| Author: Auke Latour, Kjeld v.d. Schaaf | Date of issue: 2007-03-07<br>Kind of issue: Public | Scope: CEP<br>Doc.id: LOFAR-ASTRON-ADD-012 | |
|---|---|---|---|
| | Status: Final<br>Revision nr: 3.0 | | |

# LOFAR Central Processing Facility
# Architecture description

| **Verified:** | | | |
|---|---|---|---|
| Name | Signature | Date | Rev.nr. |
| | | | |

| **Accepted:** | | |
|---|---|---|
| Work Package Manager | System Engineering Manager | Program Manager |
| | André Gunst | Jan Reitsma |

## Distribution list:

| ASTRON: | For Information: |
|---|---|
| LOFAR Project<br>    Mark Bentum<br>    Ger de Bruyn<br>    André Gunst<br>    Michiel van Haarlem<br>    Hanno Holties<br>    Auke Latour<br>    Ronald Nijboer<br>    Ruud Overeem<br>    Jan Reitsma<br>    John Romein<br>    Gijs Schoonderbeek<br>    Michael Wise | Review committee<br>    E. van den Heuvel (chair, University of Amsterdam)<br>    A. Berg (SARA)<br>    W. Brouw (University of Groningen)<br>    R. Schilizzi (SKA-ISPO/U. Leiden)<br>    C.H. Slump (University of Twente)<br>    A.B. Smolders (NXP Semiconductors)<br>    E. Stolp |

## Document history:

| Revision | Date | Section | Page(s) | Modification |
|---|---|---|---|---|
| 0.1 | 2004-05-19 | | - | Creation |
| 0.2 | 2004-05-28 | | | Updated based on intial discussion with BG/L development team |
| 0.3 | 2004-06-11 | | | Added Storage system description, MAC system and appendix C |
| 0.31 | 2004-06-13 | 1.3 | | Added reference to application definition documents |
| 1.0 | 2004-07-06 | | | First release<br>moved appendices A,B and D to draft detailed design documents |
| 2.0 | 2005-10-13 | All | All | Update for CDR review |
| 2.1 | 2007-03-22 | All | All | Auke Latour: Major Update. Process descriptions for the KSP observation modes added (Chapter 5). Many needless details removed. |
| 3.0 | 2007-03-27 | All | All | Auke Latour: Many changes after internal review:<br>• Update of observation modes<br>• Appendix about requirements coverage<br>• Many minor changes |

# Table of contents

# 1. Introduction

LOFAR is a radio telescope with antenna stations distributed over an area of about 100 km in diameter and a CEntral Processing facility (CEP) where the station data are collected and processed. The stations are located in the north of the Netherlands. The central processing facility is located in Groningen.

The stations contain many low-budget antennas and analogue and digital pre-processing hardware for processing the antenna data. A Wide Area Network (WAN) connects the stations to the central processing facility. The central processing facility combines the data of the stations and makes the data usable for the astronomer. The User Software Group (USG) develops software for additional scientific processing of the CEP output. The observations are specified with the Specification And Scheduling application. The Monitoring And Control system (MAC) coordinates the activities of the stations and the central processing facility. With the Navigator application, the LOFAR system can be inspected and interventions are possible from there (see Figure 1).

Beside antennas for astronomy, LOFAR will contain sensors for other science projects. Among others, geophones will be placed at the stations to monitor the deep soils underneath the telescope, but those aspects are not covered by this document.



**Figure 1** *Global LOFAR overview*

One of the major differences with traditional telescopes is that LOFAR needs to be extremely flexible (see for example the specifications in [9]). That is why LOFAR is referred to as a software radio telescope and the software already was one of the main focuses at the start of the project. The basic idea was to start software development at an early stage without knowing exactly on what platform it will finally run. One of the main advantages of this is that the project can buy off-the-shelf hardware at a late time in the project and

consequently does not end up with expensive obsolete hardware when the operations of the instrument are about to start.

## 1.1. Purpose of this document

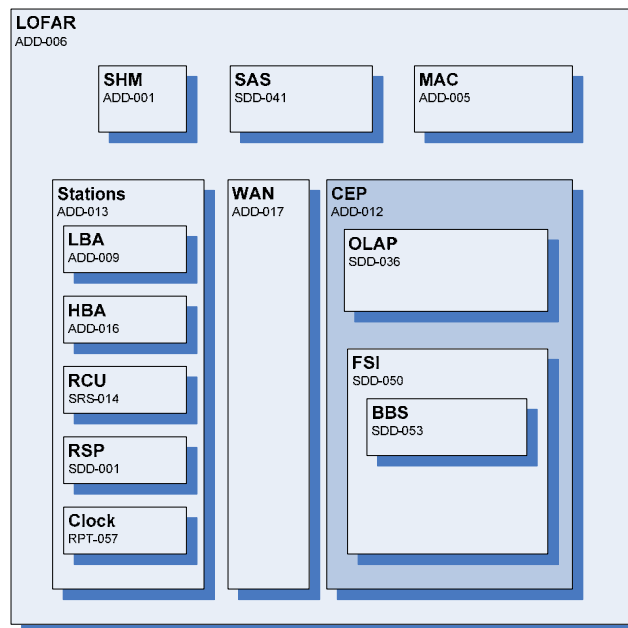This document gives a description of the CEP facility hardware and software architectures. The CEP facility is described from multiple viewpoints in order to focus on specific aspects of the system. This document should provide enough information for the detailed design and construction of the subsystems.

## 1.2. Context of this document



**Figure 2** *Hierarchical relations between the LOFAR architectural documents, and the place of the CEP ADD in this hierarchy. The document numbers refer to LOFAR-ASTRON documents*

This document describes the architecture of the CEntral Processor subsystem, which is part of the LOFAR imaging radio telescope. The top-level architecture of LOFAR is described in ADD-006 [7]. The other subsystems for LOFAR have their own Architectural Design documents (see Figure 2). The subsystems of CEP also have their own architectural document, where the design is described on a more detailed level.

## 1.3. Applicable documents

The requirements that are the basis for this architecture document can be found in the CEP Subsystem Requirements specification [8]. The architectural process and reasoning leading to this architecture is reported in [3]. The IBM Blue Gene/L architecture is described in [1] and is assumed to be known in the remainder of this document.

## 1.4. Structure of this document

This document is a 'living' document. LOFAR is a long term project, and the architecture evolves with progressing insight. Chapter 2 describes the current status of this architecture document. Chapter 3 is the actual starting point of this architecture description. In this chapter you can find a description of the interaction of CEP with its environment, and the responsibilities of the main components of the hardware architecture. Chapter 4 explains the communication between CEP and its environment. Chapter 5 specifies the deployment of the activities for each observation mode. In chapter 6, operational aspects are described. Chapter 7 contains a list of Open Ends.

## 1.5. Terms and abbreviations

| Term or abrreviation | Description |
|---|---|
| ACC | Application Configuration and (lifecycle) Control |
| AIPS++ | Astronomical Image Processing System |
| BBS | BlackBoard Selfcal system |
| BG/L | Blue Gene/L |
| CEP | CEntral Processor |
| COTS | Commercial off-the-shelf |
| DAS | Direct Attached Storage; storage hardware (e.g. RAID cabinet) 1-to-1 attached to individual systems (cluster nodes in our case) |
| DM | Dispersion Measures |
| FPGA | Field Programmable Gate Array |
| FSI | Flagger Selfcal Imaging |
| FOV | Field of View |
| GbE | 1 Gigabit Ethernet |
| 10 GbE | 10 Gigabit Ethernet |
| HBA | High Band Antenna |
| HECR | High Energy Cosmic Rays |
| KVT | Key-Value-Time (SAS metadata logging format) |
| LBA | Low Band Antenna |
| LOFAR | Low Frequency Array |
| MAC | Monitoring and Control system |
| MPI | Message Passing Interface |
| Node | individual server in a cluster computer<br>or node card in the Blue Gene machine |
| OLAP | On-Line Application Processing |
| PVSS | Prozessvisualisierungs- und Steuerungssystem; a SCADA System |
| RC | Rekencentrum (The computer centre (RC) is the centre for information technology offering high-grade IT-services to the State University of Groningen and other educational organisations) |
| RCU | Receiver Unit |
| RSP | Remote Station Processing |
| SAS | Specification, Administration and Scheduling |
| SCADA | Supervisory Control and Data Acquisition |
| SNMP | Simple Network Management Protocol |
| TBC | To be confirmed |
| TBD | To be defined |
| TAB | Tied Array Beam |
| UHEP | Ultra High Energy Particles |
| UML | Unified Modelling Language |
| USG | User Software Group |
| VHECR | Very High Energy Cosmic Rays |
| WAN | Wide Area Network |

# 2. Main drivers and status for the CEP architectural design

This section describes the main drivers for the CEP architecture design. The architectural design is guided by the CEP SRS [8], which is based on the LOFAR System Requirement Specification (SRS) [9]. However, the requirement specification documents are not maintained, and the ideas about LOFAR evolve with time.

Due to hard work of the User Software Group (USG), it is much clearer what is needed in CEP for the Key Science Projects, but the descriptions of the transient, pulsars and cosmic rays observation modes are still provisional. Input from the User Software Group has been extremely important for a complete architectural description of CEP.

The LOFAR software team has made many design decisions during initial implementation of the LOFAR CEP. These decisions are incorporated in this CEP architectural document or in de architectural description of one of the subsystems.

The LOFAR software team made some important assumptions and decisions that have guided the design effort:

- Use the Blue Gene/L for computationally intensive tasks
- Built to cost; minimise the cost of the system where possible. Use COTS components as much as possible. Use market standards wherever possible.
- Based on the current state of technology and expected rate of growth in processing power versus network bandwidth it was decided to design the system such that it is IO-bound. In other words, we provide enough processing power to all data that can be brought in and through the system.
- Off-line cluster complies with common accepted and installed standards at universities which provide portability between CEP and local resources at universities.

The following CEP requirements are fulfilled by the presented architecture (see section 2 of the CEP SRS [8]).

The goals realised by the current CEP architecture are:

- Provide enough processing and storage resources for the handling of observation data for the specified observation modes.
- Provide export capabilities to export data products (possibly raw or half-processed) to users
- Provide a rigorous framework for the development of high-performance processing pipelines. Use available hardware with high efficiency (~80% for basic processing steps)

The CEP architecture has the following main properties:

- High bandwidth and (very) low latency connections between all cluster nodes (but not between clusters)
- General purpose processing power in all parts of the system.
- Multiple observations can run simultaneously, partly sharing services or products from other applications
- Support for transient / quick switch operational modes
- Applications can be made tolerant to hardware failures using fail-over/hot swap hardware

Regarding the non-functional requirements of the CEP architecture:

- CEP is flexible in use and reconfigurable. All hardware resources can in principle be used for each observation mode.
- A scalable design allowing for both a staged installation and extension to larger sizes later on in the life of LOFAR. The current design can be extended at least a factor of 2 compared to the baseline configuration of LOFAR.

Requirements, wishes and features not met by the current design (yet):

- Inspection of data products in running applications (pipelines) is only implemented in a limited way.

- Required size of the off-line cluster is very much dependent on the Selfcal algorithm studies and implementation efficiency of Selfcal and imaging. All these are rather uncertain at the moment of writing. The current architecture does allow for quite some scaling. There may appear a need for specific hardware features such as large single-node systems or shared memory between cluster nodes.
- A framework for the development of off-line applications has to be designed and developed using parts of CEPFrame. This is addressed in the user software plan [10].

## 2.1. Top-level CEP specifications

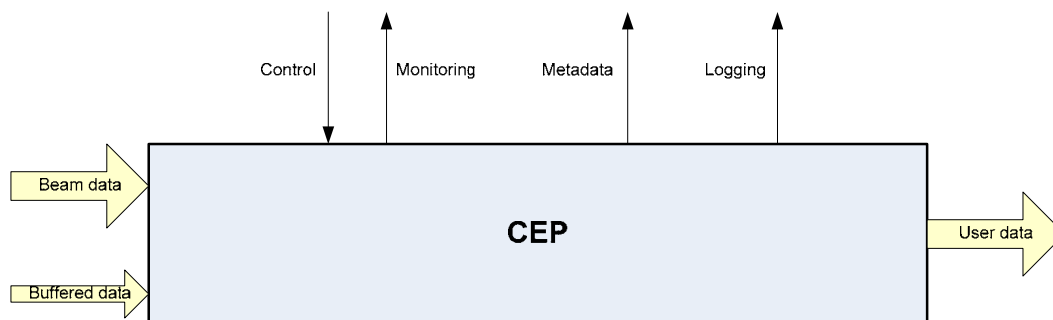| Example Applications performance | | value |
|---|---|---|
| Correlator capacity | Nominal | 77 stations, (1 FOV, 32 MHz)..( 4 FOV, 8 MHz) , 16 bit samples |
| | Peak | 32 stations, 24 FOV, 32 MHz, 4 bit samples |
| Beamformer capacity | Tied array 77 stations | 70 FOVs, 32 MHz |
| Hardware specs | | |
| Sustained data rates [Gbps] | Input section | 400 |
| | Blue Gene/L (in+out) | 400 |
| | Auxiliary section | 400 |
| | Storage in | 50 |
| | Storage out | 100 |
| | Off-line | >400 |
| Processing power (peak measured) [TFlops] | Input section | 0.4 |
| | Blue Gene/L | 27 |
| | Aux. processing | 0.2 + 15 (coprocessor) |
| | Off-line processing | 4 |
| Storage | | 1 PByte |
| internal network | in clusters | Infiniband or similar |
| | In Blue Gene/L | 3D torus and tree |
| | Between clusters and BG/L | GbE |
| External  IO | Via storage | 20 Gbps into internet |
| Software specs | | |
| Operating system | | Linux |
| Programming language | | C/C++, Java Swing |
| Compilers and IDEs | | GCC, Netbeans |
| Application programming environment | | Framework based; specialisation of Data formats and data transformations. |
| Run-time environment | | Run time definition of  user parameters, data transport, mapping onto hardware |
| Applied middleware | | MPI |
| databases | | Postgresql, AIPS++ tables |
| Scripting language | | Python |
| Modelling language/ tools | | UML, Rational (IBM) Rose, Enterprise Architect |

**Table 1** *top-level specification of the CEP facility*
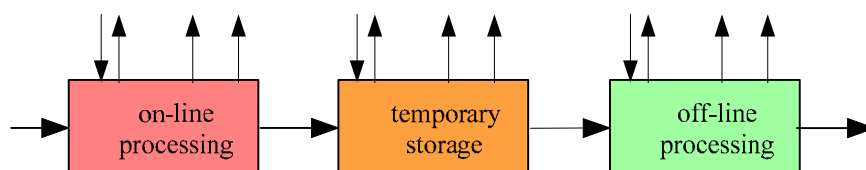
# 3. CEP and its subsystems

CEntral Processing has to deal with the following external communication (see Figure 3):

- **Incoming station beam data:** This is astronomical data which has been measured by the station antennas, and has been pre-processed / reduced at the stations. It has to be received and processed in real time by CEP. The maximum data rate is 400 Gbps.
- **Buffered data:** This is astronomical data which has been measured by the station antennas, and has been buffered in the transient buffer boards at the stations. This data are transported from the stations to CEP on (automatically generated) request. The maximum data rate is 400 Gbps.
- **The User data:** The end-products of the CEP processing are a set of standard LOFAR data products. The type of data product will depend on the specific observation mode. After creation these data products are exported to the archive and/or to the users. The User data may be processed further by software of the User Software Group [10].
- **Control:** ACC/MAC configures CEP for observations and starts / stops observations.
- **Monitoring:** During the lifetime of CEP, CEP sends out status information about the CEP hardware and the CEP processes to ACC/MAC. Among others, the monitoring information can be viewed with the Navigator application.
- **Logging:** During the lifetime of CEP, CEP sends logging to the logging system which is part of MAC.
- **Metadata:** During observations, CEP sends metadata to the metadata collection system in SAS (e.g. weather information)



**Figure 3:** *Input / output of CEP*

Because the astronomical data from the stations has to be received real-time, while the data leaves CEP on request, we chose the following division of CEP (see Figure 4):



**Figure 4:** *Division of CEP based on time / processing stage*

The online processing is a real-time subsystem. This component receives astronomical data from the stations and performs data reduction. After this stage, the data are temporarily stored. After the data of an observation is stored completely, the data are post-processed by the off-line processing system. What exactly happens in which stage, depends on the observation mode and is specified in chapter 5. Each of the stages has all aspects of the communication with ACC/MAC (control, monitoring, logging, metadata). The colours of the stages are reused in other pictures in this document. The combination of on-line processing and temporary storage is often mentioned as OLAP (On-line Application Processing).

The CEP system consists of a Blue Gene/L (BG/L) supercomputer embedded in a Linux cluster. This division is shown in Figure 5. The Linux Cluster internally communicates via an Infiniband or similar infrastructure (see Appendix A.2). The Linux cluster and the BG/L communicate via Gigabit Ethernet (GbE). The BG/L system has been incorporated in the design due to the high computational requirements of CEP.
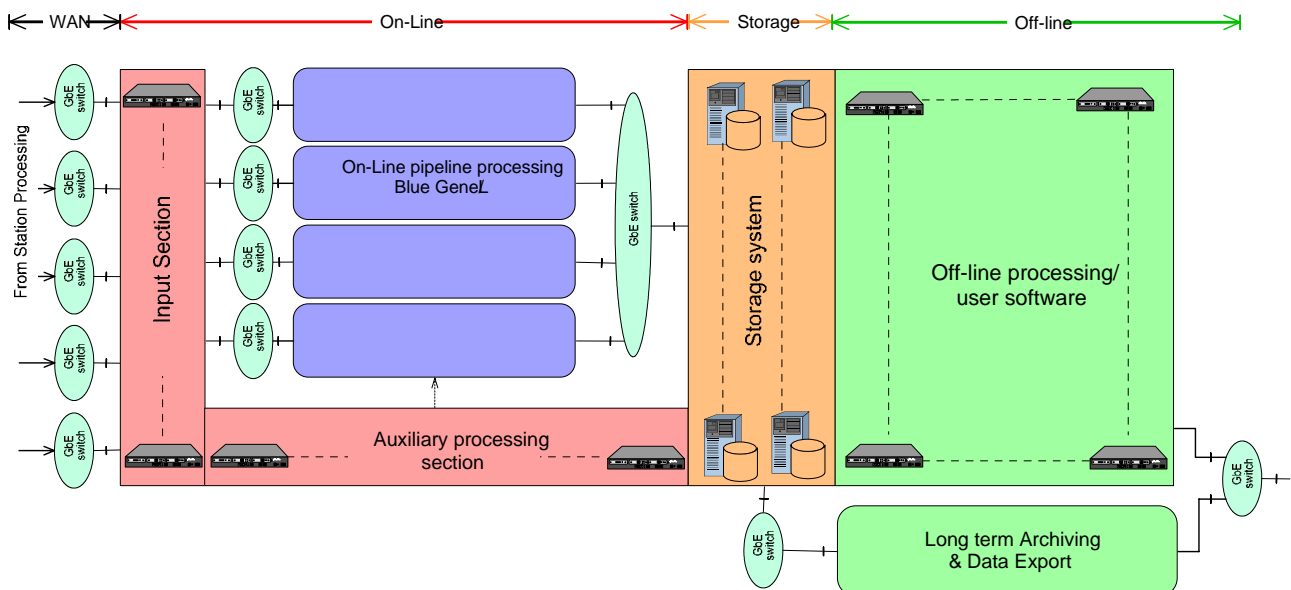


**Figure 5:** *Global overview of the CEP hardware*

The Linux Infiniband Cluster is divided in four sections (see Figure 6):

1. **The input section:** This cluster receives and processes real-time data streams from the stations. It utilises Infiniband interconnections for fast data transport.
2. **Auxilary processing section:** Parallel to the BG/L resources, this cluster is available for general-purpose processing and it is meant to perform tasks that the BG/L can't perform. The results of applications that run on this cluster may be used as control data for the processing applications running on the BG/L platform.
3. **The storage section:** This cluster collects the processed data streams and makes the resulting datasets available to the off-line section.
4. **The offline section:** This cluster post-processes the collected astronomical data, and makes it usable for the scientific research.

Four BG/L racks are connected to the input section through Gigabit Ethernet (GbE) (see Appendix A.1). Two additional BG/L racks are present, but only for EoR observations they are available for LOFAR operation. The BG/L 3D torus structure is available for internal communication between processors within the BG/L [15].



**Figure 6** *Central Processor Facility hardware architecture overview. Four Blue Gene racks are embedded in a hybrid cluster computer. The cluster is connected internally with Infiniband; the connections between the*

*cluster and the Blue Gene IO nodes is through Ethernet switches connected to a selection of the cluster nodes. The sizes of the various parts of the cluster are scalable.*

## 3.1. Input section

This part of the CEP facility contains the connections to the stations (through the WAN). Each input node will receive data from the stations and runs a data-handling application that will buffer the input data and synchronise its output stream with the other input nodes, based on the timestamps contained in the data.
In short the tasks executed on the input section are:

- Receive data streams from each station
- Validate data, replace missed data packages with dummy data, and flag faulty/replaced data
- Buffer the data stream, thus allowing asynchronous processing in the next applications
- Synchronise input streams based on time stamps
- Apply delay compensation (over complete sample periods)
- Re-pack and re-order data format
- Route data to correct output (BG/L and/or auxiliary cluster) connection for the next applications

The hardware available for these tasks can:

- Receive GbE data streams from the stations
- Provide 4 GB RAM to implement buffering (=several tens seconds)
- Provide processing power for data validation, injection etc.
- Route data via internal connections between nodes
- Connect to external clusters/systems for synchronisation and delay compensation control

The connections to the WAN are implemented in 1 GbE technology with 1-2 copper connections per node. Data output to the BG/L is over copper GbE as well. The nodes of the input section internally communicate through an Infiniband switch fabric. This interconnect system is used to route the data streams to the appropriate node for point-to-point connection with BG/L nodes.

CPU power is available for data recognition, simple validation, a single correction factor per input word and routing. Initial prototypes of the input section application indicate sufficient availability of processing power for the handling of each data package.

## 3.2. On-line pipeline processing (Blue Gene/L)

All data transformations with high computational requirements, needed before the storage are integrated in processing pipelines that are executed on the Blue Gene/L. The BG/L provides a large amount of processing power and internal bandwidth. Analysis has shown that all operational modes will be IO bound on this system. This was actually the goal of the architecture: build a system that has more than enough processing power. The internal bandwidth in BG/L is more as needed, which is good for algorithms that do require data exchange between parallel executing parts.

The BG/L system is basically a very large and integrated cluster of processing nodes. These nodes are part of a system-wide 3D torus network for internal communication. The on-line processing pipelines are parallelised in many parallel streams, where each such stream is mapped in a single "IO cell" with 1 GbE bandwidth and 8 associated dual-core compute nodes.
The on-line processing pipelines are mapped onto the available IO capacity of the BG/L in the environment of Linux clusters. A total of 512 GbE channels are available for the 4 racks that are available for nominal operation (up to 768 GbE channels for peak operation with all racks).

## 3.3. Auxiliary processing section

This system consists of a Linux cluster with connections to the input section and the storage sections, exactly like BG/L. However, these connections are high bandwidth Infiniband connections. Data transport

connections are available between the Auxiliary processing tasks and the on-line pipeline processing tasks on the BG/L through the GbE connections.

The auxiliary processing section provides processing power parallel to the BG/L. This additional processing power is needed for three reasons:

1. This processing power is available in a Linux run-time environment, complete with multi-threading etc. For some complex applications such as run-time features, it is not possible to port the needed software libraries to BG/L (e.g. AIPS++).
2. Application-to-application communication within BG/L is not possible. So shared applications providing tuning parameters to the on-line processing pipelines must always execute external from BG/L.
3. The IO bandwidth of BG/L may be too low for the most demanding observation modes (in particular the EOR). So additional resources can be installed.

Typical tasks executing on this hardware are:

- On-line ionosphere calibration; possibly providing correcting phase rotations to the correlator or beamformer
- RFI detection and mitigation; providing correction factors to the correlator
- Smaller on-line analysis applications using only a few nodes, for example the non-astronomical applications.

The availability of high-bandwidth connections to the input section and storage system makes the auxiliary processing cluster processing power limited. We can turn this part of the system back to IO-bound mode by providing more processing power for on-line processing pipelines executing on this system. Actually, this is the way to scale the processing power of the on-line section to higher values apart from the possibility to buying addition BG/L racks.

## 3.4. Storage section

The storage section provides disk space for the collection of data streams and storage of complete observation datasets for off-line processing. It has GbE connections to the BG/L system and Infiniband connections to the off-line processing section. The storage system is constructed as a cluster with an internal interconnects system. The storage is intended for temporary usage until the final data products are generated and archived, or the raw data itself is exported or archived. The system is intended to hold raw data for ~5 days on average and, in addition, data end-product for a few weeks.

The operational mode that is assumed in this architecture is that all off-line processing is scheduled soon after observation and raw data can be removed afterwards. End products are exported to other systems from this storage system, after which the data can be removed from this storage system.

Data from the on-line pipelines are sent to the storage system and are converted to the AIPS++ Measurement Set data format, but it is not yet decided whether this will be the final choice. Performance requirements may lead to a different choice for a LOFAR internal storage format. HDF5 is a candidate. The available processing power of the storage system might also be used for simple data transformations or sort actions.

It is not yet decided how the storage system cluster nodes access disks. Disks may be integrated in the storage systems nodes or directly connected to those machines (e.g. an external SCSI RAID set). The storage nodes may also have connections to a large SAN or NAS system. It is important that the chosen solution will facilitate simultaneous observation and off-line processing. This means that reading data for post-processing of an observation may not disturb the real-time writing activities of a running observation.

## 3.5. Off-line processing section

This section offers general-purpose processing power and high-bandwidth interconnections to the off-line processing applications. This cluster is used to perform processing on stored data. The largest part of this

cluster is a "normal" Linux cluster computer optimised on cost per flop. A shared-nothing architecture is used, possibly with dual CPU SMP nodes. Such a shared-nothing architecture consists of cluster nodes that only interact with each other by message passing over the network; they have no shared files space or memory space.

It has not yet been decided whether nodes with special capabilities will be added to the cluster for specific tasks. A few multi-processor nodes with shared memory might be added for special processing tasks such as the solve algorithm in Selfcal.

## 3.6. Data export section

This section comprises a Linux cluster with the following responsibilities:
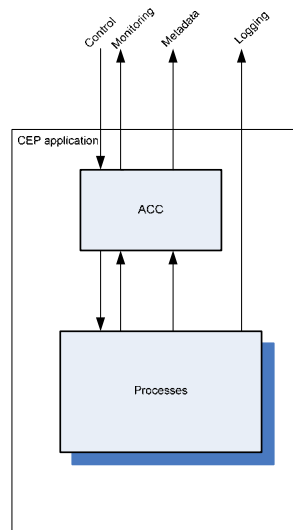
- **Storage of LOFAR astronomical end products:** These products are stored for maximal 3 months.
- **Exporting of LOFAR astronomical end products:** Once the end result of an observation has been created, the data should be exported to the data centres of the science organisations that requested the observation. The science organisations will permanently archive the end products or use the LOFAR end products for further processing (grid computing solutions like Big Grid). Export will take place using a network solution (no physical media like DVDs, etc.). A throughput of 20 Gb/s is foreseen for this goal. Exporting will be managed from SAS (via MAC).

# 4. Communication between CEP and its environment

The execution of applications on the CEP facility is always ordered by MAC, based on the schedule in the SAS system. The following communication takes place (see also Figure 7):

- **Control:** Comprises commands from MAC to a CEP application
- **Monitoring:** This exchanges information about results of the control commands
- **Metadata:** This ends up in the observation result tree in SAS
- **Logging:** This is displayed in the Navigator application



**Figure 7** *CEP internal handling of communication with CEP's environment. All interaction between the CEP processes and the environment of CEP happens via ACC, except the logging, which is directly sent to the logging processor in MAC.*

The ACC software is used to interface between MAC and the CEP processes of the CEP application. When an application of CEP (e.g. input section) has to perform a task, MAC starts an ACC process and provides task information to ACC. ACC then starts up the application to perform the task. This will usually comprise many processes on many nodes. The lifetime of ACC ends when all application processes have stopped.

## 4.1. Monitoring and control

The monitoring and control communication are mainly active at the start and the end of a task. Starting or ending a task happens in a sequence of stages. MAC dictates the transition between stages, and ACC realizes these transitions by sending the right commands to various the processes. ACC verifies whether the transition has been realized by collecting data from the processes. When all processes meet the required conditions, ACC reports this to MAC. For more detailed information see [13].

## 4.2. Logging and metadata

Each software product has the possibility to log information. This logging activity can be configured to log to a log-file, but it can also simultaneously be exported to the logging processor, which is a part of MAC. In MAC, the importance of the logging is determined, and a subset of the logging will be reported to the PVSS database, so that this logging will be available in the Navigator SCADA application (for more information see [11]).

During observation, metadata will be logged to SAS. This metadata consists of KVT triplets (=*Key, value, time)*. The processes that create metadata send their KVT triplets to ACC, after which ACC forwards the data via a central key-value manager to SAS (for more information see [12]).
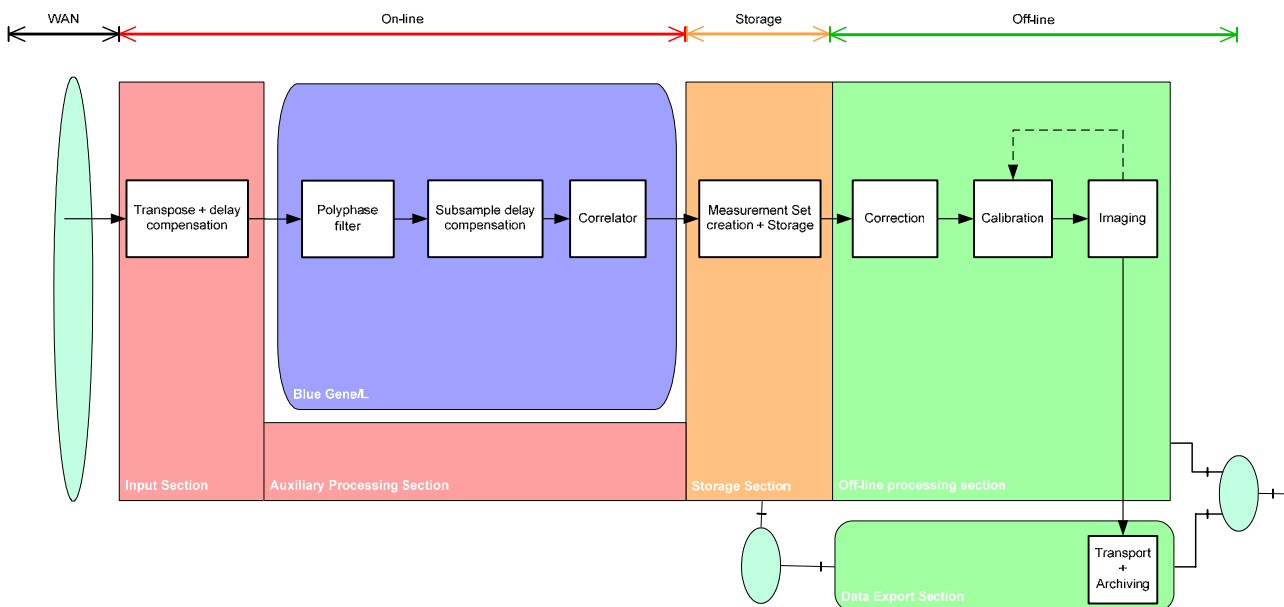
# 5. Observation modes

For LOFAR, a number of observation modes has been defined. For each mode, some processes with their own responsibilities will run somewhere in the CEP system. In this chapter, it will be clarified how these responsibilities for the different observation modes will be distributed on CEP.

## 5.1. Aperture synthesis modes

For aperture synthesis the following observation modes are needed:
- Standard imaging mode
- Surveys mode
- EoR mode

### 5.1.1. Standard imaging mode

**Figure 8** *Processing scheme for the Standard imaging mode*

In the standard imaging mode (see Figure 8), the input section of CEP receives in time changing spectral information from the stations via the WAN. The input section performs the following data transformations:
1. **Transpose:** This re-orders the data from 'subbands grouped per station' into 'stations grouped per subband'.
2. **Delay compensation:** The delay compensation only shifts the data over complete sample periods.

After the input section has performed these transformations, it sends the data to the BG/L. The BG/L transforms the data to uv-data (visibilities) by performing the following steps:
1. **Polyphase filtering:** This cuts the subbands in 256 channels to sub-kHz resolution
2. **Subsample delay compensation:** This performs the additional delay compensation on the channels over a fraction of the sample period by applying a phase rotation.
3. **Correlation:** Correlates the channel data.

The storage section receives the uv-data from the BG/L, and stores it.

The off-line processing comprises post-processing of the measured visibility data. This can only be scheduled to start after the observation has finished. The following steps can be performed:
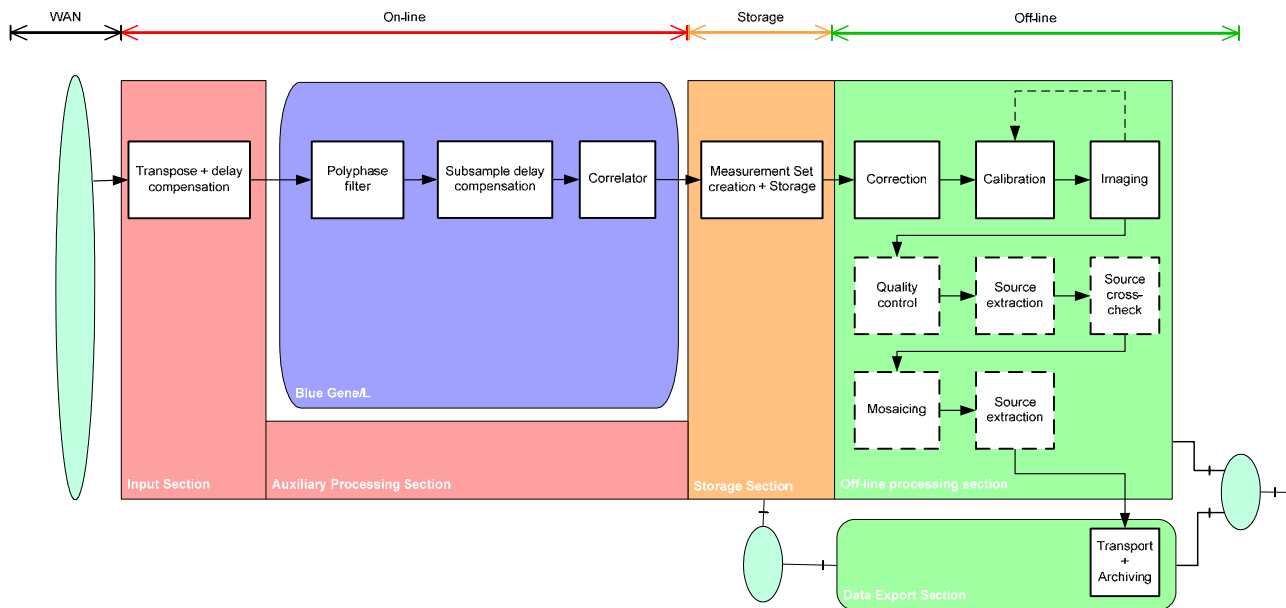1. **Correction:** Correction for known off-sets and initial flagging

2. **Calibration:** Compensating instrumental/environmental effects, subtracting known sources and additional flagging
3. **Imaging:** Converting the visibility data to image data, deconvolution, source extraction, update of the sky model

The user may perform these steps more than once (= major loop) to enhance the quality of the data.

The result of the off-line processing is temporarily stored in the data export section. From this place, the users can export the results to their own science centre for interpretation or additional processing of the astronomical data.

### 5.1.2. Surveys mode



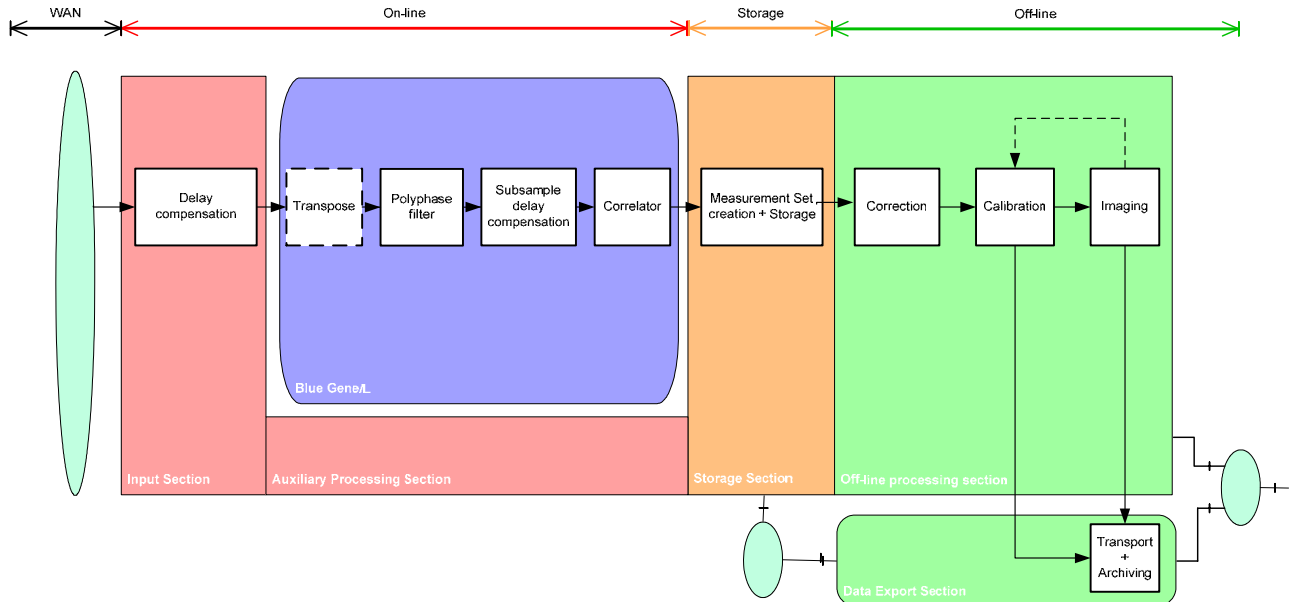**Figure 9:** *Processing scheme for the Survey mode*

The processing pipeline for the Surveys mode consists of the same steps as the standard LOFAR imaging pipeline with several additional post-processing steps. These offline processing steps are depicted in Figure 10 and consist of the following:

1. **Quality Control:** Performs a series of quality control checks on the output image and appends the results to the metadata associated with the image.
2. **Source Extraction:** Runs a source detection algorithm on the output image and creates a list of detected sources and their measured properties.
3. **Source Cross-check:** Compares resulting source list with the existing source database in order to check for consistency and appends the results to the image metadata.
4. **Mosaicing:** If the image passes all quality control and consistency checks, the image is combined with existing images to form a deeper mosaic. Associated metadata are also updated.
5. **Source Extraction:** The source detection algorithm is re-run on the combined mosaic to create a new master source list for the field. The resulting source list is then added to the master source database.

Upon completion of the pipeline, the resulting imaging data, mosaics, source lists and associated metadata will be transferred to the LOFAR archive and ultimately the Survey Key Project science centre. The individual pipeline components and associated derived data products are discussed in more detail in the Survey Key Science project plan [14]. The dashed boxes represent components which may run on either the off-line cluster or on external hardware, perhaps at the KSP science centres.

### 5.1.3.  EoR mode



**Figure 10** *Processing scheme for the EoR mode*

The EoR mode processing is basically the same as the standard imaging mode. However, EoR is both computationally and with respect to network bandwidth a much more demanding application, because 24 beams are observed simultaneously (with 32 stations, 32 MHz bandwidth).

It has not been decided yet whether the transpose will be performed on the Input Section or on the BG/L. At the moment we believe that the EoR requirements dictate the number of nodes of the Input Section, and moving the transpose to the BG/L gives a good opportunity to significantly limit the number of required input nodes, reducing the hardware cost [15].

The software in BG/L is mainly the same as for the standard imaging mode. Main differences are:

1. The stations will send 4-bit samples rather than 16-bit quantities, to reduce the amount of communication from the stations to the BG/L to something that can actually be handled.
2. To reduce output bandwidth, the integration time is 10 seconds instead of 1 second.
3. 24 beams in an observation instead of 1

For the scientific research, not only the image data are exported after off-line processing, but the visibility data are exported as well.
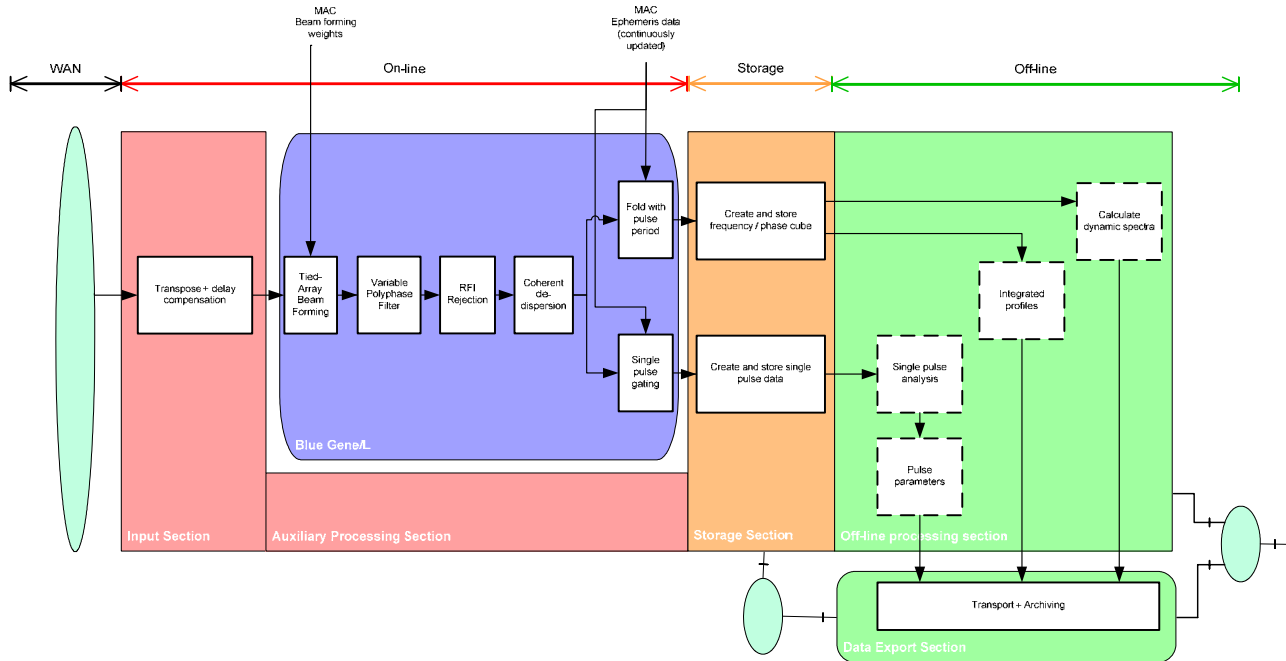
### 5.2.  Tied array modes

The following tied array modes are needed:
- Known pulsar mode
- Pulsar survey mode

### 5.2.1.  Known pulsar mode



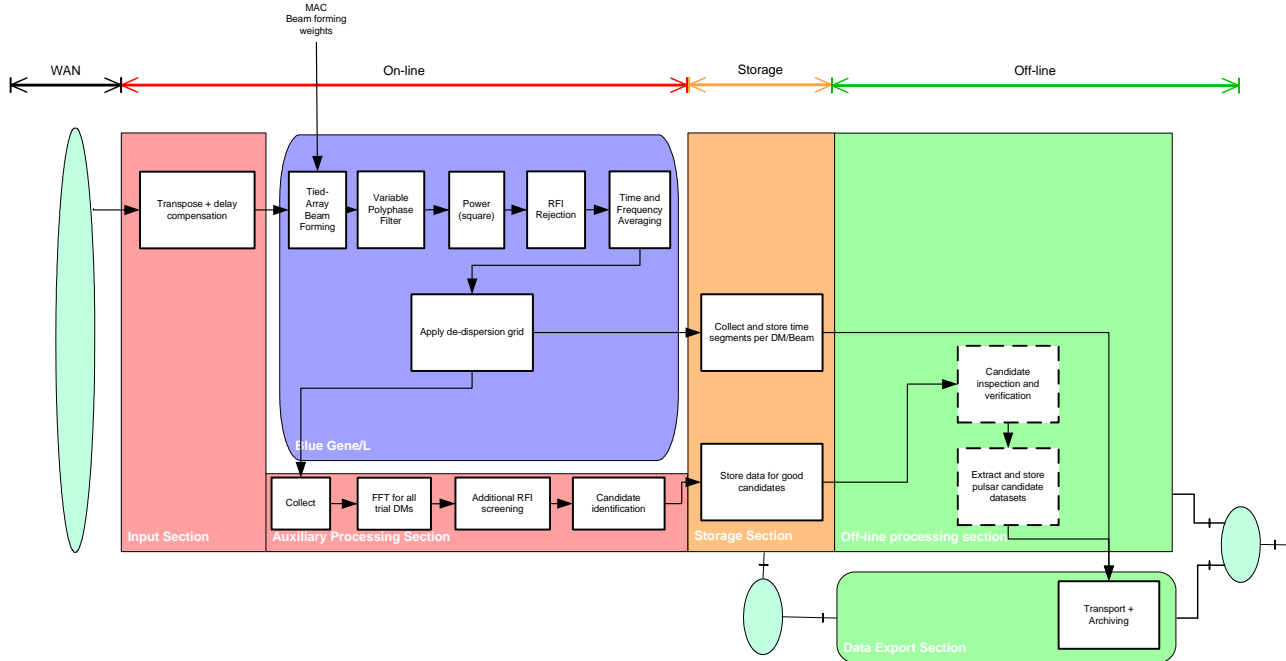**Figure 11** *Processing scheme for the Known pulsar mode*

For observations of known pulsars, data will be taken in a tied-array mode. In this mode, the signals from stations in the core are summed to form "tied-array beams" or TABs. In principle, many such TABs will be available and capable of tracking sources. This pipeline differs significantly from the standard imaging mode and consists of the following online components:

1. **Tied-Array Beamforming**:  Combines the station signals to form the tied-array beams. The necessary weights to construct the tied-array beams are supplied by MAC.
2. **Variable Poly-Phase Filter**:  Flexible filter for dividing the sub-band data into channels.
3. **RFI Rejection**: Screens incoming streaming data for potential RFI.
4. **Coherent De-dispersion**: Applies the known de-dispersion correction for the target pulsar.
5. **Period Fold**: Folds the data stream with the known pulsar period. This routine utilizes ephemeris data for the known pulsar, which is supplied by MAC.
6. **Single Pulse Gating**: Selects out the data for individual pulses. This routine also utilizes ephemeris data for the known pulsar, which is supplied by MAC.

Upon completion the data for the individual pulses as well as the phase/frequency cubes will be stored for subsequent off-line processing. The offline processing will produce various data products including dynamic spectra, integrated pulse profiles, and pulse parameters. These derived data products will be transferred to the LOFAR archive and ultimately the Key Project science centre for subsequent analysis. The dashed boxes represent components which may run on either the off-line cluster or on external hardware, perhaps at the KSP science centres.

## 5.2.2. Pulsar survey mode



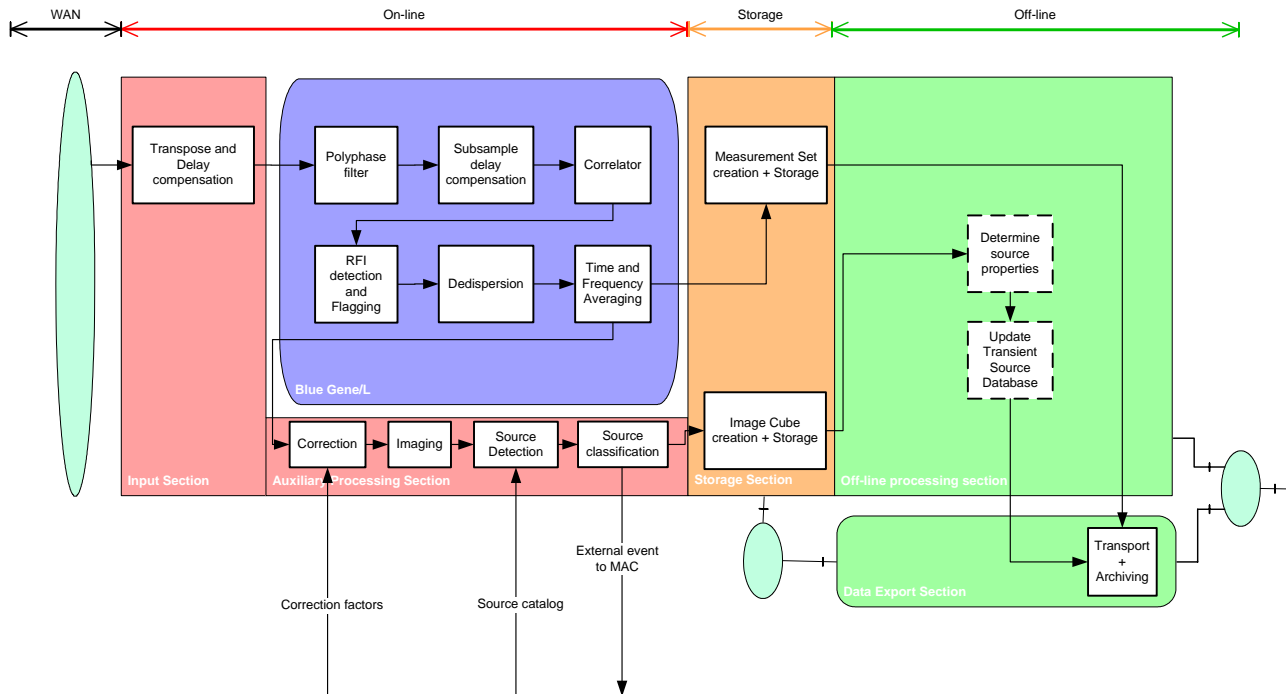**Figure 12** *Processing scheme for the Pulsar survey mode*

For survey observations to detect new pulsars, data will be taken in one of two tied-array modes. Signals from stations in the core are summed either coherently or incoherently to form "tied-array beams" or TABs. In principle, many such TABs will be available and capable of tracking sources. This pipeline has many similarities with the targeted pulsar observation mode discussed previously. The main difference is the necessity of correcting the data for a grid of possible dispersion measures (DM) in order to search for pulsar signatures. This pipeline differs significantly from the standard imaging mode and consists of the following online components:

1.  **Tied-Array Beam-forming**: Combines the station signals to form the tied-array beams (coherent or incoherent). This routine utilizes ephemeris data for the known pulsar, which is supplied by MAC.
2.  **Variable Poly-Phase Filter**: Flexible filter for dividing the sub-band data into channels.
3.  **Power:** The two polarizations are combined to form total intensity data.
4.  **RFI Rejection**: Screens incoming streaming data for potential bad spectral channels due to RFI.
5.  **Time/Frequency Averaging**: Averaging in time and/or frequency is applied as necessary.
6.  **De-dispersion Grid**: Applies de-disperison correction for a grid of possible dispersion measures.
7.  **Collect**: For each DM, a single timeseries is formed from all available bandwidth.
8.  **FFT**: Power spectra are formed for all DMs.
9.  **Candidate Identification**: The resulting power spectra are searched for significant peaks and candidate pulsars selected.

At this point, data for all candidate pulsars are stored for further offline processing. In this off-line stage, a series of diagnostic plots are created for each candidate pulsar and a final selection is made. Data for the final set of candidate pulsars are then extracted and exported to the LOFAR archive and ultimately the Key Project science centre for subsequent analysis. Pulsar candidates identified via the survey mode will then be scheduled for subsequent monitoring using the Known pulsar observation mode. The dashed boxes represent components which may run on either the off-line cluster or on external hardware, perhaps at the KSP science centres.

## 5.3. Transient detection mode



**Figure 13** *Processing scheme for the Transient detection mode*

In Transient detection mode, multiple beams from the LOFAR core stations are used to tile a large fraction of the available field of view and thereby monitor a significant fraction of the entire sky for variable low-frequency radio sources. This pipeline share similarities with the standard LOFAR imaging pipeline and Survey mode pipeline and uses some of the same software components. The primary distinction is that in order to respond rapidly once a new transient source is detected or a known source unexpectedly exhibits interesting new behaviour, the transient detection mode must operate in near real-time.

Up to the correlation stage, the transient detection pipeline mirrors the processing stages of the standard imaging pipeline. The additional online, post-correlation processing steps consist of the following:

1. **RFI Rejection**: Screens incoming streaming data for potential bad spectral channels due to RFI.
2. **De-dispersion**: Applies de-dispersion correction for a grid of possible dispersion measures (DM).
3. **Time/Frequency Averaging**: Averaging in time and/or frequency is applied as necessary.
4. **Correction**: Represents a quick calibration. Applies most recent base level calibration. No iteration is done.
5. **Imaging:** Visibility data are converted into image data.
6. **Source Detection**: Residual images are formed and a source detection algorithm is run to identify candidate sources. These candidates are compared to the existing transient and monitoring source database to identify potential new transients.
7. **Source Classification**: The new candidate sources are classified according to a limited set of criteria and, if appropriate, events are generated by MAC to initiate various responses. These responses might include reconfiguring LOFAR for more targeted observation of the source, dumping the contents of the transient buffer boards for detailed offline analysis, or even signalling other observatories (such as GLAST) to begin coordinated target-of-opportunity observations.

Once the real-time, online part of the pipeline has concluded, the residual images, compressed *uv* data, and source lists are stored. Further off-line processing will then measure various source properties and do a

more detailed source classification. These derived data products as well as a copy of the compressed UV data will shipped to the LOFAR archive and Transient KSP for later reference or possible further analysis. The dashed boxes represent components which may run on either the off-line cluster or on external hardware, perhaps at the KSP science centres.

## 5.4. Cosmic Ray detection modes

For the cosmic ray key science project the following observation modes are needed:

- **HECR:** High Energy Cosmic Rays
- **UHEP:** Ultra High Energy Particles
- **VHECR:** Very High Energy Cosmic Rays

### 5.4.1. HECR mode

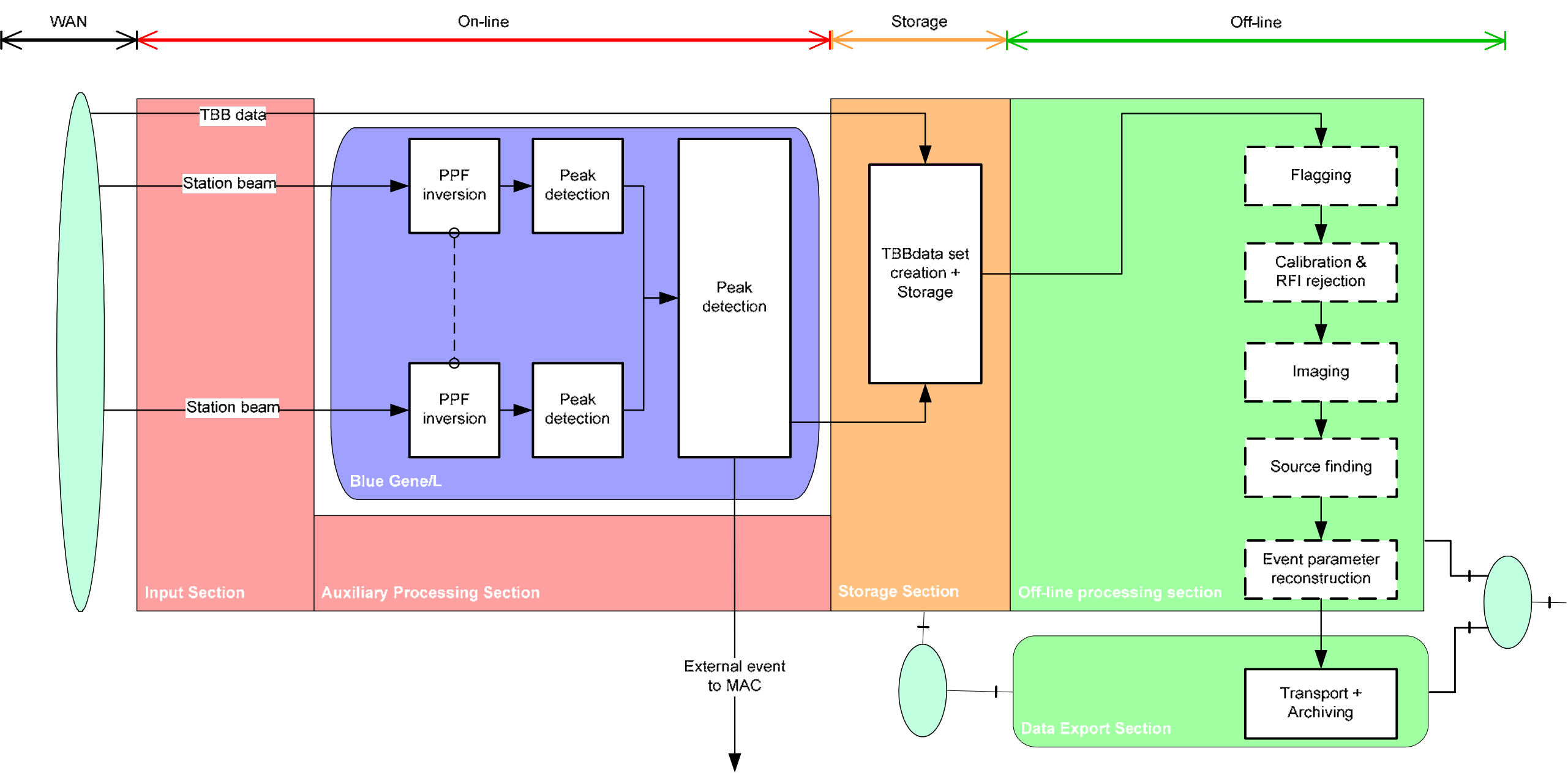


**Figure 14** *Processing scheme for the HECR mode*

Central processing for the observation of High-Energy Cosmic Rays (HECR) originally has been conceived to be run piggy-back on another observation using the LBAs (e.g. low frequency survey or transients observation). The dashed boxes represent components which may run on either the off-line cluster or on external hardware, perhaps at the KSP science centres.

The input section receives the station beamlets (a beamlet represents one subband from a beam), as they are utilized by the standard imaging mode.

The BlueGene/L will be responsible for retrieving beamformed time-series to be scanned by the peak detection:

- **PPF Inversion:** This inverts the effects of the poly-phase filter bank, which previously has split the frequency range into a number of channels. The aim of this processing step is to reconstruct a time-series from the station beams.
- **Peak detection:** This scans the previously generated time-series for intensity peaks; the first stage is working on the beams for the individual stations, whereas the second stage checks for possible coincidences between individually detected events (as some of them might produce a footprint covering multiple stations). If an event is identified, a freeze of the TBB data is initiated via MAC.
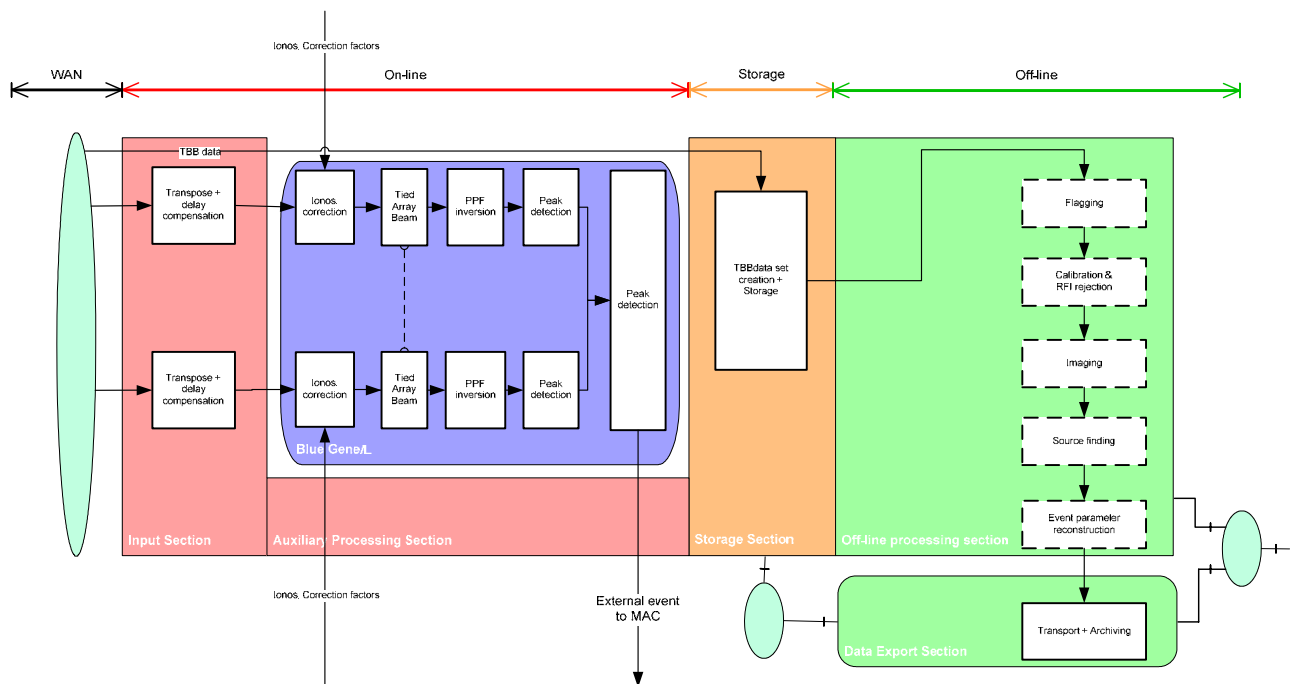
Once an event has been identified, the raw TBB data as well as the peak detection parameters are stored for in-depth analysis (everything so far essentially is triggering to decide whether or not to freeze TBB data).

- **Flagging:** This inspects the data of the individual dipoles for suspicious power levels, sudden jumps, etc.
- **Calibration & RFI Mitigation:** Calibration is taking place in the frequency domain. Here we first correct for the antenna gains, next for the antennas phases (using the signal of a narrow-band RFI source) and last remove the RFI from the data.
- **Imaging:** The cleaned and calibrated data are input to the imaging, which makes use of a number of different beamforming methods and supports near-field imaging.
- **Source finding:** The generated 5-dimensional image (3D position, time, frequency) is inspected to locate the cosmic ray air shower (EAS) signal.
- **Event parameter reconstruction:** Once the EAS signal has been located, the physical parameters of the event are reconstruced (e.g. peak height, distance from the shower core, etc.).

For the scientific research both the original TBB data as well as the event parameter lists are exported after off-line processing.

### 5.4.2. UHEP mode



**Figure 15** *Processing scheme for the UHEP mode*

The processing for the Ultra-High-Energy Particle (UHEP) mode is very similar to what is done in the HECR mode. The only significant change is the presence of two additional processing steps, which have to be performed on-line:
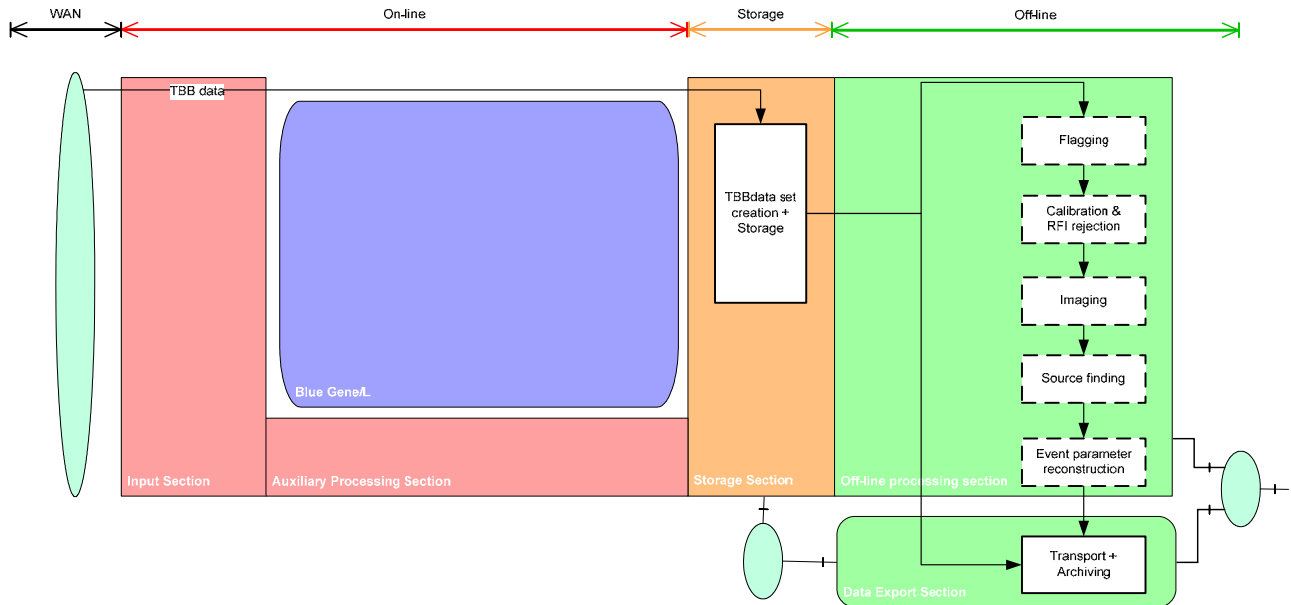
- **Ionospheric correction:** Apply solutions for the phase corruptions introduced by the Earth's atmosphere.
- **Tied-array beamforming:** Combine the corrected beamlets from a set of LOFAR stations into a tied-(sub)array beam – from this point on processing is as in the HECR mode.

The main constraint for this observation mode is that the complete processing chain – i.e. from data acquisition at dipole level, through CEP, up to dumping the contents of the TBBs – has to take place within the 1 sec time-interval for which the raw data are buffered on the TBBs.

The dashed boxes represent components which may run on either the off-line cluster or on external hardware, perhaps at the KSP science centres.

### 5.4.3. VHECR mode



**Figure 16** *Processing scheme for the VHECR mode*

For the very high-enery cosmic rays (VHECR) the radio signal produced by an extensive air shower is strong enough to be detected in the raw time-series of an individual dipole. For this reason the first two levels of signal detection is carried at station level (TBB and LCU) - the processing inside MAC is appended to search for events with a footprint covering more than a single LOFAR station.

The off-line processing is the same as for the HECR and UHEP modes. The dashed boxes represent components which may run on either the off-line cluster or on external hardware, perhaps at the KSP science centres.
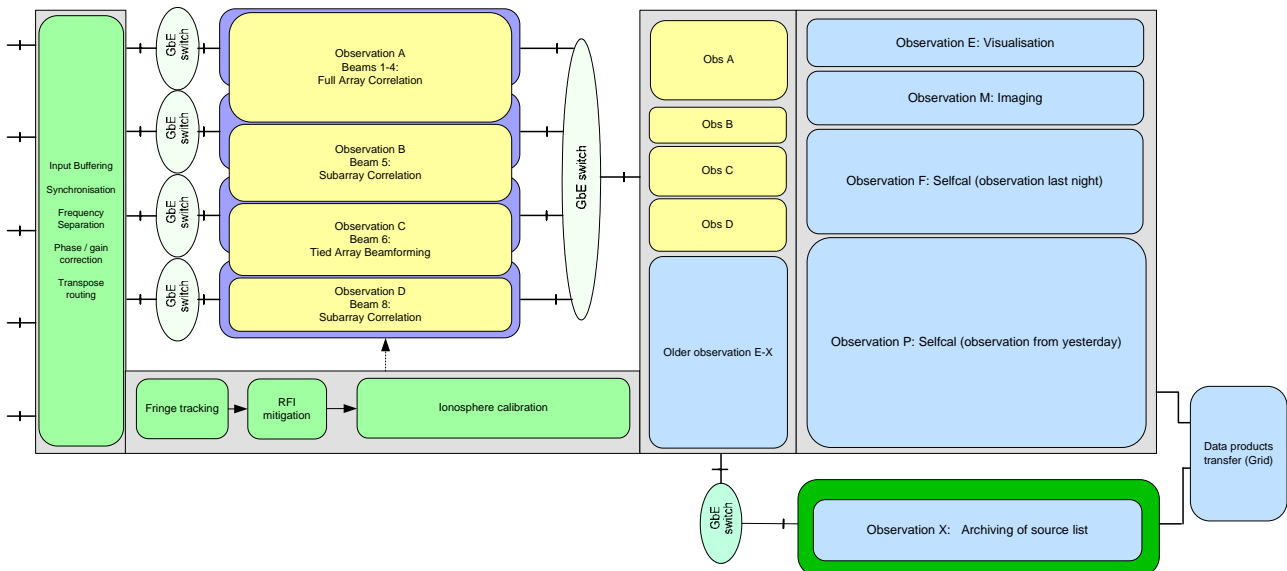
# 6. Operational aspects of CEP

## 6.1. Multiple observations

The previous chapter deployed the different responsibilities that had to be fulfilled, to the different subsystems. This paragraph describes how this will function with multiple observations running.

In Figure 17 we see an example with four concurrent on-line observations running that store their data through storage access tasks running in the storage subsystem. The input section and auxiliary processing section run continuously and are running supporting observations only. Multiple off-line applications are running concurrently on the off-line processing facility, accessing multiple datasets from past observations. Typically between 2 and 8 concurrent observations will be running, resulting in typically 15-20 concurrently running applications on the CEP resources.



**Figure 17** *Typical mapping of applications and data files on the hardware resources in the operational system. Multiple observations are processed concurrently. The green applications show basic processing services that have a lifetime longer than the running observations. The observation specific applications are shown in yellow. These applications will have different lifetimes reflecting the observation schedule. The light blue applications operate on data observed in the past.*

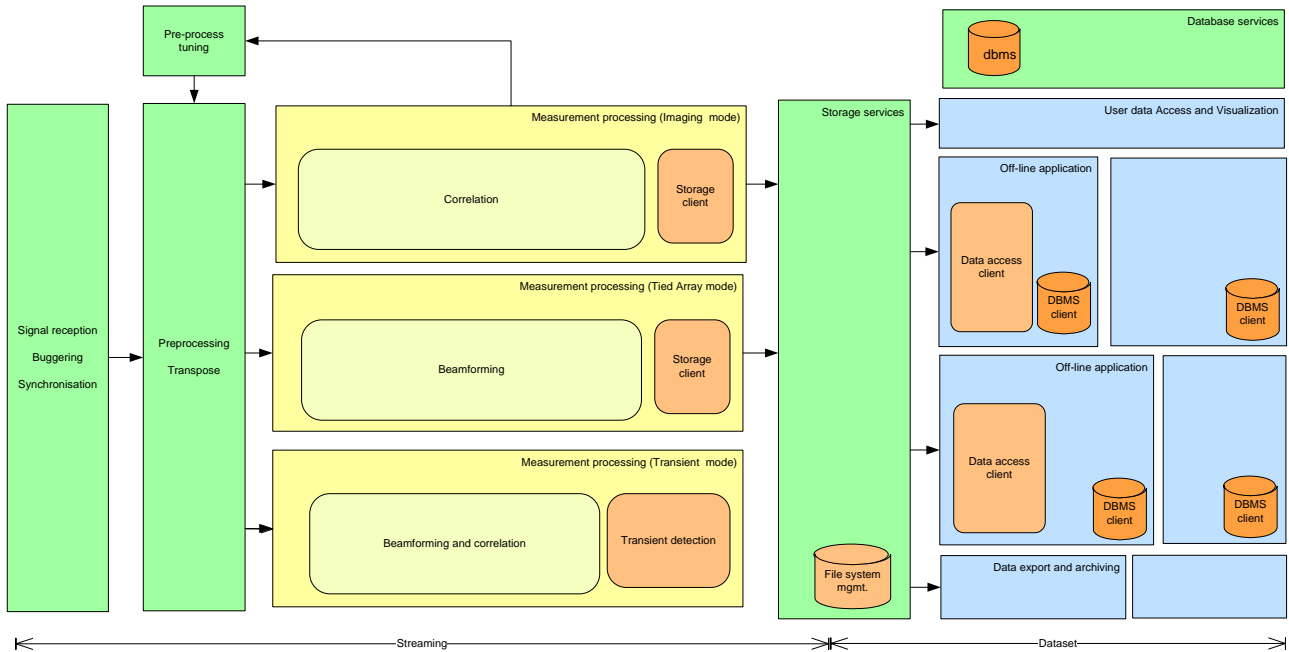## 6.2. Software Operational model

Based on the straight forward mapping of applications on the hardware discussed in the former section, we now look at the software only.

The schedule produced by SAS is based on the resource model which will prevent that applications are mapped onto incompatible resources. For instance, multi-threaded operation may be unavailable on some of the on-line platforms, or unwanted in view of real-time constraints. Moreover, the resource model will often forbid sharing resources over multiple applications in order to guarantee real time performance to on-line applications.

The overview of applications in the operational model, as shown in Figure 18, closely matches the hardware structure of the CEP. This is a result of the combined HW/SW design effort. All applications utilise multiple processing resources through parallel execution. In the figure we see a configuration during operation of

LOFAR with 15-20 applications running simultaneously. Typically between 2 and 8 observations will run concurrently; these correspond to one measurement processing pipeline application each.



**Figure 18** *Overview of applications in the CEP system during typical operation. The green applications run continuously, offering basis services to the remaining measurement or user specific applications. Multiple measurement applications run concurrently (shown in yellow, these are different processing pipelines and may have different lifetimes.*

# 7. Open ends

This section contains the items of the CEP architecture that are not clear yet.

| No. | Description | Solution |
|---|---|---|
| 1 | What will be the number of nodes of the input cluster? | |
| 2 | What will be the number of nodes of the auxiliary cluster? | |
| 3 | What will be the number of nodes of the storage cluster? | |
| 4 | What will be the number of nodes of the off-line cluster? | |
| 5 | How will the storage system cluster nodes access disks (internal to the node, external, NAS?) | |
| 6 | How do we make sure that the off-line cluster can read from storage during an observation? | |
| 7 | What will be the final format for uv-data on the storage? | |
| 8 | What will be the final format for TBB-data on the storage? | |
| 9 | Will the EoR transpose take place on the input section or on the BG/L? | |
| 10 | Are we definitely going to use Infiniband for the Linux clusters? | |
| 11 | Where will the source detection catalog live and how will it be populated | |
| 12 | Will the off-line section get identical nodes, or will there be some nodes available with special capabilities (e.g. multi-processor nodes with shared memory might be added for special processing tasks) | |

# References

[1]  N.R. Adiga *et al.*, *An Overview of the Blue Gene/L Supercomputer*, SC-2002, http://sc-2002.org/paperpdfs/pap.pap207.pdf

[2]  H.J. Pepping, *Correlator on a FPGA*, LOFAR-ASTRON-RPT-033, 2003

[3]  K. v.d. Schaaf, *CEP Requirements,Requirements Analysis,Architectural Design and Description,* LOFAR-ASTRON-MEM-35 v1.1, 2002

[4]  K. v.d. Schaaf, *IBM-ASTRON Blue Gene research collaboration,* LOFAR-ASTRON-PLN-024 v3.0

[5]  K. v.d. Schaaf, LOFAR central processor facility *hardware detailed design,* LOFAR-ASTRON-SDD-023 v0.1 (draft)

[6]  K. v.d. Schaaf, *CEPFrame detailed design,* LOFAR-ASTRON-SDD-024 V0.1 (draft)

[7]  K. v.d. Schaaf, *LOFAR Architerctural design document,* LOFAR-ASTRON-ADD-006 V4.0 (draft)

[8]  K. v.d. Schaaf, *LOFAR Central Processing Subsystem Requirement Specification,* LOFAR-ASTRON-SRS-007 Version draft1b

[9]  J. Kollen, *LOFAR System Requirements Specification,* LOFAR-ASTRON-SRS-007 v4.0

[10] M. Wise, *LOFAR User Software Overview,* LOFAR-ASTRON-USG-P1, Revision 0.5

[11] T. Müller, *Detailed design GCF*, LOFAR-ASTRON-SDD-007,

[12] R. Overeem, *Global Design SAS*, LOFAR-ASTRON-SDD-044,

[13] R. Overeem, *Monitor and Control ADD*, LOFAR-ASTRON-ADD-044,

[14] H.J.A. Rottgering, P.D. Barthel, M. van Haarlem, G.K. Miley, N. Mohan, R. Morganti, I. Snellen, *LOFAR Surveys of the Radio Sky*, May 2006

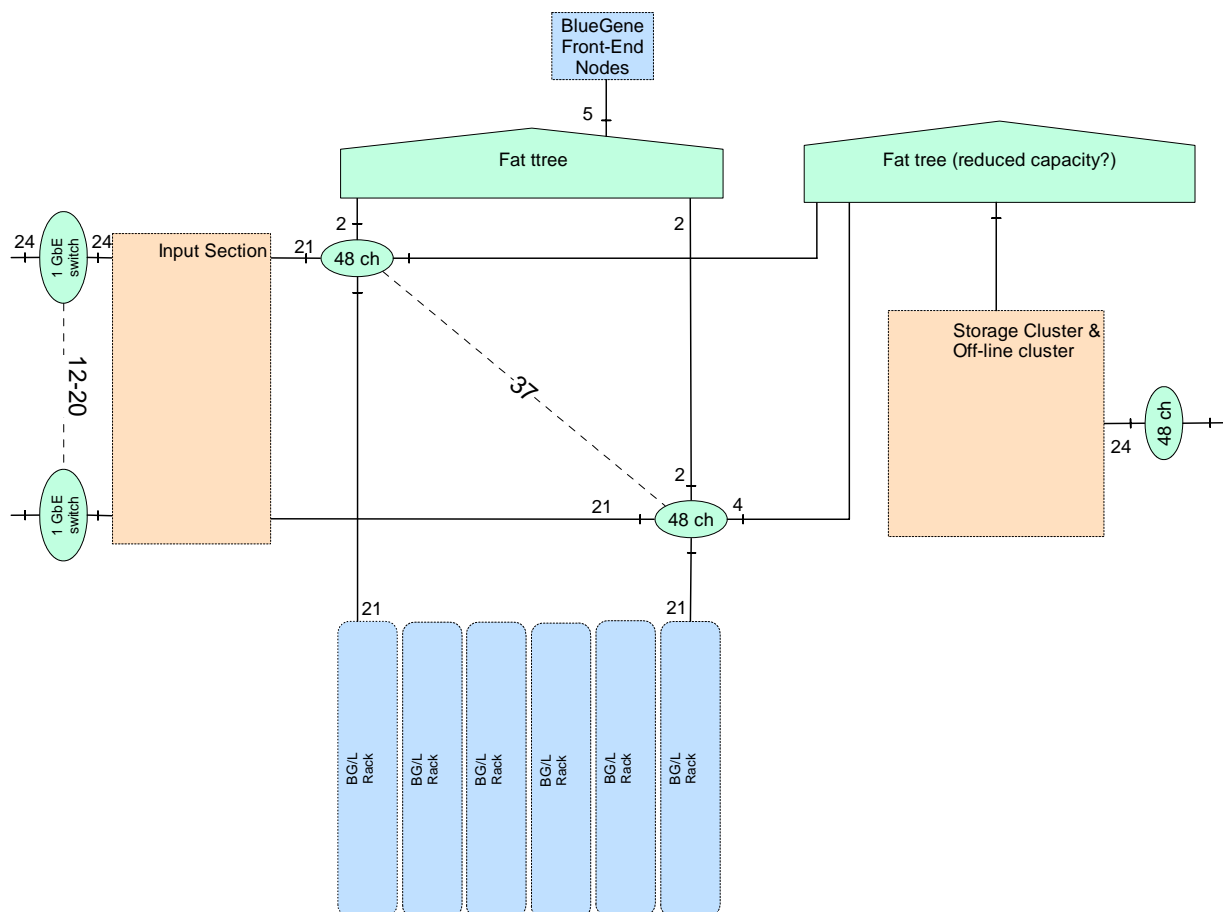[15] J. W. Romein, *OnLine Application Processing: Software Design Document*, LOFAR-ASTRON-SDD-036

# Appendix A    Network architectures

## A.1  Ethernet structure

The Ethernet network infrastructure is shown in Figure 19.  The main part of the Ethernet network is used for the connection between the input section and the BG/L IO nodes. This connection scheme is a-symmetric since the on-line processing pipelines reduce the volume of the data streams, in this topology the bandwidth into the BG/L nodes is 5 times as large as the output. A small all-to-all fat-tree network is attached to provide a small all-to-all bandwidth between the input nodes, and the front-end nodes to BG/L.

The BG/L output connections are connected to the storage section through a all-to-all tree structure. This connection tree is not completely symmetric (or "fat") but allows for the required sustained input data rate into the storage system.

The architecture shown here is based on small (cheap) Ethernet switches and limited all-to-all performance. This provides enough bandwidth and flexibility to fulfil all requirements on the network. In the detailed design we may choose to provide more all-to-all bandwidth by using for instance stackable switches.



**Figure 19** *Gigabit Ethernet network topology. Gigabit Ethernet connections are used to transport the  data streams from the stations into the input section. A specialised GbE network is used to transport data streams into the BG/L IO nodes and out to the storage section. Since the processing tasks on the BG/L platform are effectively reducing the data rates,  this network is asymmetric. Two tree networks are used to allow for some flexibility in the mapping of data streams to IO nodes.*

©ASTRON 2005-2007

## A.2 Infiniband networks

The Infiniband network topology is shown in Figure 20. Three switch fabrics are used for communication between the connected subsystems. Limited bandwidth is provided between the three switch fabrics for additional bandwidth along the typical processing line (from left to right in the figure). Note that a large part of the data transport in the on-line processing sections is provided through the GbE network connections between the input section and the BG/L system.



**Figure 20** *The Infiniband network topology consists of three connected switch fabrics. Each of the switch fabrics implements an all-to-all connection scheme with limited bandwidth compared to full symmetric "fat trees". The first switch fabric has a dedicated connection scheme optimised for the Transpose connection. The second one is used to transport data from the input section and auxiliary processing section into the storage system. Remember that the bulk data stream is transported into the storage system using the Ethernet network (see* Figure 19*). The third switch fabric is used for access to the stored data and for inter-process communication in the Off-line processing application. All node numbers and connections multiplicities are examples only.*

# Appendix B    Requirements coverage

This section shows which requirements are met by the current CEP architecture. In the requirements coverage tables, the contents of the 'Covered' field can have the following values:

| Value | Meaning |
|---|---|
| ✓ | The CEP design satisfies this requirement |
| ✗ | The CEP design does not satisfy this requirement |
| OExx | An Open End about this requirement can be found in Chapter 7. The Open End number is xx. |
| ? | This requirement is unclear |
| - | This requirement doesn't exist anymore |
| SAS | The responsibility for this requirement is not in CEP but in SAS. (Similar entries for the following work packages: FSI, MAC, SHM, WAN, Station) |

## B.1  Data handling requirements

This section shows the coverage of the LOFAR requirements that can be found in section 3.09.5 of the LOFAR System Requirements Specification [9].

| LO-3.09.5 DETAILED DATA HANDLING REQUIREMENTS | | |
|---|---|---|
| **Req. ID** | **Description** | **Covered** |
| -01 | LOFAR shall provide sufficient processing and storage functions to handle the data flow from the Remote Stations and the Virtual Core and to transform it in the specified data products. | ✓ |
| -02 | LOFAR shall provide storage capacity to keep intermediate data products of normal synthesis observations (1 sec and 1 kHz resolution, 32 MHz processed bandwidth and 77 stations) available for a limited period of time (at least 7 days = approx. 1 PetaByte). | ✓ |
| -03 | LOFAR shall be able to buffer data from Remote Station in order to correct for the arrival time difference due to the geographical positions of the Remote Stations and the Virtual Core substations, and to any delays introduced by data processing and transport equipment. | ✓ |
| -04 | LOFAR shall provide for random read access to (intermediate) data products stored for intermediate periods. | ✓ |
| -05 | LOFAR shall provide for processing capacity for analysis of the data during processing. | ✓ |
| -06 | LOFAR shall be capable of concurrent read and write access to (intermediate) data products stored for intermediate periods. | OE06 |
| -07 | LOFAR shall provide for processing capacity for data (flow) processing tasks. | ✓ |
| -08 | LOFAR shall contain processing capacity for the generation of final data products. | ✓ |
| -09 | LOFAR shall contain storage capacity for temporal storage of final data products. | ✓ |
| -10 | LOFAR shall be able to export final data products. | ✓ |
| -11 | LOFAR shall be able to export final data products via remote access. | ? |
| -12 | LOFAR shall provide for a resource model (as input for the monitoring and control process) including on-line information on:<br>• storage capacity<br>• input bandwidth<br>• data routing bandwidth<br>• processing capacity<br>• analysis capacity | SAS |

| -13 | LOFAR shall provide processing capacity for monitoring and control processes. | MAC |
|---|---|---|
| -14 | LOFAR shall provide for access to (intermediate) data products. | ✓ |
| -15 | LOFAR shall make the (intermediate) data products available for inspection by the user (scientist). | ✓ |
| -16 | LOFAR shall provide processing capacity for data access and visualization tasks. | ✓ |
| -17 | Deleted | - |
| -18 | It shall be possible to perform pulsar observation analysis applications. | ✓ |
| -19 | It shall be possible to analyze processed data after completion of the observation. | ✓ |
| -20 | It shall be possible to process incoming data streams in real-time. | ✓ |
| -21 | It shall be possible to start analysis of processed data during the observation. | ✗ |
| -22 | LOFAR shall be able to synchronize incoming data streams. | ✓ |
| -23 | It shall be possible to perform the processing tasks as listed in the previous requirements simultaneously. | ✗ |
| -24 | LOFAR shall have an on-line storage of system models. | ? |

Remarks:
- **Req. 11:** The content of the Data export section will be managed from SAS, and SAS can be executed over the internet. It is not clear if this fulfills this requirement.
- **Req. 21:** According to the current architecture, it is not guaranteed that it is possible to read from the storage cluster during observation (see OE06). At the moment, off-line processing can start after the online observation has completely finished. The current architecture does not support analysis of processed data of an observation that is currently on-line.
- **Req. 23:** This requirement is not met because some of the previous requirements are not met.

## B.2 Calibration and Post-Acquisition Processing requirements

This section shows the coverage of the LOFAR requirements that can be found in section 3.10 of the LOFAR System Requirements Specification [9].

| LO-3.10 CALIBRATION AND POST-ACQUISITION PROCESSING | | |
|---|---|---|
| Req. ID | Description | Covered |
| -01 | LOFAR shall be capable to perform the on-line or off-line calibration and processing required to generate data products as specified in section 3.9.2. | FSI |
| -02 | LOFAR shall be able to reduce the volume of the (raw) observation data by means of calibration such that it is compatible with the available processing capacity. | ✓ |
| -03 | For imaging applications the LOFAR calibration function shall be capable of correcting phase fluctuations to the accuracy of < 0.1 rad based on available bright calibration sources (Cat. I, as defined in RD. 4). | FSI |
| -04 | LOFAR shall be designed to estimate the instrumental responses in the direction of the Cat I sources continuously in all observing modes. | ? |
| -05 | LOFAR shall be capable of performing on-line selfcalibration in order to extend the integration time of acquired data to get sufficient SNR per visibility for further processing (this integration time is specified in 3.9.1-01 as 1 sec). | ✓ |
| -06 | LOFAR shall be designed to characterize the ionosphere for each beam direction in each remote station, both in synthesis imaging and in full phased array mode. | FSI |
| -07 | Inspection of intermediate calibration results shall be available to the investigator or science center for whom the observations were performed. | FSI |
| -08 | Each LOFAR station shall be able to acquire at least the following environmental data at 10 sec (TBC) intervals:<br>a) Air temperature at several places in the antenna field and inside the housing of station | Station |

| | electronics, with 1 K accuracy | |
|---|---|---|
| | b) Relative humidity inside and outside the housing of station electronics, with 10 % accuracy | |
| | c) Wind velocity and direction, with 1m/s and 5 deg accuracy | |
| | d) Lightning information, using general purpose lightning detectors | |
| | e) Rain and snow conditions, including wetness of the ground | |

Remarks:
- **Req. 4:** It is unclear if this requirement is still actual. The instrument is able to sacrifice beams for this purpose, but this is probably not relevant in all observation modes.
- **Req. 8:** This requirement contains a TBC. After confirmation of this requirement, we will check if this requirement is met.

## B.3  Archiving requirements

This section shows the coverage of the LOFAR requirements that can be found in section 3.11.1 of the LOFAR System Requirements Specification [9].

| LO-3.11.1 ARCHIVING | | |
|---|---|---|
| **Req. ID** | **Description** | **Covered** |
| -01 | LOFAR shall provide both short-term and long-term storage of intermediate and final dataproducts. | ✓ |
| -02 | Intermediate dataproducts that can be further processed based on user input shall be kept available on short-term storage for a period of one week (TBD: the selection of dataproducts to which this requirement applies needs to be defined more precisely). | ? |
| -03 | Only final dataproducts shall be kept on long-term storage. | ✓ |
| -04 | LOFAR shall provide a quick-look function to allow users to assess the quality of intermediate dataproducts over general communication networks. | MAC |
| -05 | LOFAR shall provide a catalog of all data products. | SAS |
| -06 | The data product catalogue shall contain information on the quality and usability of the observation's results. | SAS |
| -07 | Deleted | - |
| -08 | Deleted | - |
| -09 | The archive shall keep data for the lifetime (TBC) of the instrument. | ? |
| -10 | The archive shall be set up such that the chance of data loss is less than TBD | ? |

Remarks:
- **Req. 2:** This requirement contains a TBD. However, the short-term storage after on-line processing is able to store the data on average 5 days (see section 3.4)
- **Req. 9:** This requirement contains a TBC. Furthermore, the archive itself is not part of the CEP architecture.
- **Req. 10:** This requirement contains a TBD. Furthermore, the archive itself is not part of the CEP architecture.

## B.4  Export requirements

This section shows the coverage of the LOFAR requirements that can be found in section 3.11.2 of the LOFAR System Requirements Specification [9].

| LO-3.11.2 EXPORT |
|---|

| Req. ID | Description | Covered |
|---|---|---|
| -01 | It shall be possible to deliver final data products to scientists on digital storage media. | ✘ |
| -02 | It shall be possible to deliver final data products to scientists over computer networks. | ✓ |
| -03 | It shall be possible to retrieve observation data from the archive by searching on parameters in the product metadata. | SAS |

Remarks:

- **Req. 1:** This requirement is obsolete. Because of the large amount of data that LOFAR-CEP produces, it is not possible to export the data on storage media like DVD's / portable hard disks.