export by exporting the necessary environment variables before running Arm Performance Reports.

The use of the --task-prolog argument to srun (SLURM 14.03 or older) is also not supported, as the necessary libraries cannot be preloaded. You will either need to avoid using this argument, or explicitly link to the libraries.

Compiler notes

For a list of supported compiler versions, refer to section B.

AMD OpenCL compiler

Not supported by Arm Performance Reports.

Berkeley UPC compiler

Not supported by Arm Performance Reports.

Cray compiler environment

The Cray UPC compiler is not supported by Arm Performance Reports.

GNU

The -foptimize-sibling-calls optimization (used in -02, -03 and -0s) interfere with the detection of some OpenMP regions.

If your code is affected with this issue add -fno-optimize-sibling-calls to disable it and allow Arm Performance Reports to detect all the OpenMP regions in your code.

Debugging Python code compiled with GNU versions greater than 7.4.0 on Ubuntu 18.04 is currently unavailable. This is to be resolved in a future version of the product.

GNU UPC

Arm Performance Reports does not support this.

Intel compilers

Refer to section B for a list of supported versions.

Portland Group compilers

Arm Performance Reports has been tested with Portland Tools 14 onwards.

Platform notes

This page notes any particular issues affecting platforms. If a supported machine is not listed on this page, it is because there is no known issue.

Intel Xeon

Intel Xeon processors starting with Sandy Bridge include Running Average Power Limit (RAPL) counters. Performance Reports can use the RAPL counters to provide energy and power consumption information for your programs.

Enabling RAPL energy and power counters when profiling

To enable the RAPL counters to be read by Performance Reports you must load the intel_rapl kernel module.

The intel_rapl module is included in Linux kernel releases 3.13 and later.

For testing purposes Arm have backported the powercap and intel_rapl modules for older kernel releases. You may download the backported modules from:

Download backported modules

Note: These backported modules are unsupported and should be used for testing purposes only. No support is provided by Arm, your system vendor or the Linux kernel team for the backported modules.

NVIDIA CUDA

- CUDA metrics are not available for statically-linked programs.
- CUDA metrics are measured at the node level, not the card level.

Arm

Arm®v8 (AArch64) known issues

• Performance Reports may fail to finalize a profiling session if the cores are oversubscribed on AArch64 platforms. For example, this issue is likely to occur when attempting to profile a 64 process MPI program on a machine with only 8 cores. This issue appears as a hang after finishing a profile.

POWER8 and POWER9 (POWER 64-bit) known issues

Hardware watch points do not function on some POWER9 systems. Check with IBM for compatibility.

General troubleshooting

If you have problems with any of the Arm Performance Reports products, please read this section carefully.

Additionally, check the support pages on the Arm Developer website, and make sure you have the latest version of the product.

Starting a program

Problems starting scalar programs

There are a number of possible sources for problems. The most common for users with a multi-process license is that the *Run Without MPI Support* check box has not been checked.

If the software reports a problem with MPI and you know your program is not using MPI, then this is usually the cause.

If you have checked this box and the software still mentions MPI then please contact Arm support at Arm support.

Other potential problems are:

- A previous session is still running, or has not released resources required for the new session.
 Usually this can be resolved by killing stale processes. The most obvious symptom of this is a
 delay of approximately 60 seconds and a message stating that not all processes connected. You
 may also see a QServerSocket message in the terminal.
- The target program does not exist or is not executable.
- Arm Performance Reports products' backend daemon, ddt-debugger, is missing from the bin directory. In this case you should check your installation, and contact Arm support at Arm support for further assistance.

Problems starting multi-process programs

If you encounter problems while starting an MPI program, the first step is to establish that it is possible to run a single-process (non-MPI) program such as a trivial "Hello, World!", and resolve such issues that may arise. After this, attempt to run a multi-process job and the symptoms will often allow a reasonable diagnosis to be made.

In the first instance verify that MPI is working correctly by running a job, without Arm Performance Reports products applied, such as the example in the examples directory.

Verify that mpirun is in the PATH, or the environment variable ALLINEA_MPIRUN is set to the full pathname of mpirun.

Sometimes problems are caused by environment variables not propagating to the remote nodes while starting a job. The solution to these problems depend on the MPI implementation that is being used.

In the simplest case, for rsh-based systems such as a default MPICH 1 installation, correct configuration can be verified by rsh-ing to a node and examining the environment. It is worthwhile rsh-ing with the env command to the node as this will not see any environment variables set inside the .profile command.

For example, if your nodes use a .profile instead of a .bashrc for each user then you may well see a different output when running rsh node env than when you run rsh node and then run env inside the new shell.

If only one, or very few, processes connect, it may be because you have not chosen the correct MPI implementation. Please examine the list and look carefully at the options. Should no other suitable MPI be found, please contact Arm support at Arm support.

If a large number of processes are reported by the status bar to have connected, then it is possible that some have failed to start due to resource exhaustion, timing out, or, unusually, an unexplained crash. You should verify again that MPI is still working, as some MPI distributions do not release all semaphore resources correctly, for example MPICH 1 on Redhat with SMP support built in.

To check for time-out problems, set the ALLINEA_NO_TIMEOUT environment variable to 1 before launching the GUI and see if further progress is made. This is not a solution, but aids the diagnosis. If all processes now start, please contact Arm support at Arm support for further advice.

No shared home directory

If your home directory is not accessible by all the nodes in your cluster then your jobs may fail to start.

To resolve the problem open the file ~/.allinea/system.config in a text editor. Change the shared directory option in the [startup] section so it points to a directory that is available and shared by all the nodes. If no such directory exists, change the use session cookies option to no instead.

Performance Reports specific issues

MPI wrapper libraries

Arm Performance Reports wrap MPI calls in a custom shared library. One is built for your system each time you run Arm Performance Reports.

If this does not work please contact Arm support at Arm support.

You can also try setting MPICC directly:

\$ MPICC=my-mpicc-command bin/perf-report --np=16 ./wave_c

Thread support limitations

Performance Reports provides limited support for programs when threading support is set to MPI_THREAD_SERIALIZED or MPI_THREAD_MULTIPLE in the call to MPI_Init_thread.

MPI activity on non-main threads will contribute towards the MPI-time of the program, but not the more detailed MPI metrics.

MPI activity on a non-main thread may result in additional profiling overhead due to the mechanism employed by Performance Reports for detecting MPI activity.

Warnings are displayed when the user initiates and completes profiling a program which sets MPI_THREAD_SERIALIZED or MPI_THREAD_MULTIPLE as the required thread support.

Performance Reports does support calling MPI_Init_thread with either MPI_THREAD_SINGLE or MPI_THREAD_FUNNELED specified as the required thread support.

It should be noted that the requirements that the MPI specification make on programs using MPI_THREAD_FUNNELED are the same as made by Performance Reports: all MPI calls must be made on the thread that called MPI_Init_thread.

In many cases, multi-threaded MPI programs can be refactored such that they comply with this restriction.

No thread activity while blocking on an MPI call

Unfortunately Arm Performance Reports is currently unable to record thread activity on a process where a long duration MPI call is in progress.

If you have an MPI call that takes a significant amount of time (multiple samples) to complete then Arm Performance Reports will record no thread activity for the process executing that call for most of that MPI call's duration.

I am not getting enough samples

By default sampling takes place every 20ms initially, but if you get warnings about too few samples on a fast run, or want more detail in the results, you can change that.

To increase the frequency of sampling to every 10ms set environment variable ALLINEA_SAMPLER_INTERVAL=10.

Notes:

- Sampling frequency is automatically decreased over time to ensure a manageable amount of data is collected whatever the length of the run.
- Whether OpenMP is enabled or disabled in Arm MAP, the final script or scheduler values set for OMP_NUM_THREADS will be used to calculate the sampling interval per thread. When configuring your job for submission, check whether your final submission script, scheduler or the Arm MAP GUI has a default value for OMP_NUM_THREADS.
- Custom values for ALLINEA_SAMPLER_INTERVAL will be overwritten by values set from the combination of ALLINEA_SAMPLER_INTERVAL_PER_THREAD and the expected number of threads (from OMP_NUM_THREADS).

Increasing the sampling frequency is not recommended if there are lots of threads or there are deep stacks in the target program as this may not leave sufficient time to complete one sample before the next sample is started.

Performance Reports is reporting time spent in a function definition

Any overheads involved in setting up a function call (pushing arguments to the stack and so on) are usually assigned to the function definition.

Some compilers may assign them to the opening brace '{' and closing brace '}' instead. If this function has been inlined, the situation becomes further complicated and any setup time (for example, allocating space for arrays) is often assigned to the definition line of the enclosing function.

Performance Reports is not correctly identifying vectorized instructions

The instructions identified as vectorized (packed) are enumerated below. Arm also identifies the AVX-2 variants of these instructions (with a "V" prefix).

Contact Arm support at Arm support if you believe your code contains vectorized instructions that have not been listed and are not being identified in the *CPU floating-point/integer vector* metrics.

Packed floating-point instructions: addpd addps addsubpd addsubps andnpd andnps andpd andps divpd divps dppd dpps haddpd haddps hsubpd hsubps maxpd maxps minpd minps mulpd mulps rcpps rsqrtps sqrtpd sqrtps subpd subps

Packed integer instructions: mpsadbw pabsb pabsd pabsw paddb paddd paddd paddsb paddsw paddusw paddw palignr pavgb pavgw phaddd phaddsw phaddw phminposuw phsubd phsubsw phsubw pmaddubsw pmaddwd pmaxsb pmaxsd pmaxsw pmaxub pmaxud pmaxuw pminsb pminsd pminsw pminub pminud pminuw pmuldq pmulhrsw pmulhuw pmulhw pmulld pmullw pmuludq pshufb pshufw psignb psignd psignw pslld psllq psllw psrad psraw psrld psrlq psrlw psubb psubd psubd psubsb psubsw psubusb psubusw psubw

Performance Reports takes a long time to gather and analyze my OpenBLASlinked application

OpenBLAS versions 0.2.8 and earlier incorrectly stripped symbols from the .symtab section of the library, causing binary analysis tools such as Arm Performance Reports and objdump to see invalid function lengths and addresses.

This causes Arm Performance Reports to take an extremely long time disassembling and analyzing apparently overlapping functions containing millions of instructions.

A fix for this was accepted into the OpenBLAS codebase on October 8th 2013 and versions 0.2.9 and above should not be affected.

To work around this problem without updating OpenBLAS, simply run "strip libopenblas*.so"—this removes the incomplete . Symtab section without affecting the operation or linkage of the library.

Performance Reports over-reports MPI, I/O, accelerator or synchronization time

Arm Performance Reports employs a heuristic to determine which function calls should be considered as MPI operations.

If your code defines any function that starts with MPI_(case insensitive) those functions will be treated as part of the MPI library resulting in the time spent in MPI calls to be over-reported.

Starting your functions names with the prefix MPI_ should be avoided and is in fact explicitly forbidden by the MPI specification (page 19 sections 2.6.2 and 2.6.3 of the MPI 3 specification document http://www.mpi-forum.org/docs/mpi-3.0/mpi30-report.pdf#page=49):

All MPI names have an MPI_ prefix, and all characters are capitals. Programs must not declare names, for example, for variables, subroutines, functions, parameters, derived types, abstract interfaces, or modules, beginning with the prefix MPI_.

Similarly Arm Performance Reports categorizes I/O functions and accelerator functions by name.

Other prefixes to avoid starting your function names with include PMPI_, _PMI_, OMPI_, omp_-, GOMP_, shmem_, cuda_, __cuda, cu[A-Z][a-z] and allinea_. All of these prefixes are case-insensitive.

Also avoid naming a function start_pes or any name also used by a standard I/O or synchronisation function (write, open, pthread_join, sem_wait etc).

Obtaining support

To receive additional support, contact Arm support at Arm support with a detailed report of the problem you are having.

If possible, you should obtain a log file for the problem and contact Arm support at Arm support. When describing your issue, state that you have obtained a log file and the support team will be in contact.

To generate a log file, start Performance Reports with the --debug and --log arguments:

Where <log> is the name of the log file to generate.

Next, reproduce the problem using as few processors as possible. Once finished, close the program as usual.

On some systems this file may be quite large. If so, please compress it using a program such as gzip or bzip2 before sending it to support.

If your problem can only be replicated on large process counts, then omit the --debug argument as this will generate very large log files.

Arm IPMI Energy Agent

The Arm IPMI Energy Agent allows Arm MAP and Arm Performance Reports to measure the total energy consumed by the compute nodes in a job.

Note: Measuring energy with IPMI Energy agent requires Arm Forge Professional.

The IPMI Energy Agent is available to download from our website: IPMI Energy Agent.

Requirements

- The compute nodes must support IPMI.
- The compute nodes must have an IPMI exposed power sensor.
- The compute nodes must have an OpenIPMI compatible kernel module installed, such as ipmi_-devintf.
- The compute nodes must have the corresponding device node in /dev, for example /dev/ipmi0.
- The compute nodes must run a supported operating system.
- The IPMI Energy Agent must be run as root.

To list the names of possible IPMI power sensors on a compute node use the following command:

ipmitool sdr | grep 'Watts'

Index

Arm IPMI Energy Agent, 64	CPU, 42
Requirements, 64	Mean node power, 42
MAP file, 30	Peak node power, 42
Performance Reports	System, 42
Specific issues, 60	Energy metrics
•	Requirements, 43
Accelerator breakdown, 41	Environment variables, 12
Global memory accesses, 41	Example, 17
GPU utilization, 41	Compiling, 17
Mean GPU memory usage, 41	Cray, 17
Peak GPU memory usage, 41	Generating a performance report, 19
AMD	Overview, 17
OpenCL, 57	Running, 18
Arm, 51, 58	Express Launch, 26
Known issues, 58	Compatible MPIs, 26
D. II ACDI E 4	•
Bull MPI, 54	General Troubleshooting, 59
Compatibility Launch, 26	Generating a report, 27
Compute node access, 49	Getting Support, 50
Configuration, 49	GNU Compiler, 57
CPU breakdown, 33	LITMI vaparta 21
Memory accesses, 34	HTML reports, 31
OpenMP code, 34	I/O breakdown, 38
Scalar numeric ops, 34	Effective process read rate, 39
Single core code, 34	Effective process write rate, 39
Vector numeric ops, 34	Lustre metrics, 39
Waiting for accelerators, 34	Time in reads, 39
CPU metrics breakdown, 35	Time in writes, 39
Cycles per instruction, 35	Installation, 8
FLOPS scalar lower bound, 35	Linux, 8
FLOPS vector lower bound, 35	Graphical install, 8
L2 cache misses, 35	Text-mode install, 9
L3 cache miss per instruction, 35	Intel Compiler, 57
Memory accesses, 36	Intel MPI, 54
Stalled cycles, 35	Intel Xeon, 58
Cray Compiler Environment, 57	RAPL, 58
Cray MPT, 54	Interpreting, 31
Cray Native SLURM, 56	Introduction, 7
Cray X, 54	IPMI, 64
Cray X-Series, 20, 21, 24, 25	
CSV performance reports, 43	Known issues
Custom DCIM, 28	Performance Reports, 60
Custom gmetric, 29	Incorrect MPI time, 62
Casiom Sincine, 20	Insufficient samples, 61
DCIM output, 28	MPI wrapper libraries, 60
Dynamic linking	No thread activity while blocking on an MPI
Cray X-Series, 21	call, 61
•	Not correctly identifying vectorized instruc-
Enable and disable metrics, 29	tions, 62
Energy breakdown, 42	OpenBLAS application, 62

Reporting time spent in a function definition,	FLOPS vector lower bound, 35
61	Global memory accesses, 41
Thread support limitations, 60	GPU Utilization, 41
Compiler, 57	I/O breakdown, 38
General, 59	Input/Output, 33
No shared home directory, 60	L2 cache misses, 35
Problems starting multi-process programs, 59	L3 cache miss per instruction, 35
Starting a program, 59	Lustre metrics, 39
Starting scalar programs, 59	Mean GPU memory usage, 41
	Mean process memory usage, 40
Licensing	Memory accesses, 34, 36
Architecture licensing, 11	Memory breakdown, 40
Floating licenses, 10	MPI, 33
License files, 10	MPI breakdown, 37
Supercomputing and other floating licenses, 10	OpenMP breakdown, 36
Using multiple architecture licenses, 11	OpenMP code, 34
Workstation and evaluation licenses, 10	Peak GPU memory usage, 41
Linking, 20	Peak node memory usage, 40
Dynamic	Peak process memory usage, 40
On Cray X-Series using modules environ-	Physical core utilization, 36, 37
ment, 25	Scalar numeric ops, 34
Static, 22	Single core code, 34
On Cray X-Series using modules environ-	Stalled cycles, 35
ment, 25	Synchronization, 36, 37
Log file, 63	System load, 36, 37
	Threads breakdown, 37
map-link modules, 25	Time in collective calls, 38
Installation	Time in point-to-point calls, 38
Cray X-Series, 25	Time in reads, 39
Memory breakdown, 40	Time in writes, 39
Mean process memory usage, 40	Vector numeric ops, 34
Peak node memory usage, 40	Waiting for accelerators, 34
Peak process memory usage, 40	MPI
Metrics	Troubleshooting, 59
Accelerator breakdown, 41	MPI breakdown, 37
Computation, 36, 37	Effective process collective rate, 38
Compute, 33	Effective process point-to-point rate, 38
CPU breakdown, 33	Time in collective calls, 38
CPU metrics breakdown, 35	Time in conective cans, 38 Time in point-to-point calls, 38
Cycles per instruction, 35	MPI wrapper libraries, 60
Effective process collective rate, 38	
Effective process point-to-point rate, 38	MPICH 2, 55 MPICH 3, 55
Effective process read rate, 39	WFICH 5, 55
Effective process write rate, 39	NVIDIA CUDA, 58
Energy	
Accelerator, 42	Obtaining support, 63
CPU, 42	Online resources, 7
Mean node power, 42	Open MPI, 55
Peak node power, 42	OpenMP breakdown, 36
System, 42	Computation, 36
Energy breakdown, 42	Physical core utilization, 36
FLOPS scalar lower bound, 35	Synchronization, 36

System load, 36 Output locations, 28
Perf, 45
–target-host, 46
advanced configuration, 47
Command line, 46
Metrics, 45
Probe, 45
Template file, 46
Viewing, 47
perf_event_paranoid, 45
Performance reports
Energy breakdown
Accelerator, 42
Threads breakdown
Synchronization, 37
Platform MPI, 56
Portland Group Compiler, 57
POWER8 and POWER9
Known issues, 58
Profiling
Preparing a program, 20
Freparing a program, 20
Report summary, 33 Compute, 33 Input/Output, 33 MPI, 33
Requirements
Energy metrics, 43
Running, 20
SGI, 56 SLURM, 56 Static linking, 22 On Cray X-Series, 24 Supported Platforms, 51
Textual performance reports, 43
Thread support limitations, 60
Threads breakdown, 37
Computation, 37
Physical core utilization, 37
System load, 37
System roud, 57
Unified Parallel C, 57
UPC
Berkeley, 57
GNU, 57
Worked examples, 44
Code characterization and run size comparison,
44

Deeper CPU metric analysis, 44 I/O performance bottlenecks, 44