

# Fast and Reproducible LOFAR Workflows with AGLOW

A.P. Mechev<sup>a,1</sup>, J.B.R. Oonk<sup>a,b,c,2</sup>, T. Shimwell<sup>b,3</sup>, A. Plaat<sup>d,4</sup>, H.T. Intema<sup>a,5</sup>, H.J.A. Röttgering<sup>a,6</sup>

<sup>a</sup>Leiden Observatory,

Leiden University,  
Niels Bohrweg 2 NL-2333 CA  
Netherlands

<sup>1</sup>apmechev at strw.leidenuniv.nl

<sup>2</sup>oonk at strw.leidenuniv.nl

<sup>5</sup>intema at strw.leidenuniv.nl

<sup>6</sup>rottgering at strw.leidenuniv.nl

<sup>b</sup>ASTRON,

Oude Hoogeveensedijk 4,  
7991 PD Dwingeloo  
Netherlands

<sup>3</sup>shimwell at astron.nl

<sup>c</sup>SURFsara

P.O. Box 94613  
1090 GP Amsterdam  
Netherlands

<sup>d</sup>LIACS,

Leiden University,  
Niels Bohrweg 1 NL-2333 CA  
Netherlands

<sup>4</sup>a.plaat at liacs.leidenuniv.nl

**Abstract**—The LOFAR radio telescope creates Petabytes of data per year. This data is important for many scientific projects. The data needs to be efficiently processed within the timespan of these projects in order to maximize the scientific impact. We present a workflow orchestration system that integrates LOFAR processing with a distributed computing platform. The system is named Automated Grid-enabled LOFAR Workflows (AGLOW). AGLOW makes it fast and easy to develop, test and deploy complex LOFAR workflows, and to accelerate them on a distributed cluster architecture. AGLOW provides a significant reduction in time for setting up complex workflows: typically, from months to days. We lay out two case studies that process the data from the LOFAR Surveys Key Science Project. We have implemented these into the AGLOW environment. We also describe the capabilities of AGLOW, paving the way for use by other LOFAR science cases. In the future, AGLOW will automatically produce multiple science products from a single dataset, serving several of the LOFAR Key Science Projects.

**Target:** IEEE International Conference On e-Science Key-words: Workflow management software, Radio Astronomy, Distributed Computing, Big Data applications

## I. INTRODUCTION

Data sets in radio astronomy have increased 1000-fold over the past decade [1]. It is no longer feasible to move, store and process these data sizes at university clusters, nor to process these data manually. LOFAR, the Low-Frequency Array [2] is a modern and powerful radio telescope that creates more than 5 Petabytes of data per year. At present, the majority of LOFAR time is allocated to several Key Science Projects (KSPs) [3]. These projects need to process hundreds or thousands of observations. Typical observations produce approximately 14 TB of archived data. Obtaining high fidelity images from this data requires complex processing steps. To manage and automate the data processing, workflow management software is needed. This software needs to accelerate LOFAR processing on a High Throughput Computing (HTC) cluster while ensuring it is easy to prototype, test, and integrate future algorithms and pipelines.

To automate LOFAR data processing, we have worked with the LOFAR Surveys KSP (SKSP). Together, we designed a software suite that integrates LOFAR software [4] with the Dutch grid infrastructure [5]. This software, based on

Apache Airflow\*, makes it easy to add future science cases, extend and modify pipelines, include data quality checks, and rapidly prototype complex pipelines. For the SKSP use cases, AGLOW achieves a significant reduction in development time: from months to days, allowing researchers to concentrate on data analysis rather than management of processing. Additionally, and perhaps more importantly, the software versions and repositories used are defined within the workflow. This makes reproducibility an integral part of the AGLOW software. Finally, the software is built to leverage an HTC cluster by seamlessly submitting the processing jobs through the cluster's job submission system [6]. The work presented here builds on our previous work parallelizing single LOFAR jobs [7] on a distributed environment. The majority of processing was done at SURFsara at the Amsterdam Science Park [8], which is one of the sites used by the LOFAR Long Term Archive (LTA)<sup>†</sup>. Ongoing efforts include scheduling and processing data at clusters in Poznań in Poland and Jülich in Germany.

**Contributions:** The main features of the AGLOW software are the following:

- Integration of the Grid middleware with Apache Airflow, allowing us to dynamically define, create, submit and monitor jobs on the Dutch national e-infrastructure.
- Integration of the LOFAR (LTA) utilities in Airflow, facilitating pipeline developers to automate staging (moving from tape to disk) and retrieval of LOFAR data.
- Integration of the SURFsara storage with Airflow, making LOFAR pipelines aware of the storage layer available at the Dutch national e-infrastructure.
- Ease of creating simple software blocks, with which users can integrate and test their pipelines.
- Storing all software versions and script repositories as part of the workflow to make LOFAR processing reproducible and portable.

**Outline:** The organization of this manuscript is as follows: We provide background on data processing in radio astronomy and why LOFAR science requires complex workflows

\*<https://airflow.apache.org/>

<sup>†</sup><https://lta.lofar.eu>

80 and cover workflow management algorithms and capabilities  
81 (section II). We discuss related work in workflow management  
82 (section III). In section IV, we introduce our software and two  
83 use cases. Both of our use cases require acceleration at an HTC  
84 cluster and automation by a workflow orchestration software.  
85 We follow these examples with details on the integration  
86 between LOFAR software, LOFAR data and the resources at  
87 SURFsara in Amsterdam in section IV-B2. Finally, we discuss  
88 our results (sect. V) and look ahead to the demands of future  
89 LOFAR projects and upcoming telescopes in section VI.

## 90 II. BACKGROUND

91 This work lies at the intersection of Radio Astronomy  
92 and Computer Science. The goal of the study is to leverage  
93 the flexibility of an industry standard workflow management  
94 software and use CERN's Worldwide Computing Grid\* at  
95 SURFsara [9] to accelerate reproducible processing of LOFAR  
96 data.

97 A single LOFAR surveys observation is recorded in distinct  
98 frequency chunks (henceforth called 'subbands'), each of  
99 which is uploaded to the LTA as a separate file. Some of  
100 the processing steps require the entire frequency information,  
101 while others can run independently and operate on a single  
102 subband. The latter steps can be easily accelerated on an HTC  
103 cluster by taking advantage of the data level parallelism.

104 Multiple scientific projects may desire to run different  
105 processing steps on a single LOFAR observation. To minimize  
106 time spent on retrieving data from the LTA and eliminate re-  
107 processing of data, pipelines for multiple science cases need  
108 to be integrated together. This integration should be done by  
109 a software that encodes the dependencies between different  
110 steps and automatically executes processing steps once their  
111 dependencies have been met. Software packages that solve  
112 these challenges are called 'workflow management software'  
113 (see, e.g., [10]–[12].)

## 114 III. RELATED WORK

115 A workflow is described by a set of tasks. The dependencies  
116 between these tasks are encoded in a Directed Acyclic Graph  
117 (DAG) [13]. This data structure imposes a strict dependency  
118 hierarchy between the tasks [14]. This means that there exists  
119 a well-defined execution order and a well-defined list of  
120 dependencies for each task. The execution order is typically  
121 determined by algorithms such as Kahn's algorithm [15] or a  
122 depth-first search [16].

123 Workflow management software is used in various fields  
124 from research to industry. In biology, gene sequencing and  
125 analysis pipelines require automation of multiple processing  
126 steps. In gene sequencing, Toil<sup>†</sup> has been successfully used  
127 to automate RNA sequence analysis [17]. Additionally, many  
128 software teams in biotech develop their own in-house work-  
129 flow management software [18].

\*<http://wlcg.web.cern.ch/>

<sup>†</sup><https://toil.readthedocs.io>

130 Currently, we can parallelize a single processing step of the  
131 pipeline using the Grid LOFAR Tools (GRID\_LRT<sup>‡</sup>) [7]. The  
132 LOFAR Surveys science cases incorporate multiple steps with  
133 inter-linked dependencies. Resolving these dependencies can  
134 be done efficiently by a comprehensive workflow orchestration  
135 software. The purpose of such software is to resolve dependen-  
136 cies between the multiple tasks in a workflow, execute these  
137 tasks, and track the status, logs, output, and runtime of each  
138 task.

139 In astronomy, workflow systems have been developed that  
140 are telescope specific, such as ESOReflex [19] by the European  
141 Southern Observatory. Other projects, such as astrogrid<sup>§</sup> and  
142 'Workflow 4Ever'<sup>¶</sup>, have either been completed or are no  
143 longer supported. The astrogrid project, for example, was a  
144 collaboration to create standards, infrastructure, and software  
145 for distributed astronomical processing. Its operation phase  
146 spanned 2008-2010. Workflow4Ever, likewise, has been out  
147 of support since 2013. To ensure continuing support for  
148 the LOFAR workflows, we have decided to use a leading  
149 enterprise workflow management software, Airflow.

150 Airflow is an open source Python software package devel-  
151 oped by Airbnb<sup>||</sup> to manage complex workflows. It encodes  
152 workflows in Python and makes it easy to re-use, re-arrange,  
153 schedule and execute blocks in a user-defined workflow.  
154 Airflow is capable of scheduling and executing workflows by  
155 resolving the dependencies between tasks and scheduling these  
156 tasks for execution. The software uses a metadata database<sup>\*\*</sup>  
157 to retain metadata such as task state, execution date, and  
158 output. While Airflow allows building workflows easily from  
159 Python and bash functions, it can easily be extended to support  
160 custom processing scenarios. Additionally, Airflow conforms  
161 to the Common Workflow Language (CWL) [20] standard  
162 using the *cwl-airflow* package [21], meaning it can execute  
163 CWL workflows as well. Finally, Airflow is part of the Apache  
164 incubator and upon certification will receive continual support  
165 by the Apache software foundation<sup>††</sup>.

## 166 IV. AGLOW

167 Complex astronomical pipelines are time consuming to  
168 develop and operate. Furthermore, they may evolve rapidly  
169 to incorporate new processing techniques or requirements.  
170 Migrating these pipelines to a distributed, high throughput  
171 environment is often justified, or even required, in order to  
172 meet the timelines set by scientific projects. The time saved  
173 by running on a cluster must be balanced by the flexibility  
174 and development time required to implement or update the  
175 scientific pipelines. To address these concerns, we have devel-  
176 oped a software package, Automated Grid-enabled LOFAR  
177 Workflows (AGLOW)<sup>‡‡</sup>. AGLOW is based on Airflow and

<sup>‡</sup>[https://github.com/apmechev/GRID\\_LRT](https://github.com/apmechev/GRID_LRT)

<sup>§</sup><http://www.astrogrid.org>

<sup>¶</sup><http://wf4ever.github.io/ro/>

<sup>||</sup><https://www.airbnb.com/>

<sup>\*\*</sup>In our case implemented by PostgreSQL

<sup>††</sup><https://www.apache.org/>

<sup>‡‡</sup><https://github.com/apmechev/AGLOW>

178 the LOFAR software and addresses issues with automation  
179 and acceleration of LOFAR processing.

180 With AGLOW, we can translate LOFAR pipelines into  
181 DAGs. We provide the tools that enable users to easily imple-  
182 ment their LOFAR science pipelines and execute them on a  
183 distributed architecture. Using these tools, the data processing  
184 required by various LOFAR science cases is automated and  
185 accelerated.

### 186 A. AGLOW: Case Study

187 For our case-study, we have chosen two ways to process  
188 LOFAR Surveys data: coverage and depth. The Surveys Key  
189 Science Project (SKSP) [3] is an ambitious project to map  
190 the northern sky at low frequencies using the Dutch LOFAR  
191 stations. These maps will help understand the formation and  
192 evolution of massive black holes, galaxies, clusters of galaxies  
193 and large-scale structure of the Universe.

194 The LOFAR surveys observations consist of several tiers  
195 with the widest Tier (Tier 1) covering the whole sky visible  
196 from the Northern Hemisphere with 3168 observations of 8  
197 hours each [3]. The other tiers (Tier 2, Tier 3) consist of  
198 much longer observations of smaller sections of the sky and  
199 can collect hundreds of hours of data for a single direction.  
200 The deepest single field being analyzed, in collaboration with  
201 the EoR group, is the North Celestial Pole (NCP) field which  
202 has  $\sim 1700$  hrs of observations to date. Processing this data  
203 will create an image with an unprecedented resolution and  
204 sensitivity. Here we have implemented processing pipelines  
205 for both the Tier 1 data and Tier 2 data into AGLOW.

206 The scientific importance of these two examples, as well as  
207 the large processing requirements, make them ideal candidates  
208 for acceleration and automation with AGLOW.

209 1) *Surveys Project: All Sky Survey:* The main driver for  
210 the development of AGLOW and its constituent packages has  
211 been the LOFAR SKSP Project. A typical 8-hour observation  
212 produces 14 TB of data. This data is eventually reduced to  
213 several hundred gigabytes. Data needs to be processed by two  
214 pipelines: first by the Direction Independent (DI) pipeline, and  
215 then by the Direction Dependent (DD) pipeline.

216 We have split the DI pipeline into four stages, and the DD  
217 pipeline into two subsequent stages. Splitting up the pipelines  
218 in stages allows speedup through parallelization for the stages  
219 that can benefit from data-level parallelism. Additionally, this  
220 setup allows fault tolerance and easy re-processing. Current  
221 SKSP processing is easily started by launching a new DAG-  
222 run in Airflow. Importantly, with AGLOW, adding new func-  
223 tionality to the pipeline is easy and can be done at any time  
224 without disrupting current processing.

225 2) *Deeper Surveys Fields:* To create deep images of a  
226 single field, minor modifications were made to the process-  
227 ing pipeline described in the previous section. Scripts were  
228 included to re-align the data in the frequency axis, and the  
229 DD processing steps include an extra final combination step  
230 that stacks multiple observations. Being able to rapidly test  
231 alternative processing strategies is crucial to creating a deep  
232 image of the NCP field. With the success of this project,

233 future deep LOFAR observations will be processed with these  
234 pipelines.

### 235 B. AGLOW: Implementation

236 AGLOW combines the LOFAR software, the Grid LOFAR  
237 Tools (GRID\_LRT), and Airflow to allow automation and  
238 makes large-scale LOFAR processing easily reproducible. The  
239 components of the AGLOW software are shown in Fig. 1.  
240 In this section, we will discuss these components and their  
241 functions.

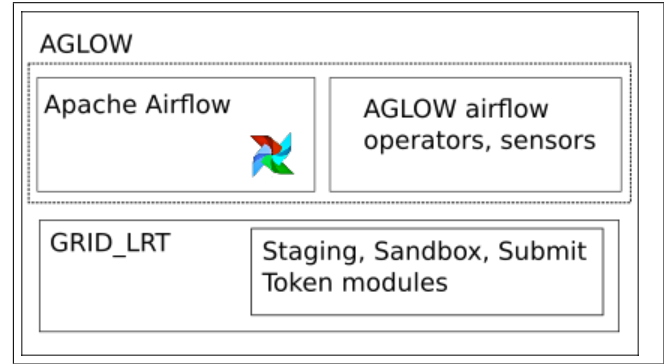


Fig. 1. Design of the AGLOW software, incorporating Airflow, the GRID\_LRT package [7] and custom operators designed to integrate LOFAR software, Grid middleware and dCache storage. GRID\_LRT is a software package developed to parallelize single LOFAR jobs at SURFsara. It contains several modules to help set up, and launch jobs on an HTC cluster at SURFsara. Airflow is a stand-alone package by Airbnb, which is extended with several classes that couple Airflow with the Grid infrastructure. These classes are collectively named the AGLOW operators/sensors.

242 1) *GRID LOFAR Tools and LOFAR software:* We have pre-  
243 viously developed tools to create LOFAR jobs and launch them  
244 on a distributed infrastructure [7]. These tools have matured to  
245 a point where it is easy to both plug and play existing scripts  
246 and extend the framework to add more complex pipelines.  
247 These steps make it possible for a user to batch execute bash  
248 or Python scripts on their LOFAR data in parallel. After the  
249 scripts are executed, the results are uploaded to shared dCache  
250 storage [22] at SURFsara [8].

251 More complex steps use additional Github repositories, such  
252 as the *prefactor\** for direction independent calibration or *DDF-*  
253 *pipeline†* for direction-dependent calibration and imaging. The  
254 sequence of steps is encoded in parameter-set files (*parsets*),  
255 which can be modified and dropped into AGLOW depending  
256 on the processing requirements.

257 With AGLOW, we can easily include the *DDF-pipeline* and  
258 *prefactor* repositories, as well as any other scripts. Since these  
259 scripts are tracked by git [23], a full commit and branch  
260 history of the scripts is available. We use this history to make  
261 processing reproducible, by using the same git-commit for all  
262 LOFAR datasets.

263 In addition to these script repositories, we have integrated  
264 the most common software packages used to process LOFAR  
265 data with AGLOW. These are the Default Pre-Processing  
266 Pipeline (DPPP) [4], the LOFAR Solutions Tool (LoSoTo),

\*<https://github.com/lofar-astron/prefactor>

267 WSclean [24], AOflagger [25], CASA [26], pyBDSF [27],  
 268 DDFacet [28] and KillMS [29], [30].

269 2) *Extending Airflow*: Two types of modifications were  
 270 made to Airflow to allow processing on a Grid environment.  
 271 First, functions were added to check the number of files  
 272 located in intermediate grid storage. We use this to decide  
 273 whether to stage files, or whether enough files have been  
 274 successfully processed by a previous task.

275 Second, more complex tasks were implemented as Airflow  
 276 operators or sensors (Figure 2). These tasks include creating  
 277 job description files, setting up job scripts, launching jobs  
 278 using the gLite workload management system and monitoring  
 279 the status of these jobs. Future additions will include operators  
 280 that evaluate the current cluster workload and make decisions  
 281 on location to launch the data processing. With the AGLOW  
 282 package, such tasks are easy to implement without modifying  
 283 or interrupting processing. This leads to an easily reproducible,  
 284 intelligent scientific processing that is also efficiently executed  
 285 and requires minimal interaction. The operators and sensors  
 286 added to Airflow are shown in Fig. 2.

287 Using AGLOW to accelerate the execution of a pipeline  
 288 requires deciding how to split the processing to benefit from  
 289 parallelization. Once the steps to be parallelized are selected,  
 290 users can add git repositories of scripts to the configuration  
 291 file. Next, each step is added to a Python script called the DAG  
 292 file. This file is placed in the Airflow’s dags folder, which adds  
 293 it to AGLOW. To migrate LOFAR workflows to a new server,  
 294 the DAG and configuration files need to be transferred to the  
 295 new AGLOW instance.

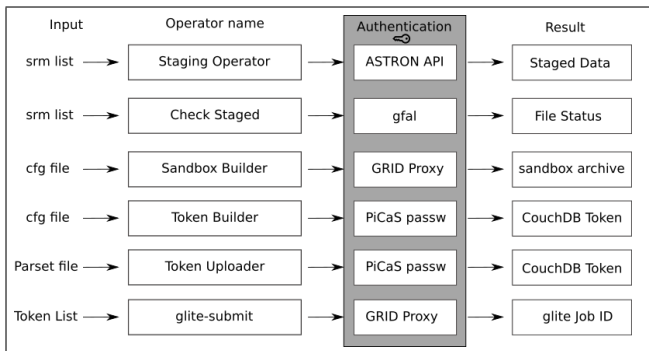


Fig. 2. Airflow Operators for Staging LOFAR data, creating job descriptions and submitting jobs to the Dutch grid. On the left is the input given to each operator. ‘SRM lists’ are lists of links to data at the LOFAR LTA or located on the SURFsara dCache storage. Parsets are files specific to ‘prefactor’ and ‘DDF-pipeline’ and define the processing for each pipeline step. Finally, the ‘Sandbox’ and ‘Token’ operators read their parameters from a configuration file. The use of a scripts sandbox and job description tokens is detailed in our previous work [7].

### 296 C. AGLOW: Jobs

297 Once LOFAR observations are downloaded from the LTA,  
 298 they are typically processed with several packages before  
 299 producing a science ready dataset. We have integrated these

<sup>†</sup><https://github.com/mhardcastle/ddf-pipeline>. DDF-pipeline is a leading example of a Direction Dependent calibration pipeline used for LOFAR data. It uses DDFacet [28], KillMS [29] and to create high quality images.

packages with Airflow to make it easy to create complex  
 LOFAR workflows.

Each of the processing steps above requires extra set-up  
 to process on the Dutch Grid infrastructure. The job scripts  
 setup, job description, and job submission are done by the  
 GRID\_LRT package [7]. With AGLOW, we automate this  
 setup, enabling users to focus on developing more compre-  
 hensive data processing pipelines. Below we outline several  
 possible steps a user can use in their pipeline.

1) *DPPP Parset*: The DPPP software is used extensively  
 in LOFAR data processing. It has many capabilities such as  
 flagging bad data, averaging data in time and frequency, and  
 calibrating the data with a sky-model.

The input parameters of this software are stored in a text  
 file called a parset. The input data and the DPPP parset  
 are sufficient to define a DPPP execution step. As noted in  
 section II, LOFAR data is split in frequency into subbands.  
 Much of the DPPP processing, such as averaging and flagging,  
 can be done independently for each subband, thus they can  
 be processed on independent machines. This parallelization  
 makes these steps a perfect candidate for an HTC cluster. For  
 a dataset that is split into 244 subbands, 244 jobs are launched  
 concurrently.

In Airflow, the DPPP parset task is encoded in a DAG (Fig.  
 3). The DPPP DAG is a linear workflow that consists of the  
 ‘sandbox’ setup, creation of the job-description documents, up-  
 loading of the DPPP parset and job launching and monitoring.

2) *WSclean Job*: The WSclean [24] package is used to  
 create an image from a LOFAR dataset. This software has a  
 very wide range of parameters options, however, it cannot take  
 a parset file as an input. Instead, the parameters are specified  
 in the command line. In the AGLOW implementation, we parse  
 all the command-line parameters from a text file, referred to as  
 the ‘wsclean parset’. This file is added to the jobs in the same  
 way as the DPPP parset, i.e. using the *Token Uploader* Opera-  
 tor. The DAG for the wsclean software uses the same blocks as  
 the DPPP DAG, with different configuration and parset files.  
 The reuse of Airflow operators makes maintainability of all  
 tasks easier.

### D. Shell/Python Script

Users that require the run of multiple software packages on  
 a single dataset can craft a custom shell or Python script that  
 executes these steps using the LOFAR tools during a single  
 distributed job. This option increases flexibility and minimizes  
 the overhead associated with scheduling and running multiple  
 jobs in sequence. At the workflow orchestration level, we use  
 the same Airflow operators as the above tasks. The script  
 is uploaded to the job description database using the *Token  
 Uploader* Operator. It is executed once the jobs are launched.

Currently only the LOFAR Spectroscopy project uses cus-  
 tom shell scripts to process LOFAR data. A recent study  
 of carbon recombination lines used a custom bash script to

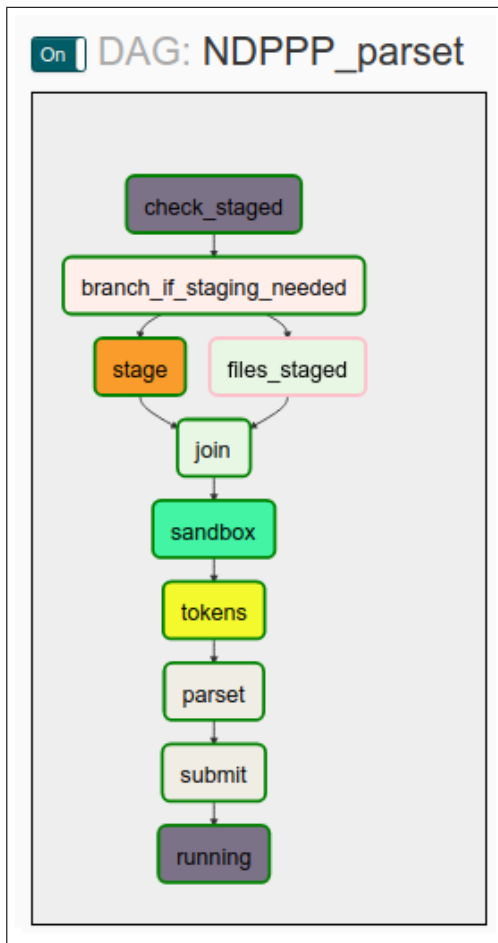


Fig. 3. Render of the DPPP parset DAG in the Airflow User Interface. This view shows the setup and submission steps. Even this simple DAG can include branching options such as the `branch_if_staging_needed` task which checks if the data is not staged and stages it. All of the operators in this figure are part of the AGLOW software. Their inputs and outputs are shown in Fig. 2. Using configuration files, the NDPPP DAG can be used by different users for different science cases making it portable and maintainable. These features make reproducible science with LOFAR data easy.

352 calibrate and image LOFAR data on the SURFsara GINA\*  
 353 cluster [31].

#### 354 E. Prefactor parset

355 The input to the prefactor pipeline software is a parset  
 356 file which describes a linear workflow. The description of  
 357 this workflow consists of a list of processing steps and  
 358 their associated parameters. The ‘prefactor’ package uses the  
 359 LOFAR software to do the direction-independent calibration  
 360 of the archived LOFAR datasets. Prefactor steps are executed  
 361 by the generic pipeline framework [4]. While this framework  
 362 can run a sequential pipeline, it is not capable of conditional  
 363 branching nor parallelization on all cluster architectures. The  
 364 original goal of the GRID\_LRT software was to tackle the

\*The GINA cluster is an HTC cluster located at SURFsara integrated with the Dutch Grid initiative. It supports massively parallel processing which is required to efficiently process LOFAR data with *prefactor*.

365 parallelization challenge while AGLOW solves the additional  
 366 challenge of pipeline management.

367 We have already processed more than 50 datasets through  
 368 the ‘prefactor’ DAG using AGLOW. The full ‘prefactor’  
 369 pipeline is shown in figure 4. This DAG shows the four  
 370 processing steps as well as additional Python operators that  
 371 manage the staging and result verification.

#### 372 F. DDF-pipeline

373 The final AGLOW DAG is the implementation of the DDF-  
 374 pipeline repository which is a pipeline that is extensively used  
 375 by the LOFAR surveys KSP and is described in detail in [3].  
 376 This pipeline operates on the products of the prefactor pipeline  
 377 and consists of a series of calibration and imaging loops with  
 378 the objective of creating a final science quality image. For each  
 379 of these loops the majority of the processing time is spent in  
 380 DDFacet [28] and KillMS [29], [30] steps that perform the  
 381 direction-dependent imaging and calibration respectively.

382 In total, DDF-pipeline takes  $\sim 4$  days of processing to  
 383 complete. As DDF-pipeline creates large intermediate files we  
 384 have so far not divided the pipeline into too many steps to  
 385 avoid filling the storage on the GINA cluster. However, we  
 386 have split the pipeline into two steps and there is further  
 387 potential for parallelization that will be implemented in the  
 388 future.

#### 389 G. Linking Multiple Jobs

390 Pre-processing of LOFAR SKSP data can be done by a  
 391 single DPPP task, with 244 jobs running in parallel. More  
 392 complex LOFAR pipelines will include multiple processing  
 393 tasks as well as tasks responsible for job setup. Therefore,  
 394 it is important to facilitate running multi-step pipelines with  
 395 AGLOW.

396 Creating workflows by defining dependencies between tasks  
 397 is a core Airflow capability. We use this functionality to link  
 398 multiple steps of a LOFAR pipeline together. In the SKSP  
 399 pipeline, we take advantage of the data level parallelism for  
 400 the initial processing steps for the calibrator and target. The  
 401 other two steps are run as a single grid job. Switching the  
 402 parallelization for each step is done by changing the number  
 403 of datasets per node parameter in the configuration file for  
 404 each step.

## 405 V. RESULTS AND DISCUSSIONS

406 The implementation of AGLOW makes it possible to effi-  
 407 ciently process LOFAR data with minimal user interaction.  
 408 The scheduling algorithm automatically launches pipelines,  
 409 meaning that there is little time spent between runs. Addition-  
 410 ally, controlling/fixing the version of the scripts is done by  
 411 specifying the commit of each script repository. This makes  
 412 data processing easily reproducible. Once the dependencies  
 413 of multiple science pipelines have been encoded in a DAG,  
 414 Airflow efficiently executes this DAG, running tasks in parallel  
 415 where possible.

416 The first LOFAR processing pipeline integrated with  
 417 AGLOW was a single linear workflow, with only one sub-  
 418 mission to the compute cluster. This workflow is used to

419 reduce the data size making data retrieval to research institutes  
420 less time consuming. We offer this workflow as a service to  
421 LOFAR users who do not have a high-bandwidth connection  
422 to the LOFAR Archive.

423 A more complex pipeline was implemented: the LOFAR  
424 direction independent calibration pipeline (‘prefactor’). The  
425 scientific importance and complexity of this pipeline make  
426 it a good case study for the capabilities of the AGLOW  
427 software. We show that AGLOW’s design allows integration  
428 of more complex data processing workflows with the Dutch  
429 Grid resources. These workflows can be either used by PIs  
430 of LOFAR projects or offered as a processing service to the  
431 wider astronomical community.

432 An important feature of AGLOW is the loose coupling  
433 between pipeline logic, software versions, pipeline param-  
434 eters, and datasets. The goal of this decoupling is to give  
435 users complete control over all the processing variables. With  
436 AGLOW, one can develop the pipeline logic independently of  
437 the LOFAR software versions and conversely update the LO-  
438 FAR software and script repositories independently from the  
439 pipeline logic. Finally, the Airflow operators are themselves  
440 decoupled from the scientific pipelines. As these operators  
441 are reused, this decoupling makes them easy to maintain and  
442 extend.

443 In large part thanks to their flexibility, automation, and Grid  
444 integration, AGLOW and GRID\_LRT have become a standard  
445 part of the Direction Independent processing for the LOFAR  
446 SKSP project.

## 447 VI. CONCLUSIONS

448 In this work, we have detailed a comprehensive workflow  
449 management software for processing radio astronomy data on  
450 a distributed infrastructure. We leverage an industry standard  
451 workflow management software, Airflow. Using its capabili-  
452 ties, we make it possible to build, test, automate and deploy  
453 LOFAR pipelines on short timescales, generally from months  
454 to days. With the flexibility of Airflow’s Python and Bash  
455 operators, users can design their own workflows, as well as  
456 co-ordinate more complex science cases. In this way, AGLOW  
457 facilitates reproducible processing of scientific data. In the  
458 future, AGLOW will support additional LOFAR science cases  
459 including Long Baselines and Spectroscopy. In this article,  
460 we have described our implementation of the data processing  
461 pipelines used by the LOFAR Surveys Key Science Project.

462 Future work includes further de-coupling of the Grid-setup  
463 and pipeline logic. We will do this by creating ‘sub-dags’  
464 (details in VI-A) for each type of LOFAR jobs. Using these  
465 sub-dags will reduce the complexity of scientific workflows  
466 while also making the code even more reusable and thus easier  
467 to maintain and upgrade. Efforts to integrate processing at  
468 the other two LTA sites, (Jülich and Poznań) have already  
469 started with ‘prefactor’ runs being performed on Jülich using  
470 a modified version of the SKSP workflow. The software also  
471 currently works on the Eagle cluster at Poznań. Combining  
472 the Jülich and SURFsara workflows will be done in the future

so that AGLOW can track and start processing at multiple  
clusters.

473  
474  
475 Finally, AGLOW can be used as a ‘LOFAR As A Service’  
476 model. In this model, users only provide an observation ID  
477 and processing parameters and receive the final results upon  
478 job completion. This model will build upon previous success  
479 offering LOFAR processing to users without login to the  
480 GINA cluster [32]. This previous work was already useful  
481 for studying radio absorption in Cassiopeia A [33] and a  
482 ‘data-to-images’ service will be valuable to the whole LOFAR  
483 community.

484 Our experience with automating LOFAR scientific work-  
485 flows on a distributed architecture will be valuable when  
486 setting up data processing for future Radio Telescopes such  
487 as the Square Kilometer Array [34].

## 488 APPENDIX

489 The LOFAR SKSP workflow is shown in Figure 4. This  
490 figure shows how reuse of the staging, setup operators, and  
491 glite-wms sensors makes maintainability easy and allows rapid  
492 prototyping of complex pipelines.

493 This workflow additionally takes advantage of Airflow’s  
494 *PythonOperator* to check if the LOFAR data is on disk at the  
495 archive and whether all final products were uploaded by each  
496 step. AGLOW also allows for staging the calibrator and target  
497 files concurrently. When the data is staged, Airflow continues  
498 with the processing of that data.

### 499 A. Sub-DAG

500 Airflow allows developers to include entire DAGs as a single  
501 task in their workflow. Airflow can trigger a DAG execution  
502 based on parameters provided by the parent DAG. This feature  
503 makes it possible to concatenate short, commonly used tasks  
504 into DAGs and call them in a parent workflow. Using sub-  
505 DAGS makes the code more maintainable and easy to use,  
506 while it makes workflows simpler. For LOFAR, Sub-DAGs  
507 are used to automate job submission, making the resulting  
508 scientific workflows simpler.

## 509 ACKNOWLEDGEMENT

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

513 The processing and storage functionality that has made this  
514 project possible was enabled by SURF Cooperative through  
515 grant e-infra 160022 & 160152. The LOFAR software and  
516 dedicated reduction packages on <https://github.com/apmechev/>  
517 GRID\_LRT were deployed on the e-infrastructure by the  
518 LOFAR e-infra group, consisting of J.B.R. Oonk (SURFsara  
519 & Leiden Observatory), A.P. Mechev (Leiden Observatory).

520 This paper is based on data obtained with the International  
521 LOFAR Telescope (ILT) under project codes LC2\_038 and  
522 LC3\_008. LOFAR (van Haarlem et al. 2013) is the Low-  
523 Frequency Array designed and constructed by ASTRON. It  
524 has observing, data processing, and data storage facilities in  
525 several countries, which are owned by various parties (each

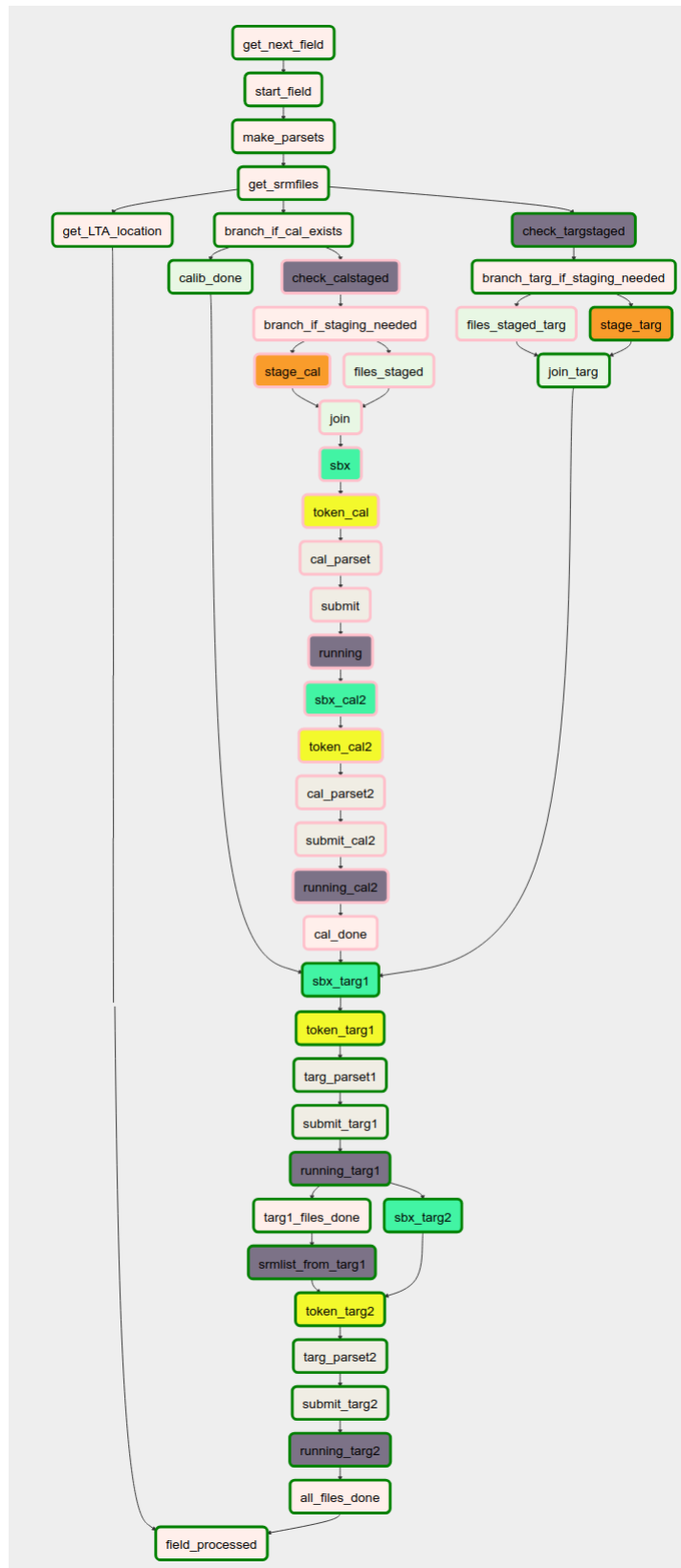


Fig. 4. Workflow for the prefactor pipeline. Here we show the reuse of AGLOW operators for the four prefactor steps. In addition to the LOFAR processing, we also have conditional operators to skip processing of the calibrator if it has been previously processed. This is done by the 'branch\_if\_cal\_exists' task. We also have a final step that checks if all the results have been uploaded, done by the 'all\_files\_done' task. Likewise, quality checks can be added in this workflow wherever needed.

526 with their own funding sources), and which are collectively  
527 operated by the ILT foundation under a joint scientific policy.  
528 The ILT resources have benefited from the following recent  
529 major funding sources: CNRS-INSU, Observatoire de Paris  
530 and Université d'Orléans, France; BMBF, MIWF-NRW, MPG,  
531 Germany; Science Foundation Ireland (SFI), Department of  
532 Business, Enterprise and Innovation (DBEI), Ireland; NWO,  
533 The Netherlands; The Science and Technology Facilities  
534 Council (STFC), UK

## 535 REFERENCES

- 536 [1] J Sabater, S Sánchez-Expósito, J Garrido, JE Ruiz, PN Best, and  
537 L Verdes-Montenegro. Calibration of radio-astronomical data on the  
538 cloud. LOFAR, the pathway to SKA. In *Highlights of Spanish Astro-*  
539 *physics VIII*, pages 840–843, 2015.
- 540 [2] MP Van Haarlem, MW Wise, AW Gunst, George Heald, JP McKean,  
541 JWT Hessels, AG De Bruyn, Ronald Nijboer, John Swinbank, Richard  
542 Fallows, et al. LOFAR: The low-frequency array. *Astronomy &*  
543 *astrophysics*, 556:A2, 2013.
- 544 [3] TW Shimwell, HJA Röttgering, Philip N Best, WL Williams, TJ Di-  
545 jkema, F De Gasperin, MJ Hardcastle, GH Heald, DN Hoang, A Horne-  
546 fer, et al. The LOFAR Two-metre Sky Survey-I. Survey description and  
547 preliminary data release. *Astronomy & Astrophysics*, 598:A104, 2017.
- 548 [4] Dijkma Tammo Jan. Lofar imaging cookbook. Available at [http://www.](http://www.astron.nl/sites/astron.nl/files/cms/lofar_imaging_cookbook_v19.pdf)  
549 [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).
- 550 [5] Jeff Templon and Jan Bot. The dutch national e-infrastructure. In  
551 *International Symposium on Grids and Clouds (ISGC)*, volume 13, 2016.
- 552 [6] Erwin Laure, A Edlund, F Pacini, P Buncic, S Beco, F Prelz,  
553 A Di Meglio, O Mulmo, M Barroso, Peter Z Kunszt, et al. Middleware  
554 for the next generation grid infrastructure. Technical report, CERN,  
555 2004.
- 556 [7] A. Mechev, J. B. R. Oonk, A. Danezi, T. W. Shimwell, C. Schrijvers,  
557 H. Intema, A. Plaat, and H. J. A. Rottgering. An Automated Scal-  
558 able Framework for Distributing Radio Astronomy Processing Across  
559 Clusters and Clouds. In *Proceedings of the International Sympo-*  
560 *sium on Grids and Clouds (ISGC) 2017, held 5-10 March, 2017 at*  
561 *Academia Sinica, Taipei, Taiwan (ISGC2017)*. Online at [https://pos.](https://pos.sissa.it/cgi-bin/reader/conf.cgi?confid=293,id.2)  
562 [sissa.it/cgi-bin/reader/conf.cgi?confid=293, id.2](https://pos.sissa.it/cgi-bin/reader/conf.cgi?confid=293,id.2), page 2, March 2017.
- 563 [8] SURF. Grid at SURFSara. [https://www.surf.nl/en/services-and-products/](https://www.surf.nl/en/services-and-products/grid/index.html)  
564 [grid/index.html](https://www.surf.nl/en/services-and-products/grid/index.html), 2018.
- 565 [9] Jamie Shiers. The worldwide lhc computing grid (worldwide lcg).  
566 *Computer physics communications*, 177(1-2):219–223, 2007.
- 567 [10] Ilkay Altintas, Chad Berkley, Efrat Jaeger, Matthew Jones, Bertram  
568 Ludascher, and Steve Mock. Kepler: An extensible system for design and  
569 execution of scientific workflows. In *Scientific and Statistical Database*  
570 *Management, 2004. Proceedings. 16th International Conference on*,  
571 pages 423–424. IEEE, 2004.
- 572 [11] David Churches, Gabor Gombas, Andrew Harrison, Jason Maassen,  
573 Craig Robinson, Matthew Shields, Ian Taylor, and Ian Wang. Pro-  
574 gramming scientific and distributed workflow with triana services.  
575 *Concurrency and Computation: Practice and Experience*, 18(10):1021–  
576 1037, 2006.
- 577 [12] Ji Liu, Esther Pacitti, Patrick Valduriez, and Marta Mattoso. A survey  
578 of data-intensive scientific workflow management. *Journal of Grid*  
579 *Computing*, 13(4):457–493, 2015.
- 580 [13] Lin Wang. Directed acyclic graph. In *Encyclopedia of Systems Biology*,  
581 pages 574–574. Springer, 2013.
- 582 [14] David J Pearce and Paul HJ Kelly. A dynamic topological sort algorithm  
583 for directed acyclic graphs. *Journal of Experimental Algorithmics (JEA)*,  
584 11:1–7, 2007.
- 585 [15] A. B. Kahn. Topological sorting of large networks. *Commun. ACM*,  
586 5(11):558–562, November 1962.
- 587 [16] Jianjun Zhou and Martin Müller. Depth-first discovery algorithm for  
588 incremental topological sorting of directed acyclic graphs. *Information*  
589 *Processing Letters*, 88(4):195–200, 2003.
- 590 [17] John Vivian, Arjun Arkal Rao, Frank Austin Nothaft, Christopher  
591 Ketchum, Joel Armstrong, Adam Novak, Jacob Pfeil, Jake Narkizian,  
592 Alden D Deran, Audrey Musselman-Brown, et al. Toil enables repro-  
593 ducible, open source, big biomedical data analyses. *Nature biotechnol-*  
594 *ogy*, 35(4):314, 2017.
- [18] Paolo Di Tommaso, Maria Chatzou, Evan W Floden, Pablo Prieto Barja,  
Emilio Palumbo, and Cedric Notredame. Nextflow enables reproducible  
computational workflows. *Nature biotechnology*, 35(4):316–319, 2017.
- [19] W. Freudling, M. Romaniello, D. M. Bramich, P. Ballester, V. Forchi,  
C. E. García-Dabó, S. Moehler, and M. J. Neeser. Automated data  
reduction workflows for astronomy. The ESO Reflex environment.  
*Astronomy & Astrophysics*, 559:A96, November 2013.
- [20] Peter Amstutz, Michael R Crusoe, Nebojša Tijanić, Brad Chapman, John  
Chilton, Michael Heuer, Andrey Kartashov, Dan Leehr, Hervé Ménager,  
Maya Nedeljkovich, et al. Common workflow language, v1. 0. 2016.
- [21] Michael Kotliar, Andrey Kartashov, and Artem Barski. Cwl-airflow: a  
lightweight pipeline manager supporting common workflow language.  
*bioRxiv*, page 249243, 2018.
- [22] Patrick Fuhrmann and Volker Güllow. dcache, storage system for the  
future. In *European Conference on Parallel Processing*, pages 1106–  
1113. Springer, 2006.
- [23] Linus Torvalds and Junio Hamano. Git: Fast version control system.  
*URL http://git.scm.com*, 2010.
- [24] A. R. Offringa, B. McKinley, N. Hurley-Walker, F. H. Briggs, R. B.  
Wayth, D. L. Kaplan, M. E. Bell, L. Feng, A. R. Neben, J. D. Hughes,  
J. Rhee, T. Murphy, N. D. R. Bhat, G. Bernardi, J. D. Bowman, R. J.  
Cappallo, B. E. Corey, A. A. Deshpande, D. Emrich, A. Ewall-Wice,  
B. M. Gaensler, R. Goeke, L. J. Greenhill, B. J. Hazelton, L. Hindson,  
M. Johnston-Hollitt, D. C. Jacobs, J. C. Kasper, E. Kratzenberg, E. Lenc,  
C. J. Lonsdale, M. J. Lynch, S. R. McWhirter, D. A. Mitchell, M. F.  
Morales, E. Morgan, N. Kudryavtseva, D. Oberoi, S. M. Ord, B. Pindor,  
P. Procopio, T. Prabu, J. Riding, D. A. Roshi, N. U. Shankar, K. S.  
Srivani, R. Subrahmanyam, S. J. Tingay, M. Waterson, R. L. Webster,  
A. R. Whitney, A. Williams, and C. L. Williams. WSCLEAN: an  
implementation of a fast, generic wide-field imager for radio astron-  
omy. *Monthly Notices of the Royal Astronomical Society*, 444:606–619,  
October 2014.
- [25] A. R. Offringa, J. J. van de Gronde, and J. B. T. M. Roerdink. A  
morphological algorithm for improving radio-frequency interference  
detection. *Astronomy & astrophysics*, 539:A95, March 2012.
- [26] JP McMullin, B Waters, D Schiebel, W Young, and K Golap. Casa  
architecture and applications. In *Astronomical data analysis software*  
*and systems XVI*, volume 376, page 127, 2007.
- [27] Niruj Mohan and David Rafferty. Pybdsf: Python blob detection and  
source finder. *Astrophysics Source Code Library*, 2015.
- [28] C Tasse, B Hugo, M Mirmont, O Smirnov, M Atemkeng, L Bester,  
E Bonnassieux, MJ Hardcastle, R Lakhoo, J Girard, et al. Facetting  
for direction-dependent spectral deconvolution. *arXiv preprint*  
*arXiv:1712.02078*, 2017.
- [29] OM Smirnov and Cyril Tasse. Radio interferometric gain calibration  
as a complex optimization problem. *Monthly Notices of the Royal*  
*Astronomical Society*, 449(3):2668–2684, 2015.
- [30] Cyril Tasse. Nonlinear kalman filters for calibration in radio interfer-  
ometry. *Astronomy & Astrophysics*, 566:A127, 2014.
- [31] P. Salas, J. B. R. Oonk, R. J. van Weeren, F. Salgado, L. K. Morabito,  
M. C. Toribio, K. Emig, H. J. A. Röttgering, and A. G. G. M. Tielens.  
LOFAR observations of decameter carbon radio recombination lines  
towards Cassiopeia A. *Monthly Notices of the Royal Astronomical*  
*Society*, 467:2274–2287, May 2017.
- [32] J.B.R. Oonk, A.P. Mechev, A Danezi, C Schrijvers, and T.W. Shimwell.  
Radio astronomy on a distributed shared computing platform: The  
LOFAR case.
- [33] M Arias, J Vink, F De Gasperin, P Salas, JBR Oonk, RJ Van Weeren,  
AS Van Amesfoort, J Anderson, R Beck, ME Bell, et al. Low-frequency  
radio absorption in cassiopeia a. *Astronomy & Astrophysics*, 612:A110,  
2018.
- [34] Peter E Dewdney, Peter J Hall, Richard T Schilizzi, and T Joseph LW  
Lazio. The square kilometre array. *Proceedings of the IEEE*,  
97(8):1482–1496, 2009.