

# Pipeline Collector: gathering performance data for distributed astronomical pipelines

Alexandar P. Mechev<sup>a</sup>, Aske Plaat<sup>b</sup>, J.B. Raymond Oonk<sup>a,c</sup>, Huib T. Intema<sup>a</sup>, Huub J.A. Röttgering<sup>a</sup>

<sup>a</sup>Leiden Observatory, Niels Bohrweg 2, 2333 CA Leiden, the Netherlands

<sup>b</sup>Leiden Institute of Advanced Computer Science, Niels Bohrweg 1, 2333 CA Leiden, the Netherlands

<sup>c</sup>ASTRON, Oude Hoogeveensedijk 4, 7991 PD Dwingeloo, the Netherlands

---

## Abstract

Modern astronomical data processing requires complex software pipelines to process ever growing datasets. For radio astronomy, these pipelines have become so large that they need to be distributed across a computational cluster. This makes it difficult to monitor the performance of each pipeline step. To gain insight into the performance of each step, a performance monitoring utility needs to be integrated with the pipeline execution. In this work we have developed such a utility and integrated it with the calibration pipeline of the Low Frequency Array, LOFAR, a leading radio telescope. We tested the tool by running the pipeline on several different compute platforms and collected the performance data. Based on this data, we make well informed recommendations on future hardware and software upgrades. The aim of these upgrades is to accelerate the slowest processing steps for this LOFAR pipeline. The *pipeline\_collector* suite is open source and will be incorporated in future LOFAR pipelines to create a performance database for all LOFAR processing.

*Keywords:* Radio Astronomy, Performance Analysis, Profiling, High Performance Computing

---

## 1. Introduction

Astronomical data often requires significant processing before it is considered ready for scientific analysis. This processing is done increasingly by complex and autonomous software pipelines, often consisting of numerous processing steps, which are run without user interaction. It is necessary to collect performance statistics for each pipeline step. Doing so will enable scientists to discover and address software and hardware inefficiencies and produce scientific data at a higher rate. To identify these inefficiencies, we have extended the performance monitoring package *tcollector*<sup>1</sup> (Apache, 2017). The resulting suite, *pipeline\_collector*, makes it possible to use *tcollector* to record data for complex pipelines. We have used a leading radio telescope as the test case for the *pipeline\_collector* suite. The discoveries made with our software will help remove bottlenecks and suggest hardware requirements for current and future processing clusters. We summarize our findings in Table 1 in Section 3.

Over the past two decades, processing data in radio astronomy has increasingly moved from personal machines to large compute clusters. Over this time, radio telescopes have undergone upgrades in the form of wide band receivers and upgraded correlators (Broekema et al., 2018;

Gupta et al., 2017). In addition, several aperture synthesis arrays such as the Low Frequency Array (LOFAR, Van Haarlem et al., 2013), Murchison Widefield Array (MWA Lonsdale et al., 2009; Tingay et al., 2013) and MeerKAT (Jonas, 2009) have begun observing the radio sky, leading to an increase of data rates by up to 3 orders of magnitude (Wu et al., 2013; Davidson, 2012).

As the data acquisition rate has increased, data size has entered the Petabyte regime, and processing requirements increased to millions of CPU-hours. In order for processing to match the acquisition rate, the data is increasingly processed at large clusters with high-bandwidth connections to the data. An important case where data processing is done at a high throughput cluster is the LOFAR radio telescope.

The LOFAR telescope is a European low frequency aperture synthesis radio telescope centered in the Netherlands with stations stretching across Europe. This aperture synthesis telescope requires significant data processing before producing scientific images (Van Weeren et al., 2016; Williams et al., 2016; Smirnov and Tasse, 2015; Oonk et al., 2014). In this work, we will use our performance monitoring utility, *pipeline\_collector*<sup>2</sup>, to study the first half of the LOFAR processing, the Direction Independent (hereafter DI) pipeline.

One major project for the LOFAR telescope is the Sur-

---

Email address: [apmechev@strw.leidenuniv.nl](mailto:apmechev@strw.leidenuniv.nl) (Alexandar P. Mechev<sup>a</sup>)

<sup>1</sup><https://github.com/OpenTSDB/tcollector>

<sup>2</sup>[https://gitlab.com/apmechev/pipeline\\_collector.git](https://gitlab.com/apmechev/pipeline_collector.git)

veys Key Science Project (SKSP) (Shimwell et al., 2017). This project consists of more than 3000 observations of 8 hours each, 600 of which have been observed. These observations need to be processed by a DI pipeline, the results of which are calibrated by a Direction Dependent (DD) pipeline. The DI pipeline is implemented in the software package *prefactor*<sup>3</sup>. The *prefactor* pipeline is itself split into four stages and implemented at the SURFsara Grid location at the Amsterdam e-Science centre (SURF, 2018; Mechev et al., 2017). The automation and simple parallelization has decreased the run time per dataset from several days to six hours, making it comparable to the observation rate. To better understand and optimize the performance of the *prefactor* pipeline, we require detailed performance information for all steps of the processing software. We have developed a utility to gather this information for data processing pipelines running on distributed compute systems.

In this work, we will use the *pipeline\_collector* utility to study the LOFAR *prefactor* pipeline and suggest optimization based on our results. To test the software on a diverse set of hardware, we will set up the monitoring package on four different computers and collect data on the pipeline’s performance. Using this data, we discuss several aspects of the LOFAR software which we present in Table 1. Finally we discuss the broader context of these optimizations in relation to the LOFAR SKSP project and touch on the integration of *pipeline\_collector* with the second half of the data processing pipeline, the DD calibration and imaging.

### 1.1. Related Work

Scientific fields that need to process large data sets employ some type of data processing pipelines. Such pipelines include e.g. solar imaging (Centeno et al., 2014), neuroscience imaging (Strother et al., 2004) and infrared astronomy (Ott, 2010). While these pipelines often log the start and finishing times of each step (using tools such as pegasus-kickstart (Vöckler et al., 2006)), they do not collect detailed time series performance data throughout the run.

At a typical compute cluster the performance of every node in a distributed systems is monitored using utilities, such as Ganglia (Massie et al., 2004). These tools only monitor the global system performance. If one is interested in specific processes, then the Linux procfs (Bowden, 2009) is used. The procfs system can be used to analyse the performance of individual pipeline steps. Likewise, the Performance API (PAPI, Mucci et al., 1999) is a tool which collects detailed low level information on the CPU usage per process. Collecting detailed statistics at the process level is required to understand and optimize the performance of the LOFAR pipeline and we will integrate PAPI into *pipeline\_collector* in the future. Finally, DTrace (Gregg and Mauro, 2011) is a Sun Microsystems

tool which makes it possible to write profiling scripts that access data from the kernel and can be used to monitor process or system performance at run time with minimal overhead. As DTrace was not installed on either of the processing clusters, we have not used it to monitor the pipeline’s performance.

The Linux procfs system and PAPI record data which is already made available by the Linux kernel. This option incurs insignificant overhead as it uses data the kernel and processor already log. Likewise PAPI reads performance counters that the CPU automatically increments during processing. These profiling utilities can run concurrently with the scientific payload without using more than 1-2% of system resources. Their low overhead is why we choose to use them to collect performance data.

Other tools for performance analysis such as Valgrind (Nethercote and Seward, 2007) collect very detailed performance information. This comes at the expense of execution time: running with Valgrind, the processing time slows by up to two orders of magnitude. As such, we do not use Valgrind along the LOFAR software.

## 2. Measuring LOFAR Pipeline performance with pipeline\_collector

We developed the package *pipeline\_collector* as an extension of the performance collection package *tcollector*. *pipeline\_collector* makes it possible to collect performance data for complex multi-step pipelines. Additionally, it makes it easy to record performance data from other utilities. A performance monitoring utility that we plan to integrate in the future are the PAPI tools described in section 1.1. The resulting performance data was recorded in a database and analyzed. For our tests, we used the LOFAR *prefactor* pipeline, however with minor modifications, any multi-step pipeline can be profiled.

*tcollector* is a software package that automatically launches ‘collector’ scripts. These scripts are sample the specific system resource and send the data to the main *tcollector* process. This process then sends the data to the dedicated time series database. We created custom scripts to monitor processes launched by the *prefactor* pipeline (Appendix A.1).

In this work, we use a sample LOFAR SKSP data set as a test case. A particular focus was to understand the effect of hardware on the bottlenecks of the LOFAR data reduction. To gain insight into the effect of hardware on *prefactor* performance, the data was processed on four different hardware configurations (Table 2). As typical upgrade cycle for cluster hardware is five years, our results will be used to select optimal hardware for future clusters tasked with LOFAR processing.

### 2.1. Prefactor Pipeline

The LOFAR *prefactor* pipeline (Van Weeren et al., 2016) is a software pipeline that performs direction inde-

<sup>3</sup>available at <https://github.com/lofar-astron/prefactor>

Result #	Description
<b>R1</b>	Native compilation of the software performs comparably to pre-compiled binaries on two test machines.
<b>R2</b>	The processing steps do not appear to accelerate significantly on a faster processor or with larger cache size.
<b>R3</b>	Both calibration steps ( <i>calib_cal</i> and <i>gsmcal_solve</i> ) show linear correlation between speedup and memory bandwidth.
<b>R4</b>	Disk read/write speed does not affect the completion time of the slowest steps.
<b>R5</b>	Both calibration steps do not use large amounts of RAM despite processing data on the order of Gigabytes.
<b>R6</b>	The <i>calib_cal</i> step can suffer up to 20% of Level 1 Instruction Cache misses, while <i>gsmcal</i> only has 5% of these misses.
<b>R7</b>	Both calibration steps are impacted by Level 2 Cache eviction at comparable rates.
<b>R8</b>	The <i>calib_cal</i> step stalls on resources 70% of cycles while the <i>gsmcal</i> step only 30% of them.
<b>R9</b>	The <i>calib_cal</i> uses the CPU at full efficiency for only 10 % of the CPU cycles.

Table 1: A table of all the results presented in Section 3.

pendent calibration using the LOFAR software. The LOFAR software stack is a software package containing commonly used processing software used by LOFAR pipelines (Dijkema, 2017; Offringa et al., 2013). These tools are built and maintained by ASTRON<sup>4</sup>.

The *prefactor* pipeline performs a sequence of four stages, namely the calibrator and target calibration. The first half of *prefactor* processes data from a calibration source and the second half processes a science target. Altogether, this processing takes six hours on a high-throughput cluster. The final result is a data-set ready for creating images of the sky at radio wavelengths. Figure 1 shows a graphical view of the *prefactor* pipeline’s Calibrator and Target stages.

The Calibrator stage consists of the *ndppp\_prep\_cal* and the *calib\_cal* step. The former flags radio interference and averages the data, and the latter performs gain calibration on a bright calibration source. It is followed by the *fitclock* step which fits a clock-TEC model to the calibration solutions (Van Weeren et al., 2016).

The Target stage consists of a *ndppp\_prep\_targ* step, *predict\_ateam*, *gsmcal\_solve* and *gsmcal\_apply* steps. The first two of these steps flag and average the target data and calculate contamination by bright off-axis radio sources. The *gsmcal\_solve* step determines phase solutions for each antenna using a model of the target and the results of the *ndppp\_prep\_targ* step. Finally, the *gsmcal\_apply* step applies these solutions to the target data. Figure 2 shows the percentage of time spent by these steps for the four *prefactor* stages.

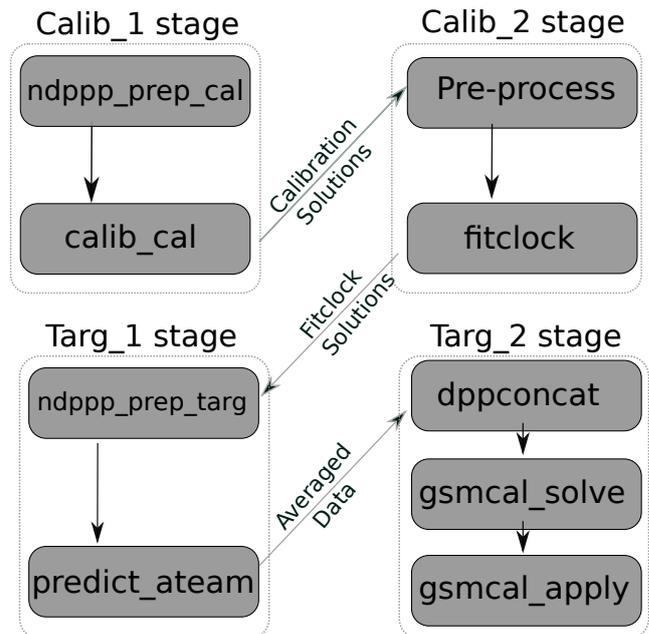


Figure 1: The four processing stages that make up the *prefactor* pipeline. The Calibrator stages (top) process a known bright calibrator to obtain the gain for the LOFAR antennas. The Target stages (bottom) process the scientific observation to remove Direction Independent effects. The *pref\_cal1* and *pref\_targ1* stages are massively parallelized across nodes without the need of an interconnect. The *pref\_cal2* step runs only on one node, while *pref\_targ2* is parallelized on 25 nodes. As the LOFAR software does not use MPI, we can run each processing job in isolation.

<sup>4</sup>ASTRON: Netherlands Institute for Radio Astronomy, url: <https://www.astron.nl/>

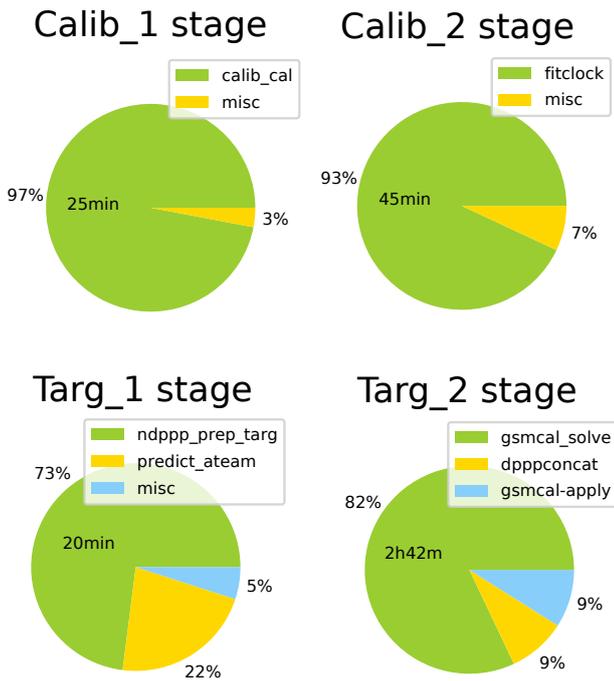


Figure 2: Portion of processing time taken by each step for the four prefactor stages, as reported by the Prefactor software. For each stage, the majority of processing time is spent during one or two steps. This is due to the fact that each prefactor stage also has intermediate book-keeping steps explicitly included in the pipeline. For each pipeline stage, the mean processing time for the longest running steps at SURFsara is also indicated. It should be noted that while faster, `pref.cal1` runs 10 times as many jobs as `pref.targ2`.

## 2.2. Performance suite

Cluster performance is frequently monitored using utilities such as Ganglia (Massie et al., 2004, discussed in Section 1.1). These tools cannot access individual processes and thus cannot collect data on a per-process basis. To collect such data, each process launched by the active pipeline step needs to be profiled individually. Our utility is designed to gather such performance data.

Our monitoring package, `pipeline_collector` adds custom performance collectors (Appendix A.1) to the performance collection framework `tcollector`. We use these collectors to monitor individual pipeline steps as defined by the user<sup>5</sup>. The tools attach to processes launched by the pipeline and record performance data at a one second interval. This sampling frequency is at high enough resolution to detect trends and anomalies in hardware utilization, and still result in a reasonable database size. The performance data is uploaded to a remote time series database, OpenTSDB (Sigoure, 2012). Details on the data collection can be found in Appendix A.

### 2.2.1. Performance API

The time-series database is also used to collect low-level CPU information for each process. This information is collected by the PAPI interface (discussed in Section 1.1). This was done through the `papiex` utility<sup>6</sup> (Ahn, 2008). This utility records the CPU’s internal performance counters. A CPU’s performance counters record information on how efficiently the software uses the CPU’s resources. The results from this test are detailed in Section 4.

## 2.3. Test Hardware

In order to study the effect of different hardware configurations on the performance of LOFAR processing, the `prefactor` pipeline was run on four different sets of hardware. The four machines tasked with processing LOFAR data comprised nodes at three computational clusters and a personal computer. The tests were run while the systems were idle to make sure there is no interference of other software with the LOFAR processing. Table 2 details the specifics of the four test machines.

## 3. LOFAR Prefactor Test Case

With the test set described in Section 2, we aim to understand processing bottlenecks in the `prefactor` pipeline and make informed decisions on future hardware and software upgrades. To do so, we processed a sample observation at institutes that typically process LOFAR data.

From the data collected by processing the sample observation, we determined the slowest pipeline steps. These steps were the `calib_cal` and `gsmcal_solve`, seen in Figure 2. The `calib_cal` step is implemented by the software

<sup>5</sup>[https://gitlab.com/apmechev/pipeline\\_collector.git](https://gitlab.com/apmechev/pipeline_collector.git)

<sup>6</sup>Available at <https://bitbucket.org/minimalmetrics/papiex-oss>

Location	CPU Speed (MHz)	CPU Model	Micro-architecture	Cache Size	RAM Speed <sup>7</sup>	Disk Speed <sup>8</sup>
Leiden	2200	E5-4620	Sandy Bridge	16 MB	1.4 GB/s	99.7 MB/s
SURFsara	2500	E5-2680	Sandy Bridge	30 MB	2.5 GB/s	65.4 MB/s
Hatfield	2900	E5-2660	Sandy Bridge	20 MB	2.4 GB/s	155 MB/s
Laptop	3300	E3-1505M	Skylake	8 MB	4.7 GB/s	822 MB/s

Table 2: CPU, Cache, RAM and Storage specifications of the four test machines. The tested machines span a factor of 1.5x in CPU speed, 4x in cache and RAM Speed and 10x in Disk speed.

`bbs-reducer` (Dijkema, 2017; Loose, 2008) and the `gsmcal.solve` step is implemented by NDPPP (Dijkema, 2017; Offringa et al., 2013). Both `bbs-reducer` and NDPPP are part of the LOFAR software suite.

We collected performance statistics using the `pipeline_collector` suite as discussed in Section 2. The runtime of the slowest `prefactor` steps on the four machines is shown in figure 3. The results discovered using `pipeline_collector` are listed in Table 1 and discussed in Section 3.2. Using the PAPI interface (discussed in Section 2.2.1) CPU performance data was collected. The results from this test are detailed in Section 4.

We will present a number of insights into the performance of the LOFAR software collected by the profiling suite. The results are presented in Table 1 and are grouped in three main areas. The effect of compilation on the runtime was result **R1**. The set of results **R2**, **R3**, **R4** and **R5** were obtained using the `pipeline_collector` package. Results **R6**, **R7**, **R8** and **R9** were collected with the PAPI package, which will be integrated into `pipeline_collector` in the future.

### 3.1. Pre-compiled vs native compilation

The performance trade-off between pre-compiled and native compilation was studied first. The majority of the processing for the LOFAR SKSP Project (Shimwell et al., 2017) is done at the SURFsara GINA cluster in Amsterdam. This location is part of the European Grid Initiative (EGI)(SURF, 2018). At this location, software is deployed by compiling on a virtual machine and mounting it on all worker nodes through the CernVM File System (CVMFS) service (Blomer et al., 2011). The CVMFS server allows any client to mount a fully compiled LOFAR installation, making it easy to distribute and version control the software within and outside of SURFsara. An alternative is to locally compile the LOFAR packages on each cluster. The performance of the natively compiled<sup>9</sup> vs CVMFS in-

stallations was compared on the laptop test machine using `pipeline_collector`. In order to validate this result, the two compilations of the same software were also tested at the Data Science Lab at the Leiden Institute of Advanced Computer Science (LIACS)<sup>10</sup>.

An interesting discovery is that the LOFAR software did not process data faster when compiled natively. This is despite the fact that the local install was compiled with advanced processor instructions available on the host machine. Figure 4 shows a histogram of its processing time with the two different compilation options for the `calibcal` software running on the sample dataset. The same test was done for the software performing the gain calibration (`gsmcal.solve`), seen in Figure 5. The result of this experiment is shown in Figures 4a and 5a. The software compiled at SURFsara showed a minor improvement for the `calibcal` step on the laptop machine, however this improvement is not seen on the computational cluster node.

Overall, the software for both steps show no significant improvement when compiled natively. This is result **R1** in Table 1. The second run at LIACS also confirms this result for both steps (Figures 4b and 5b). This result suggests that the slowest `prefactor` steps are not optimized for modern processors.

### 3.2. Prefactor Runtime and Hardware Parameters

Next, we studied the dependence of runtime on different hardware parameters. With software that collects per-step performance statistics for the LOFAR pipeline, the dependence of the pipeline processing on hardware performance can be easily profiled and studied. Using `pipeline_collector` we determined the pipeline’s slowest steps with respect to different hardware parameters.

The system parameters studied here are the CPU speed, memory throughput, cache size and disk speed. Modern computers can have a complex memory hierarchy as demonstrated in Figure 6 (Katz and Patterson, 2001). This is due to the cost trade-off between memory size and memory speed. Because of this trade-off, the full dataset is stored on disk, while the working set is placed in RAM. This is the data that the processor needs to access at the current time (Denning, 1968). The most frequently accessed parts of the data are stored in the CPU cache,

<sup>7</sup>benchmarked using `dd`

<sup>8</sup>sequential disk read, benchmarked using `fio` - flexible I/O tester:

```
fio --randrepeat=1 --ioengine=libaio --direct=1
--gtod_reduce=1 --name=test --filename=test --bs=4k
--iodepth=256 --size=4G --readwrite=read --ramp_time=4
```

<sup>9</sup>The software was compiled using `-march=native` and `-O3` compilation flags. On the laptop, gcc resolves `-march=native` as `broadwell`. The CVMFS installation resolves `-march=native` as `core-avx-i`.

<sup>10</sup><https://www.universiteitleiden.nl/en/science/computer-science/about-us/our-facilities>

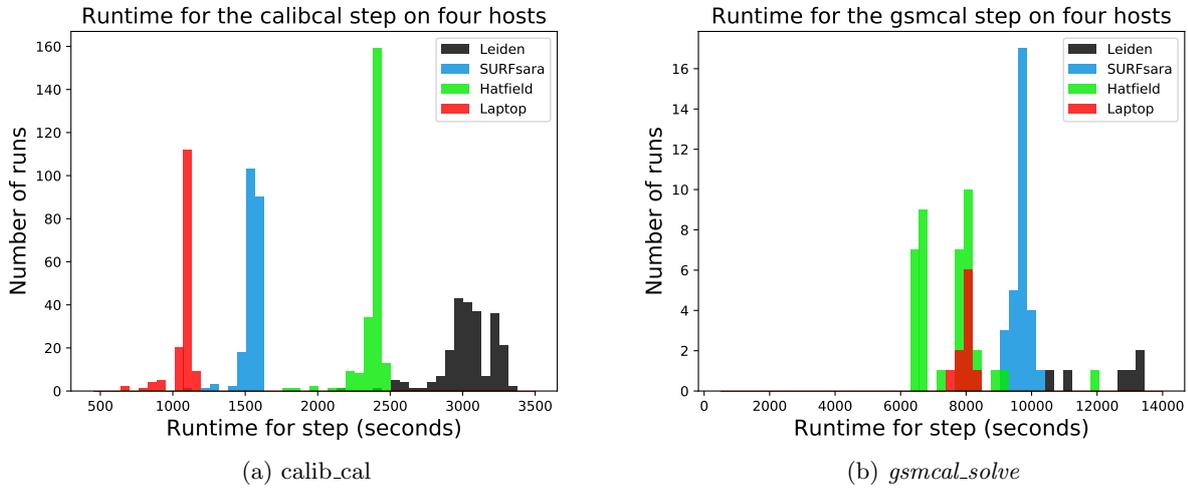


Figure 3: Job completion times for *calib\_cal* and *gsmcal\_solve* steps tested on four hardware setups. The *calib\_cal* step ran 244 times. The *gsmcal\_solve* ran 24 times as the data is concatenated from 244x1 to 24x10 sets. The step with the longest latency is *gsmcal\_solve* while *calib\_cal* consumes a comparable number of core-hours over 244 jobs.

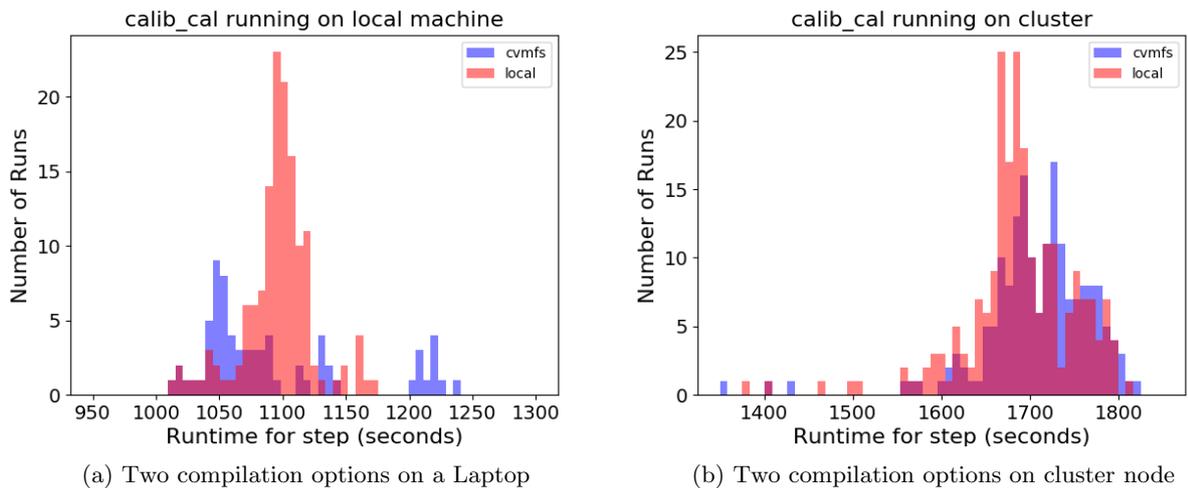


Figure 4: Difference in processing time for *calib\_cal* when compiled remotely and natively. *calib\_cal* was run 244 times with the native software and 40 times with the CVMFS compilation. Two tests were done, one on the personal laptop (4a) and one on a cluster node at the LIACS Data Science Lab (4b). The test on a cluster node shows no significant difference in runtime between compilation options. The laptop test suggests that the remotely compiled software may run 5% faster than the local compilation.

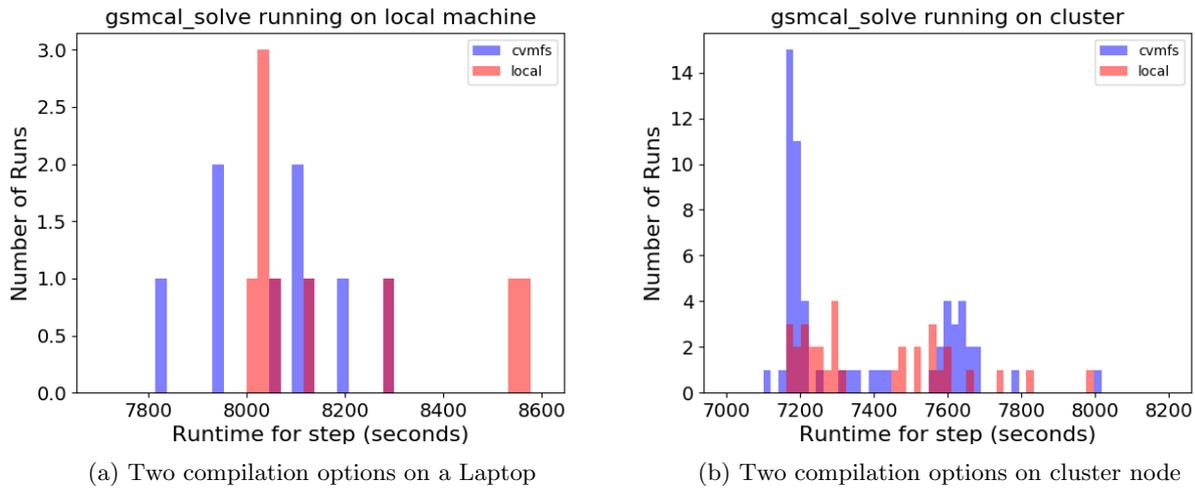


Figure 5: Difference in processing time for *gsmcal\_solve* when compiled remotely and natively. *gsmcal\_solve* was run 50 times with the native software and 120 times with the CVMFS compilation. Two tests were done, one on the personal laptop (5a) and one on a cluster node at the LIACS Data Science Lab (5b). Just like with the *calib\_cal* step, the *gsmcal\_solve* step also doesn't accelerate significantly when natively compiled.

which evicts the oldest data when full (Hazelwood and Smith, 2004).

The CPU processing speed is faster than the RAM latency, so a hierarchy of caches exist. Caches store small subsets of the working set and have a fast connection to the processor. The fastest data link is between the CPU and the L1 Cache, with the link to RAM being slower and the disk read speed slower still. The limited memory capacity of the different levels of the memory hierarchy as well as the throughput between them will lead to performance bottlenecks. These bottlenecks will lead to the processor waiting on memory. Such stalls lead to longer processing times.

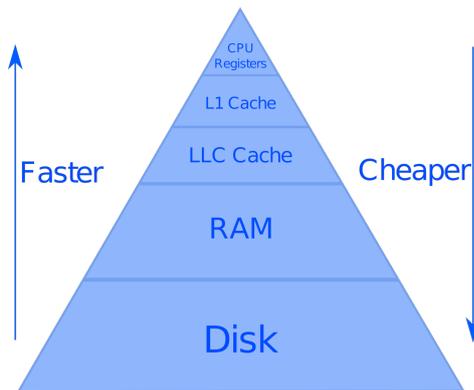


Figure 6: A model of the memory hierarchy, as described in (Goto and Geijn, 2008).

### 3.2.1. CPU

The CPU speed is usually the primary factor determining how fast computations can be made. In general, a faster CPU will result in faster data processing.

However, Fig. 7a shows that the runtime of the calibration of the calibrator does not strongly depend on the CPU frequency. While the test nodes at SURFsara and Leiden run at the same CPU frequency, running on a cluster node at SURFsara takes half the time as on a node at Leiden. Even more surprisingly, the *gsmcal\_solve* step does not benefit significantly from a faster CPU, despite being the most computationally heavy *prefactor* step (**R2**). This step does the gain calibration on the target field using the StEFCal algorithm (Salvini and Wijnholds, 2014). Figure 7b shows only a slight improvement over faster CPU clock speeds for both steps. The correlation between completion time and CPU speed is similar for both steps.

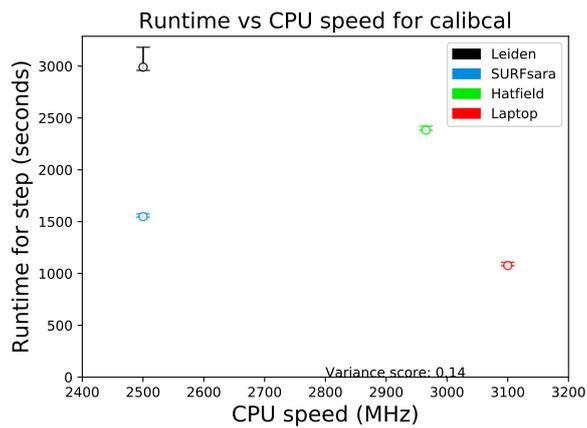
### 3.2.2. Cache

The CPU has a hierarchy of caches consisting of Level 1, Level 2 Cache and LLC Cache. For the four processors tested, the Level 1 and 2 caches were all the same size, thus the only difference is the Last Level Cache (LLC or just Cache in Figure 6). This cache stores data needed by the CPU, so the larger it is, the less the processor needs to wait for RAM to return data.

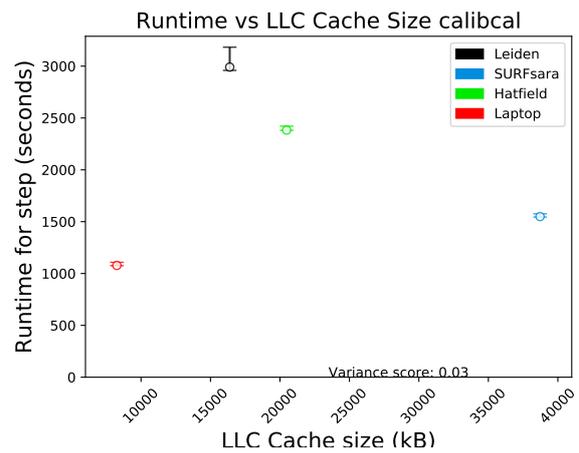
In general, numerical codes benefit from larger cache sizes (Skadron et al., 1999; Goto and Geijn, 2008). Interestingly, figure 8b suggests that the *gsmcal\_solve* step does not exclusively depend on larger cache **R3** (Table 1). On the machines with a larger cache, the *gsmcal\_solve* step completed processing as quickly as on the machines with smaller cache, even down to 8MB.

### 3.2.3. RAM Bandwidth

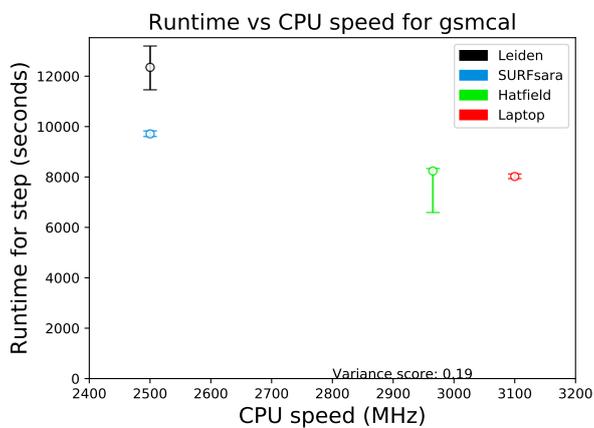
If the entire data set does not fit into cache, the software needs to transfer data from RAM to the CPU. In these cases, *prefactor* benefits from a fast bandwidth between the cache and RAM. For this study, the RAM



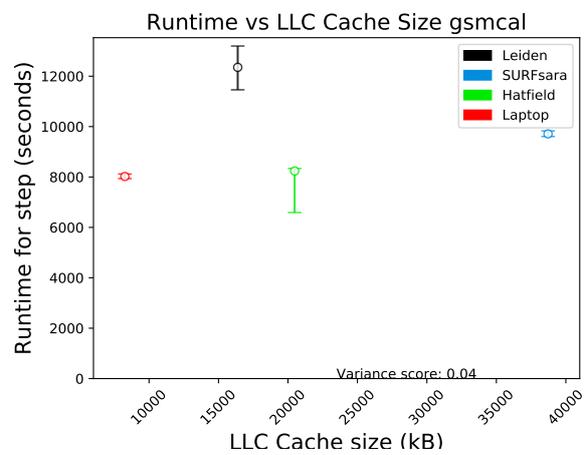
(a) calib\_cal



(a) calib\_cal



(b) gsmcal\_solve



(b) gsmcal\_solve

Figure 7: Performance of the bottleneck steps compared with the CPU speeds of the four test machines. The values are the mean of 244 runs (Standard prefactor run) and the error bars show the 1-sigma of the distribution of the run time.

Figure 8: Performance of the two bottleneck steps with respect to Last Level Cache size. The *gsmcal\_solve* step shows no trend between cache size and completion time. The *calib\_cal* step runs the fastest on the machine with the smallest cache.

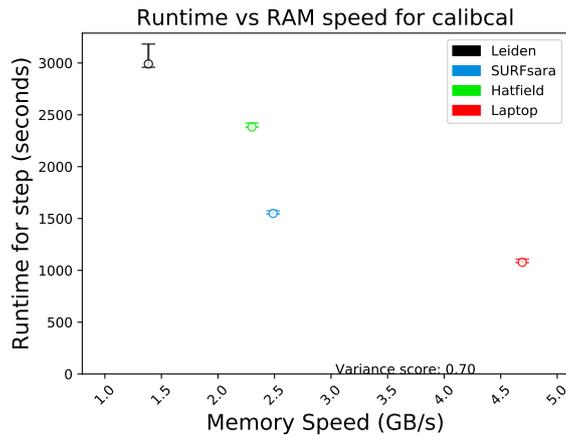
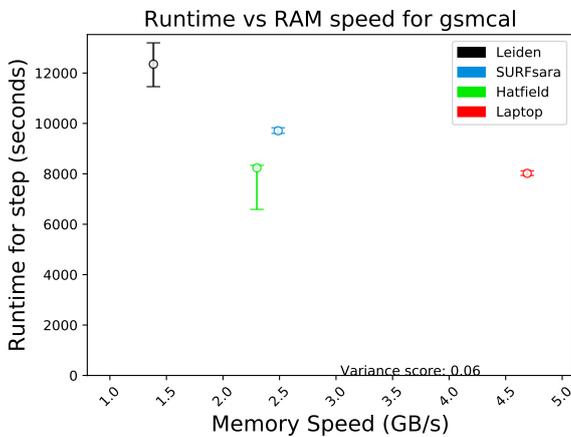
(a) *calib\_cal*(b) *gsmcal\_solve*

Figure 9: Performance of the two bottleneck steps and RAM bandwidth in GB/s. Both the *calib\_cal* and *gsmcal\_solve* steps show a trend of faster processing times on machines with higher RAM bandwidth. Both steps show a trend of decreasing processing time with increasing RAM throughput.

throughput was benchmarked<sup>11</sup>. This command copies dummy data into system memory. As this utility exists on all Unix systems, this is a standardized benchmark of the RAM performance.

Figure 9a showed that higher bandwidth is correlated with a faster completion time for the *calib\_cal* and *gsmcal\_solve* steps (**R4**). The result is to be expected as the working set of these steps is 200MB and 1.0GB respectively, and cannot fit into cache readily, however it is loaded into RAM within the first 5 seconds of the run (Figure 10), and is streamed from memory throughout the run.

<sup>11</sup>Using the command `$> dd if=/dev/zero of=/dev/shm/test bs=1M count=2048`

### 3.2.4. Disk Read speeds

The slowest link in the memory hierarchy is the disk read speed. For the *calib\_cal* step, the entire data is loaded into memory during the first few seconds of the run, after which the disk only becomes important when the results need to be written out. The *gsmcal\_solve* step streams data from the disk to memory throughout the entire run. The plot of disk read speeds (Fig. 11b) also shows that a faster disk does not speed up the slowest step **R5**. To verify that disk throughput was not the limiting factor, the entire dataset (25 GB) was moved to main memory (using `/dev/shm`). The resulting runtime for both bottleneck steps did not change.

The calibration steps both stored less than 200MB of data in memory throughout their run. Figure 10 shows the time-series of the total memory used by these steps. The *calib\_cal* step uses only 200MB of memory and *gsmcal\_solve* only 35MB. While the *gsmcal\_solve* step works on a 1GB dataset, it streams the data in memory and thus does not require 1GB of RAM. Alternatively, the *calib\_cal* step loads the entire (200MB) dataset into memory for the entire duration of the run. The RAM usage time-series in Figure 10 show that the RAM is filled for the first 5 seconds of the run, further confirming that the processing is effectively independent from disk speed.

## 4. CPU Utilization Tests with PAPI

To gain more fine grained data on the CPU utilization, the *calib\_cal* and *gsmcal\_solve* steps were tested with the PAPI package. We ran this package as a test, to determine whether collecting PAPI data is helpful in understanding pipeline performance. PAPI can record data such as cache performance, branch prediction rate, fraction of memory/branch instructions and others. This data is complementary to the `procs` information, which is collected by the Linux kernel. As the collected data was useful in understanding the *prefactor* pipeline, we will include PAPI in the *pipeline\_collector* suite in the future. In the following sections we will discuss the results obtained for the *calib\_cal* and *gsmcal\_solve* steps.

### 4.1. Level 1 Data Misses

The Level 1 Cache is split into cache for instructions and data. For all our test hardware the L1 Data cache is 32 kB, and has a direct link to the processor’s computational units (Jain and Agrawal, 2013). The processor collects information logging how many times data requested by the CPU is not located into the L1 Data cache. This counter is called the Level 1 Data Cache Miss rate. To resolve this type of cache miss, the data needs to be fetched from L2 Cache. When this happens, the processor has to wait for the requested data. **R7**: The recorded L1 data misses in Figure 12a, show that the software performing the *calib\_cal* step misses 20% of its L1 data cache requests, while the software implementing the *gsmcal\_solve* step misses less

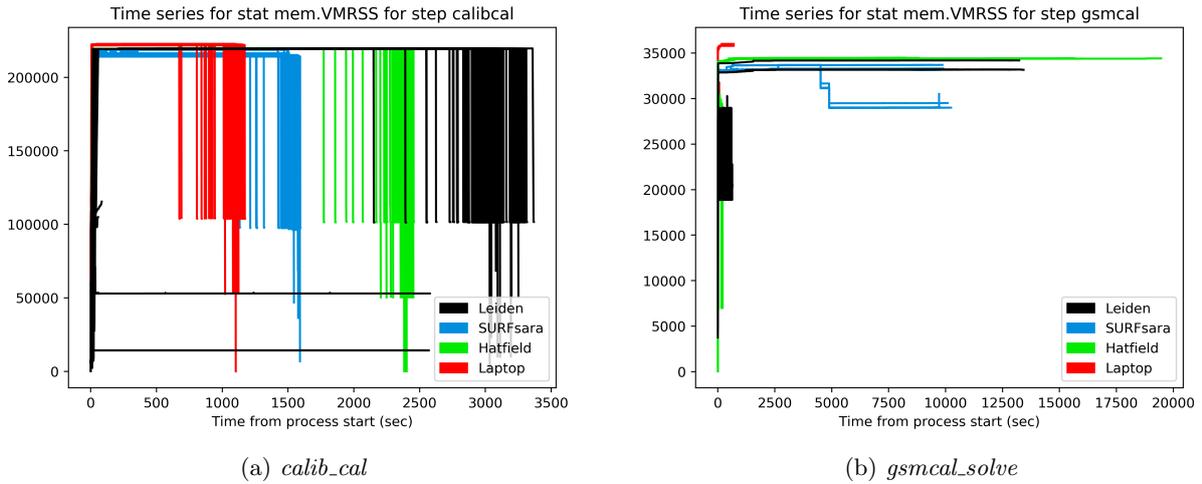


Figure 10: Time series of the Virtual Memory Resident Set Size . This is the amount of data stored in RAM (in kB) during the *calib\_cal* and *gsmcal\_solve* steps. Both steps show the same amount of memory use on all test machines. Additionally, after a brief loading of data, the memory usage remains constant until processing is finished.

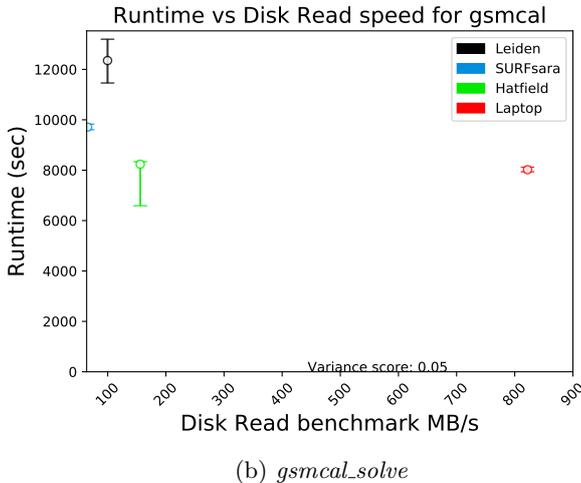
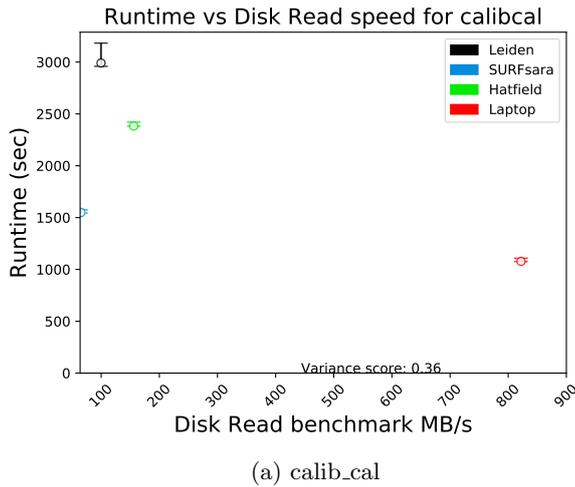


Figure 11: Performance of the two bottleneck steps and Disk bandwidth in MB/s. There is no correlation between the Disk read speed and the Runtime of the steps.

than 5% of L1 Cache requests. These cache misses often happens in multi-threaded applications where there are instructions shared by multiple threads on the same cache line (Sarkar and Tullsen, 2008).

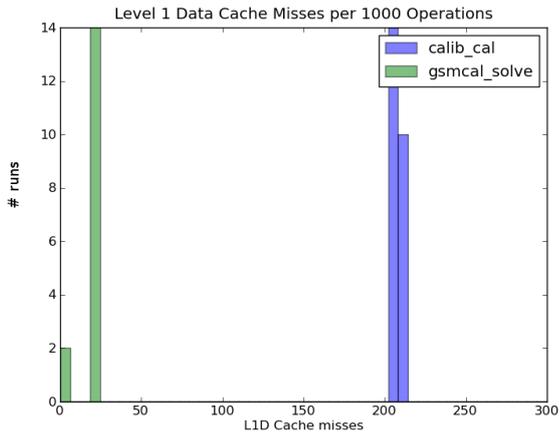
#### 4.2. Level 2 Instruction Misses

Unlike the Level 1 cache, Level 2 cache stores data and instructions in the same location. When the cache is full, it evicts the last used element in order to make space for newly requested data. PAPI also counts these eviction events. Figure 12b shows that for both steps, between 50 and 70% of L2 requests for an instruction do not match the contents of L2 Cache. This is significantly more than the applications benchmarked in (Lebeck et al., 2002, Table 2). Because both steps process data of considerable size, the large amount of data required can evict instructions from the L2 cache (insight number **R7** in table 1).

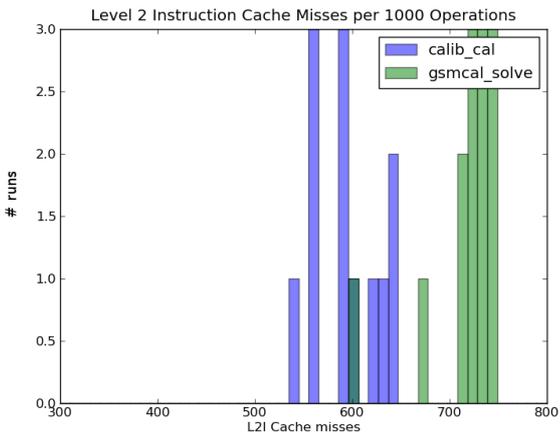
#### 4.3. Resource Stalls

Modern processors have multiple computational pipelines on chip, in order to process data in parallel (Hu et al., 2006). There are times when the processor's internal pipeline needs to wait for other instructions to finish. When this happens, it flags that it has 'stalled on a resource'. These resource stall cycles are also recorded by PAPI and represented as a percentage of total cycles. From figure 13a, it can be seen that *calib\_cal* stalls on 70% of the processor cycles, while *gsmcal\_solve* only on 33% of cycles (**R8**).

The Full Issue Cycles counter indicates the percentage of processor cycles, in which the theoretical maximum number of instructions are executed. During these cycles, the software uses the CPU optimally. The full issue cycles counter (Fig. 13b) also shows the difference in efficiency between the *calib\_cal* and *gsmcal\_solve* step (**R9**), with

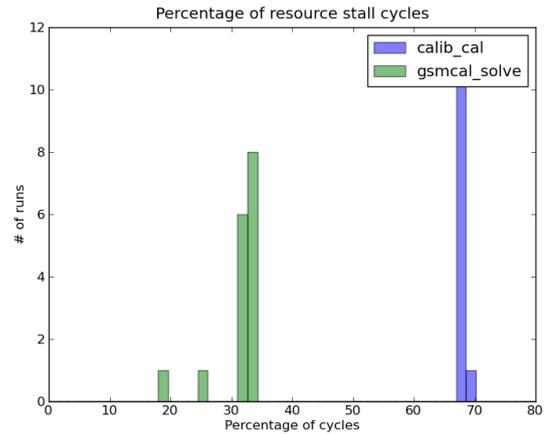


(a) Level 1 Data cache misses

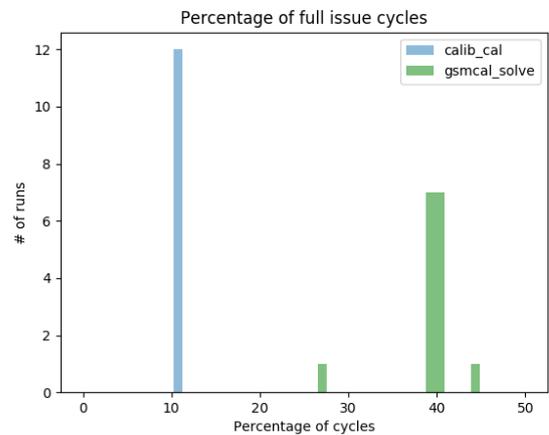


(b) Level 2 Instruction cache misses

Figure 12: Cache miss rates for *calib\_cal* and *gsmcal\_solve*, executed on the SURFsara gina cluster. The cache is split into instruction and data caches. The figures above show the difference in number of cache misses for both instruction cache and data cache for the slowest *prefactor* steps. *calib\_cal* suffers significantly more Data cache misses than *gsmcal\_solve* while the two steps undergo similar instruction Cache misses.



(a) Resource Stall Cycles



(b) Full Issue cycles

Figure 13: Resource stall cycles and Full Instruction Issue cycles. The two steps were executed on the SURFsara gina cluster.

the former only working at peak efficiency for 10% of the processor cycles.

The plots in Figures 13a and 13b indicate that the *calib\_cal* step does not use the internal CPU pipelines efficiently leading to waiting on resources and sub-optimal use of the CPU's Computational Units.

## 5. Discussions and Recommendations

With an increase of data acquisition rates and data complexity in radio astronomy, it is becoming important to thoroughly understand and optimize the performance of processing pipelines. Using *pipeline\_collector*, data can be collected for each pipeline step without altering the processing software. We store this data in a time-series database. The collected data can be studied to help researchers understand the pipeline performance for different processing parameters, datasets, and on different hardware. The *pipeline\_collector* suite is easy to deploy for mature pipelines and has minimal impact on pipeline per-

formance. Typical CPU usage is  $<0.2\%$  with a memory footprint of  $\sim 1\text{-}10$  MB.

Creating a performance model with the collected data will allow us to optimize future clusters for LOFAR data processing. Doing so is necessary given the current data throughput, number of observations and time-line of the SKSP project. Similar issues will be encountered with upcoming radio telescopes (Broekema et al., 2015).

To showcase the power of the *pipeline\_collector* suite, the LOFAR *prefactor* pipeline was run through a single data set on three clusters and a personal machine. A number of insights were made using the high resolution timing data collected from this package (such as in Figure 10) and are listed in Table 1. In the future, we'll apply the *pipeline\_collector* software to the more complex LOFAR DD pipeline, *ddf-pipeline*<sup>12</sup>.

The slowest processing steps for the *prefactor* pipeline were identified as the *calib\_cal* and *gsmcal\_solve* steps. While the data can fit into the RAM for all of the processing machines, it is much larger than the processor's internal cache (Figure 6). The discoveries made concerned the memory hierarchy in Figure 6. Results labeled **R2**, **R8** and **R9** related to the CPU performance; **R2**, **R6** and **R7** related to the Cache performance; **R3** and **R5** related to the Memory usage and **R4** discussed the Disk speed.

Faster processors did not accelerate the *gsmcal\_solve* step significantly, as this step streams data between the RAM and CPU. As the CPU speed increases, streaming applications become bottlenecked by the throughput of data into the CPU from RAM. As the *gsmcal\_solve* algorithm iteratively calibrates chunks of the data, these chunks need to be loaded from disk once, however they are moved from RAM to CPU multiple times during calibration.

Similarly, the *calib\_cal* step is more dependent on memory throughput than on CPU speed as this step moves data to and from memory frequently. This step also does minimization looping over the dataset. As the dataset does not fit in the cache, parts of it need to be constantly moving from memory and back. Figure 8a shows that the machine with the smallest LLC cache runs the *calib\_cal* step the fastest. This is likely a combination of the benefit of faster RAM and poor cache optimization for this software. The same effect is much less pronounced in Figure 8b, suggesting that software optimization at least plays a part in the outliers for the laptop machine.

### 5.1. Recommendations

Based on these results, the top hardware recommendation is that *prefactor*'s slowest steps can be accelerated by running on machines with faster memory or upgrading the memory of the current machines. The two slowest *prefactor* steps showed improvements on machines with faster RAM.

One software recommendation is to improve the efficiency of the *calib\_cal* step through refactoring or by replacing the software package used. Unfortunately, the software used for the *gsmcal\_solve* step cannot be used for the *calib\_cal* step as it is not yet able to correct for Faraday Rotation (Salvini and Wijnholds, 2014), making it impossible to currently use the software used by the *gsmcal\_solve* step. Faraday Rotation has recently been implemented in a development version of the *prefactor* pipeline and is currently undergoing testing. This version of the pipeline will be implemented by September 2018.

Additionally, the large number of data cache misses recorded for the *calib\_cal* step suggests that its source code is not optimized for multi-threaded processing. Data cache misses are often encountered when multiple threads have instructions on the same cache line<sup>13</sup>, forcing the memory controller to move this cache line between cores (Lebeck et al., 2002). This can also explain the large number of stalled cycles (Fig. 13a) and low number of full issue cycles (Fig. 13b) for the *calib\_cal* step. It is recommended to further study the inefficiencies of *calib\_cal* or to replace it with a newer software. If the software processing for this step is updated, analysing the cache and CPU performance of the new software will be necessary to determine whether it efficiently uses the available computational resources.

Finally, we discovered that compiling the software on a virtual machine did not lead to a processing slowdown. This means that the current slowest *prefactor* steps are not optimized to use advanced processor instructions. Nevertheless, the resulting cross-compatibility is an encouraging result as it will allow to easily distribute pre-compiled versions of the software without increasing the processing time. We recommend continuing CVMFS deployment of LOFAR software.

## 6. Conclusions

In this paper, we present a novel system for automated collection of performance data for complex software pipelines. We use this suite to study the LOFAR *prefactor* pipeline. The results are discussed aiming to understand the effect of different hardware parameters on the data processing. To do so, we run the pipeline on four different machines.

The software automatically collects performance data at the operating system level without impacting processing time. Data for each pipeline step is extracted using the OpenTSDB API, plotted and analyzed. Additionally, the *pipeline\_collector* suite is easy to extend with new collectors that record more detailed time-series data for each pipeline step. The performance data is stored in the time series database OpenTSDB.

<sup>13</sup>A cache line is a row of cache memory which is loaded into CPU as a single unit (David and John, 2005)

<sup>12</sup><https://github.com/mhardcastle/ddf-pipeline>



file holds the sample interval, executables to monitor and the location where `pipeline_collector` can read the current pipeline step

The `pipeline_collector` suite reads the current pipeline step from a file, the location of which is specified in the configuration. This file needs to be updated each time the pipeline begins a new step. For LOFAR we have a script running with the pipeline, and determining the current step using the pipeline logs. As each pipeline has a unique sequence of steps, the current step needs to be recorded in a file in order for `pipeline_collector` to report it to the time series database. The location of the file recording the current pipeline step is read from the configuration file.

Next, the names of the specific processes need to be included in the configuration file. In the case of LOFAR, we select the `NDPPP`, `bbs-reducer` and `losoto` processes. The `pipeline_collector` searches the running processes for the current user for these process names and launches a collector for each new process launched by the current step.

## Acknowledgements

APM would like to acknowledge the support from the NWO/DOME/IBM programme “Big Bang Big Data: Innovating ICT as a Driver For Astronomy”, project #628.002.001.

HJR gratefully acknowledge support from the European Research Council under the European Unions Seventh Framework Programme (FP/2007-2013)/ERC Advanced Grant NEWCLUSTERS- 321271.

JBRO acknowledges financial support from NWO Top LOFAR-CRRL project, project No. 614.001.351.

This work was carried out on the Dutch national e-infrastructure with the support of SURF Cooperative through grant e-infra 160022 & 160152.

## References

- D. H. Ahn. Measuring flops using hardware performance counter technologies on lc systems. Technical report, Lawrence Livermore National Laboratory (LLNL), Livermore, CA, 2008.
- S. F. Apache. Tcollector: OpenTSDB documentation. Available at [http://opentsdb.net/docs/build/html/user\\_guide/utilities/tcollector.html](http://opentsdb.net/docs/build/html/user_guide/utilities/tcollector.html), 2017.
- T. Apache HBase. Apache hbase reference guide. *Apache*, version, 2(0), 2015.
- J. Blomer, C. Aguado-Sánchez, P. Buncic, and A. Harutyunyan. Distributing LHC application software and conditions databases using the cervm file system. In *Journal of Physics: Conference Series*, volume 331, page 042003. IOP Publishing, 2011.
- T. Bowden. THE /proc FILESYSTEM v1.3. <https://www.kernel.org/doc/Documentation/filesystems/proc.txt>, 2009.
- P. C. Broekema, R. V. van Nieuwpoort, and H. E. Bal. The Square Kilometre Array science data processor. Preliminary compute platform design. *Journal of Instrumentation*, 10(07):C07004, 2015.
- P. C. Broekema, J. J. D. Mol, R. Nijboer, A. van Amelsfoort, M. Brentjens, G. M. Loose, W. Klijin, and J. Romein. Cobalt: A gpu-based correlator and beamformer for lofar. *Astronomy and Computing*, 23:180 – 192, 2018. ISSN 2213-1337. doi: <https://doi.org/10.1016/j.ascom.2018.04.006>. URL <http://www.sciencedirect.com/science/article/pii/S2213133717301439>.
- R. Centeno, J. Schou, K. Hayashi, A. Norton, J. Hoeksema, Y. Liu, K. Leka, and G. Barnes. The Helioseismic and Magnetic Imager (HMI) vector magnetic field pipeline: optimization of the spectral line inversion code. *Solar Physics*, 289(9):3531–3547, 2014.
- A. P. David and L. H. John. Computer organization and design: the hardware/software interface. *San mateo, CA: Morgan Kaufmann Publishers*, 1:998, 2005.
- D. B. Davidson. Potential technological spin-offs from MeerKAT and the South African Square Kilometre Array bid. *South African Journal of Science*, 108(1-2):01–03, 2012.
- P. J. Denning. The working set model for program behavior. *Communications of the ACM*, 11(5):323–333, 1968.
- T. J. Dijkema. LOFAR Imaging Cookbook. Available at [http://www.astron.nl/sites/astron.nl/files/cms/lofar\\_imaging\\_cookbook\\_v19.pdf](http://www.astron.nl/sites/astron.nl/files/cms/lofar_imaging_cookbook_v19.pdf), 2017.
- K. Goto and R. A. Geijn. Anatomy of high-performance matrix multiplication. *ACM Transactions on Mathematical Software (TOMS)*, 34(3):12, 2008.
- B. Gregg and J. Mauro. *DTrace: Dynamic Tracing in Oracle Solaris, Mac OS X and FreeBSD*. Prentice Hall Professional, 2011.
- Y. Gupta, B. Ajithkumar, H. Kale, S. Nayak, S. Sabhapathy, S. Sureshkumar, R. Swami, J. Chengalur, S. Ghosh, C. Ishwara-Chandra, et al. The upgraded GMRT: opening new windows on the radio Universe. *Current Science*, 113(4):707, 2017.
- K. Hazelwood and J. E. Smith. Exploring code cache eviction granularities in dynamic optimization systems. In *Code Generation and Optimization, 2004. CGO 2004. International Symposium on*, pages 89–99. IEEE, 2004.
- S. Hu, I. Kim, M. H. Lipasti, and J. E. Smith. An approach for implementing efficient superscalar CISC processors. In *High-Performance Computer Architecture, 2006. The Twelfth International Symposium on*, pages 41–52. IEEE, 2006.
- T. Jain and T. Agrawal. The haswell microarchitecture-4th generation processor. *International Journal of Computer Science and Information Technologies*, 4(3):477–480, 2013.
- J. L. Jonas. MeerKAT-The South African array with composite dishes and wide-band single pixel feeds. *Proceedings of the IEEE*, 97(8):1522–1530, 2009.
- R. H. Katz and D. A. Patterson. Memory hierarchy, CM-PUT429/CMPE382 Winter 2001. University of Calgary. Available at <https://webdocs.cs.ualberta.ca/~amaral/courses/429/webslides/Topic4-MemoryHierarchy/sld003.htm>, 2001.
- A. R. Lebeck, J. Koppanalil, T. Li, J. Patwardhan, and E. Rotenberg. A large, fast instruction window for tolerating cache misses. In *Computer Architecture, 2002. Proceedings. 29th Annual International Symposium on*, pages 59–70. IEEE, 2002.
- C. J. Lonsdale, R. J. Cappallo, M. F. Morales, F. H. Briggs, L. Benkevitch, J. D. Bowman, J. D. Bunton, S. Burns, B. E. Corey, S. S. Doleman, et al. The murchison widefield array: Design overview. *Proceedings of the IEEE*, 97(8):1497–1506, 2009.
- G. Loose. LOFAR self-calibration using a blackboard software architecture. In *Astronomical Data Analysis Software and Systems XVII*, volume 394, page 91, 2008.
- M. L. Massie, B. N. Chun, and D. E. Culler. The ganglia distributed monitoring system: design, implementation, and experience. *Parallel Computing*, 30(7):817–840, 2004.
- A. Mechev, J. B. R. Oonk, A. Danezi, T. W. Shimwell, C. Schrijvers, H. Intema, A. Plaat, and H. J. A. Rottgering. An Automated Scalable Framework for Distributing Radio Astronomy Processing Across Clusters and Clouds. In *Proceedings of the International Symposium on Grids and Clouds (ISGC) 2017, held 5-10 March, 2017 at Academia Sinica, Taipei, Taiwan (ISGC2017)*. Online at <https://pos.sissa.it/cgi-bin/reader/conf.cgi?confid=293>, id.2, page 2, Mar. 2017.
- P. J. Mucci, S. Browne, C. Deane, and G. Ho. PAPI: A portable interface to hardware performance counters. In *Proceedings of the department of defense HPCMP users group conference*, volume 710, 1999.
- N. Nethercote and J. Seward. Valgrind: a framework for heavy-weight dynamic binary instrumentation. In *ACM Sigplan notices*, volume 42, pages 89–100. ACM, 2007.

- A. Offringa, A. De Bruyn, S. Zaroubi, G. van Diepen, O. Martinez-Ruby, P. Labropoulos, M. A. Brentjens, B. Ciardi, S. Daiboo, G. Harker, et al. The LOFAR radio environment. *Astronomy & astrophysics*, 549:A11, 2013.
- J. B. R. Oonk, R. J. van Weeren, F. Salgado, L. K. Morabito, A. G. G. M. Tielens, H. J. A. Rottgering, A. Asgekar, G. J. White, A. Alexov, J. Anderson, I. M. Avruch, F. Batejat, R. Beck, M. E. Bell, I. van Bemmel, M. J. Bentum, G. Bernardi, P. Best, A. Bonafede, F. Breitling, M. Brentjens, J. Broderick, M. Brügggen, H. R. Butcher, B. Ciardi, J. E. Conway, A. Corstanje, F. de Gasperin, E. de Geus, M. de Vos, S. Duscha, J. Eislöffel, D. Engels, J. van Enst, H. Falcke, R. A. Fallows, R. Fender, C. Ferrari, W. Frieswijk, M. A. Garrett, J. Griebmeier, J. P. Hamaker, T. E. Hassall, G. Heald, J. W. T. Hessels, M. Hoeft, A. Horneffer, A. van der Horst, M. Iacobelli, N. J. Jackson, E. Juette, A. Karastergiou, W. Klijn, J. Kohler, V. I. Kondratiev, M. Kramer, M. Kuniyoshi, G. Kuper, J. van Leeuwen, P. Maat, G. Macario, G. Mann, S. Markoff, J. P. McKean, M. Mevius, J. C. A. Miller-Jones, J. D. Mol, D. D. Mulcahy, H. Munk, M. J. Norden, E. Orru, H. Paas, M. Pandey-Pommier, V. N. Pandey, R. Pizzo, A. G. Polatidis, W. Reich, A. M. M. Scaife, A. Schoenmakers, D. Schwarz, A. Shulevski, J. Sluman, O. Smirnov, C. Sobey, B. W. Stappers, M. Steinmetz, J. Swinbank, M. Tagger, Y. Tang, C. Tasse, S. t. Veen, S. Thoudam, C. Toribio, R. van Nieuwpoort, R. Vermeulen, C. Vocks, C. Vogt, R. A. M. J. Wijers, M. W. Wise, O. Wucknitz, S. Yatawatta, P. Zarka, and A. Zensus. Discovery of carbon radio recombination lines in absorption towards Cygnus A. *MNRAS*, 437:3506–3515, Feb. 2014. doi: 10.1093/mnras/stt2158.
- S. Ott. The Herschel Data Processing System — HIPE and Pipelines — Up and Running Since the Start of the Mission. In *Astronomical Data Analysis Software and Systems XIX*, volume 434, page 139, Dec. 2010.
- S. Salvini and S. J. Wijnholds. StEFCalAn Alternating Direction Implicit method for fast full polarization array calibration. In *General Assembly and Scientific Symposium (URSI GASS), 2014 XXXIth URSI*, pages 1–4. IEEE, 2014.
- S. Sarkar and D. Tullsen. Compiler techniques for reducing data cache miss rate on a multithreaded architecture. *High Performance Embedded Architectures and Compilers*, pages 353–368, 2008.
- T. Shimwell, H. Röttgering, P. N. Best, W. Williams, T. Dijkema, F. De Gasperin, M. Hardcastle, G. Heald, D. Hoang, A. Horneffer, et al. The LOFAR Two-metre Sky Survey-I. Survey description and preliminary data release. *Astronomy & Astrophysics*, 598: A104, 2017.
- B. Sigoure. OpenTSDB scalable time series database (TSDB), 2012.
- K. Skadron, P. S. Ahuja, M. Martonosi, and D. W. Clark. Branch prediction, instruction-window size, and cache size: Performance trade-offs and simulation techniques. *IEEE Transactions on Computers*, 48(11):1260–1281, 1999.
- O. Smirnov and C. Tasse. Radio interferometric gain calibration as a complex optimization problem. *Monthly Notices of the Royal Astronomical Society*, 449(3):2668–2684, 2015.
- S. Strother, S. La Conte, L. K. Hansen, J. Anderson, J. Zhang, S. Pualapura, and D. Rottenberg. Optimizing the fMRI data-processing pipeline using prediction and reproducibility performance metrics: I. A preliminary group analysis. *Neuroimage*, 23:S196–S207, 2004.
- SURF. Grid at SURFsara. <https://www.surf.nl/en/services-and-products/grid/index.html>, 2018.
- S. J. Tingay, R. Goeke, J. D. Bowman, D. Emrich, S. Ord, D. A. Mitchell, M. F. Morales, T. Boller, B. Crosse, R. Wayth, et al. The Murchison widefield array: The square kilometre array precursor at low radio frequencies. *Publications of the Astronomical Society of Australia*, 30, 2013.
- M. Van Haarlem, M. Wise, A. Gunst, G. Heald, J. McKean, J. Hessels, A. De Bruyn, R. Nijboer, J. Swinbank, R. Fallows, et al. LOFAR: The low-frequency array. *Astronomy & astrophysics*, 556:A2, 2013.
- R. Van Weeren, W. Williams, M. Hardcastle, T. Shimwell, D. Rafferty, J. Sabater, G. Heald, S. Sridhar, T. Dijkema, G. Brunetti, et al. LOFAR facet calibration. *The Astrophysical Journal Supplement Series*, 223(1):2, 2016.
- J.-S. Vöckler, G. Mehta, Y. Zhao, E. Deelman, and M. Wilde. Kick-starting remote applications. In *2nd International Workshop on Grid Computing Environments*, pages 1–8, 2006.
- W. Williams, R. Van Weeren, H. Röttgering, P. Best, T. Dijkema, F. de Gasperin, M. Hardcastle, G. Heald, I. Prandoni, J. Sabater, et al. LOFAR 150-MHz observations of the Boötes field: catalogue and source counts. *Monthly Notices of the Royal Astronomical Society*, 460(3):2385–2412, 2016.
- C. Wu, A. Wicenec, D. Pallot, and A. Checcucci. Optimising NGAS for the MWA Archive. *Experimental Astronomy*, 36(3):679–694, 2013.