

Calibration of LOFAR ELAIS-N1 data in the Amazon cloud.

Jose Sabater¹, Philip Best¹, Susana Sánchez Expósito², Julián Garrido², Lourdes Verdes-Montenegro², Huub Rottgering³ and the ELAIS-N1 consortium.

¹ Institute for Astronomy, University of Edinburgh. Royal Observatory, Blackford Hill, Edinburgh EH9 3HJ, U.K.

² Instituto de Astrofísica de Andalucía (CSIC). Apdo. 3004, 18080, Granada, Spain.

³ Sterrewatch Leiden. P.O. Box 9513 2300 RA Leiden, The Netherlands.

EXECUTIVE SUMMARY

The International LOFAR Telescope [\[1\]](#) is a next-generation software-driven telescope operating in the poorly explored 30–240 MHz frequency range. With an unprecedented field of view, and multiple beams, LOFAR is opening up a completely new phase of radio astrophysics, as well as being both a scientific and technical pathfinder instrument for the SKA. However, the calibration of LOFAR data presents several challenges: a) the effects of the ionosphere introduce systematic errors that can not be overcome with the standard calibration pipelines and prevent continued integration to reach deep fields, b) the huge amount of data and processing power involved requires the use of a parallelizable solution, and, c) the management of the specialized software is not trivial. All of these challenges will also apply to SKA data.

Recently, an innovative calibration strategy, able to reach the theoretical noise limit, has been developed. The cloud infrastructure of Amazon Web Services (AWS) provides a parallelizable (distributed) solution ideal to solve the problems with the data management, the computing constraints and the way the specialized software is dealt with. With this proposal we aim to:

1. Adapt the state of the art calibration pipeline to the Amazon cloud and release it.
2. Identify and publish the best strategy to handle data and resources in the cloud for the calibration.
3. Process the existing data for the ELAIS-N1 field (LOFAR's First Public Deep Field) to produce a wide-field image down to the theoretical limit, facilitating a vast array of science.
4. Investigate the feasibility of AWS to process all LOFAR Survey Fields.

SCIENTIFIC AND TECHNICAL CASE

An important goal that has driven the development of LOFAR since its inception is to explore the low-frequency radio sky through several surveys, in order to advance our understanding of the formation and evolution of galaxies and active galactic nuclei (AGN). The LOFAR Surveys Key Science Project has planned a wedding-cake survey strategy, with three tiers of observations. Tier-1 is the widest tier, with observations across the whole northern sky to a rms noise level of around 100 $\mu\text{Jy}/\text{beam}$ at 150 MHz, while deeper Tier-2 and Tier-3 observations over smaller areas will focus on fields with the highest quality multi-wavelength datasets available.

The ELAIS-N1 field is proposed to be the primary Tier-3 (deepest) field, with an ultimate aim of reaching 10 $\mu\text{Jy}/\text{beam}$ rms noise at 150 MHz. This depth is sufficient to detect Milky-Way-like galaxies out to redshift $z>1$, and starbursting galaxies with star formation rates of 100 solar masses per year back to the earliest epochs of galaxy formation, as well as detecting essentially all radio-AGN across cosmic time. Data of this depth are also of interest for the study of the Epoch of Reionisation, the identification of variable and transient sources, and the study of the origin and evolution of cosmic magnetic fields. For these reasons, beginning in May 2013, a joint effort was begun across four of LOFAR's Key Science Projects to produce a first “LOFAR Public Deep Field”.

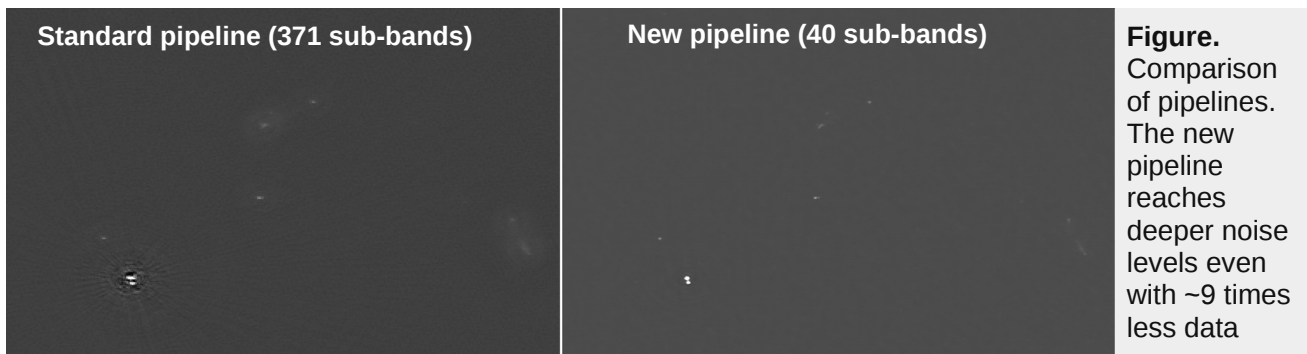
These data will be available to researchers around the world for an extremely broad range of scientific goals. ELAIS-N1 has excellent multi-wavelength data available. This will enable source identification and classification, and the determination of photometric redshifts and other host galaxy properties, critical for maximising the scientific potential of the dataset. To date, 160 hours of data (in 18 datasets) have been taken in this field, with a further 100 hours scheduled for the coming months. These should be sufficient to reach a depth of around 20 $\mu\text{Jy}/\text{beam}$.

The calibration of these LOFAR data requires powerful computing and storage resources and can be considered as an example of data-driven astronomical computing [2]. After the correlation and preprocessing, the data is stored in the Long Term Archive (LTA) and available to the final user for calibration and imaging, but this calibration process still present some challenges:

- The biggest scientific-technical challenge is to deal with the direction-dependent calibration effects due to the ionosphere. Although pipelines exist for global calibration, ultimately this limits the depth achievable, and the systematics prevent continued integration to reach deep fields. We have recently developed the techniques to overcome this, but not yet adapted these to a production-scale. As a side effect to that, the computational needs can change as the pipeline is optimized.
- The size of the data can be considerable: a single observation of 8–10 hours can reach up to 20 TB, and a calibration round of this observation requires several CPU-years of computing. Hence, some kind of parallelization is required.
- The specialized software required for the calibration is under continuous development (8 releases last year) and its installation is not trivial (improving quickly though). This is a problem if the software has to be installed in several nodes or in different operating systems.

A state of the art calibration pipeline [van Weeren et al. in prep.] is under development [3] and the first version is already working. The data are calibrated in projected facets to mitigate the direction dependent effect of the ionosphere. The thermal noise limit is reached with the current version (see Figure). The data is divided in frequency chunks that are calibrated independently (ideal for a cloud solution) and the pipeline is almost fully adapted to run in a distributed environment.

We have investigated the usage of cloud infrastructures to tackle the problems shown: a) a cloud infrastructure provides the high-throughput computing needed for the data calibration; b) the software can be installed once in an image that is instantiated seamlessly in identical virtual machines; c) the size of the infrastructure can be adapted on-demand to the problem to be solved. We installed a preliminary version of the pipeline on [Ibercloud](#) and the [European Grid Initiative Federated Cloud](#) [4] solving the problems with the [software packaging](#) and the creation of images (e.g. <https://appdb.egi.eu/store/vappliance/compss.lofar>). We performed several tests that help us to understand better how to adapt a cloud infrastructure to our problem [5]. However, the lack of an homogeneous implementation of block storage prevented us from reaching a production stage.¹



¹ AWS provides a robust implementation of block storage (EBS).

PLAN AND RESOURCES REQUESTED

We request AWS resources to calibrate the 18 ELAIS-N1 datasets that are already observed.

Pre-configuration of the base infrastructure [1 to 2 months]

First, we will configure the base virtual machine images needed to control and monitor the running of the task, to download the data from the LOFAR LTA and to run the worker instances. They will be based on the Ubuntu 14.04 LTS Amazon Machine Images (AMI) and provisioned using [Ansible](#). All the provisioning steps have already been tested in other infrastructures. We will also test the creation of our own AMIs with all the required software bundled to speed up the deployment process. Finally, we will configure the identities and network, and Route53 will be used to facilitate the identification of the resources.

Resources:

- * ~20 GB of S3 storage space for the AMIs
- * Route53 domain zone

Ingestion of the data in AWS [1 to 2 months; in parallel to other tasks]

We will upload the ELAIS-N1 datasets to the public storage area offered in the call for proposals. If this resource is not available, the data will be sequentially uploaded and calibrated. The final data products will be saved to S3.

Resources:

- * 36 TB of public storage if provided, otherwise 4TB reduced redundancy S3.
- * ~80 GB of S3 standard storage

Calibration runs [10 months]

First, we will test the integration of the pipeline within the AWS infrastructure and our approach to tolerate spot interruptions. The orchestration of the tasks will be done with [IPython Parallel](#) or a job queue system. Once the calibration is integrated, we will test the best strategy for the storage and management of data (instance storage, EBS, etc).² A calibration run is composed of the following steps:

1 - Preprocessing of the calibrator data	This task is already adapted to a distributed computing environment. Each of the 370 data chunks can be processed separately.	* Equivalent to 370 x 12 hour m3.xlarge instances per run.
2 - Preprocessing of the target data	The data are cleaned, joined in 37 data chunks and a first step of the calibration is applied. This task is also adapted to a distributed computing environment.	* Equivalent to 37 x 50 hour m3.2xlarge instances per run.
3 - Calibration of the target data	The facet calibration scheme is applied to 6 final chunks of data in parallel. This last step, at the moment, requires more memory and is still not fully prepared to resume from an intermediate step if the temporary data is deleted. Hence, the data can be stored in EBS volumes to allow to resume after a system shut down.	* Equivalent to 6 x 100 hour r3.2xlarge instances (EBS optimized) per run. * 2 TB of EBS volumes

Control and monitorization of the tasks [In parallel]

The calibration runs will be controlled from an instance that will also monitor the performance and resource consumption of the nodes into a database. This will be crucial to the profiling and optimization of the processes. The instance will also be used to download the data and to gather the final data products to allow their download or storage.

Resources:

- * One m3.large instance
- * Several GB of data out

² The size of resources requested are based on the size of the nodes where the pipeline has been previously run, the recorded resource consumption and the typical length of the tasks on these systems.

Pricing

Table: Maximum price of the resources and a reasonable lowest price that can be achieved.

	Resource	Max.	Low	<p>Cost improvements reflected in the low price column:</p> <p>* Use of spot instances. In regions like N. Virginia it will be possible to run the instances at about 7 to 8 times less price.</p> <p>* Reduce the size of the instances. As the pipeline is improved the usage of memory is reduced and the speed increased.</p> <p>* Storage of the data in the public area. The raw data will be already public by July 2015. It could be stored in the 1 PB public storage area described in the call.</p>
Fixed costs	S3 reduced redundancy; 4 TB	1140	0	
	S3; 100 GB	36	36	
	Data out; ~ 1 TB	100	100	
	Route53; 12 months, 100M requests	50	50	
	Control instance (m3.large) - 12 months	1353	1353	
	EBS persistent volumes - 2 TB	1320	1320	
	Total fixed	3999	2859	
Costs per run	370 m3.xlarge instances - 12 hours	1369	172	
	7 m3.2xlarge instances - 50 hours	1140	143	
	6 r3.2xlarge instances (EBS) - 100 hours	498	110	
	Total per run	3007	425	

We request **\$ 12000** enough to calibrate the 18 datasets available based on the reasonable lower price (adding a ~15% overhead due to the experimentation). Even in the worst case scenario (maximum price), we will be able to calibrate at least 2 datasets, enough to achieve the main technical aims of the proposal.

EXPECTED OUTCOME AND FUTURE

1. A facet calibration work-flow adapted to run in the AWS cloud infrastructure. The software will be released with an open source licence and collaboratively developed in Github.
2. Knowledge on the best strategy to deal with the large amount of data (block storage, transport of data, etc.). The results will be published within our ongoing study of the calibration of LOFAR data in the cloud and a training tutorial will be generated.
3. The final ELAIS-N1 calibrated dataset, images and catalogue will be published, fulfilling our goal of the public release. These resources are expected to lead to the publication of several papers not only from us but also from the wider astronomy community.

The skills and knowledge acquired will be directly relevant for SKA; the new calibration pipeline is especially applicable to SKA-LOW. If the pipeline is successfully implemented, we will aim to calibrate the rest of the ELAIS-N1 data to reach the maximum depth, including the development of new algorithms. The pipeline has also the potential to calibrate the Tier-1 surveys fields (~3000 single 8-hr pointings covering the entire northern sky) down to the theoretical noise level.

Suitability

JS is leading the calibration of the ELAIS-N1 surveys data. He uses AWS since 2010 (personal web and e-mail servers, for [Kaggle competitions](#), etc). JS, SSE and JG have been testing the integration of the LOFAR calibration pipeline in Ibergrid and the EGI Federated cloud. SS, JG and LVM participate in the SKA Science Data Processor consortium. PB is PI of LOFAR UK and the ELAIS-N1 collaboration, LVM coordinates the participation of Spain in SKA and HR is overall PI of the LOFAR Surveys KSP.

References

- [1] [van Haarlem et al. 2013, A&A, 556, 2](#)
- [2] [Berriman et al. 2012, Philosophical Transactions of the Royal Society of London A, 371](#)
- [3] <https://github.com/tammojan/facet-calibration>
- [4] [Sabater et al. 2014; Proceedings of the SEA 2014 meeting in press.](#)
- [5] [Sanchez-Exposito et al. 2015; poster presented at EGI 2015](#)