

# A Catalog of Cosmic Ray Detections in Gaia CCDs

Prepared by : Christian Kirsch (and Asier Abreu) (ESAC)

March 9, 2018

Document Version : v1.0

## Contents

<b>1</b>	<b>Objective</b>	<b>3</b>
<b>2</b>	<b>Background</b>	<b>3</b>
2.1	Intro . . . . .	3
2.2	Acronyms . . . . .	3
2.3	Conventions . . . . .	4
2.3.1	Definitions . . . . .	4
2.3.2	Reference System . . . . .	4
2.4	Data Source . . . . .	4
<b>3</b>	<b>Algorithm Description</b>	<b>5</b>
3.1	Feature Extraction . . . . .	6
3.1.1	SM-SIF: Laplacian Edge Detection . . . . .	8
3.1.2	BAM-OBS: Boxcar filtering / Stacking . . . . .	11
3.1.3	BAM-SIF: Thresholding . . . . .	15
3.1.4	Extraction Pipelines . . . . .	17
3.2	Post-Processing . . . . .	19
3.2.1	Observed Flux . . . . .	19
3.2.2	Track Geometries . . . . .	23
3.2.3	Edge Tracks . . . . .	28
3.2.4	Post-Processing Pipeline . . . . .	29
<b>4</b>	<b>Data Model</b>	<b>31</b>
4.1	Main Dataset Description . . . . .	31
4.2	Flux Dataset Description . . . . .	32
4.3	Auxillary Datasets Description . . . . .	33

4.3.1	Spacecraft Spin Phase . . . . .	33
4.3.2	Spacecraft Location . . . . .	33
<b>5</b>	<b>Persistence Considerations</b>	<b>34</b>
5.1	Persistence Method(s) . . . . .	34
5.2	Size Estimations . . . . .	34

## History of versions

Version	Date	Author(s)	Changes
0.1	29.10.2017	A. Abreu	First Data Model. Initial Definition
1.0	09.03.2018	C. Kirsch	First full draft, including algorithms and Data Model.

# 1 Objective

This document describes in detail a relational data model for a catalog of cosmic rays as detected in the Gaia CCDs.

## 2 Background

### 2.1 Intro

The Gaia cosmic rays catalog is the main output product of particle event detection processing performed at ESAC. This off-line processing is aimed to generate a catalogue of cosmic rays that allows to analyze the long term behavior of the radiation environment at L2 as seen by Gaia.

### 2.2 Acronyms

AC	Across Scan
AL	Along Scan
AF	Astrometric Field
SM	Sky Mapper
BAM	Basic-Angle Monitoring device
CCD	Charge-Coupled Device
TDI	Time Delay Integration
SIF	Service Interface Function
FITS	Flexible Image Transport System
FoV	Field of View
FPA	Focal-Plane Assembly
OBMT	On-Board Mission Time
UTC	Coordinated Universal Time
°	degree; unit of angle
'	arcminute, minute of arc; unit of angle
''	arcsecond, second of arc; unit of angle

## 2.3 Conventions

The following set of definitions and conventions are simply a subset of those (thoroughly) defined in [1]. We have adopted those referring to CCD characteristics and relevant for the present scope.

### 2.3.1 Definitions

- A ‘pixel’ is the elementary charge generation and storage element in the light-sensitive area of the CCD.
- A ‘column’ is the set of all pixels having the same across-scan coordinate.
- A ‘line’ is the set of all pixels having the same along-scan coordinate.
- The ‘summing register’ is a special pixel line following the light-sensitive pixels. It is used to combine (add or bin) the charges from several lines into one.
- The ‘read-out register’ is the special pixel line which is used to transfer the accumulated charges from the CCD chip into the read-out amplifier (and thus further into the further amplification and digitisation electronics).
- The ‘read-out amplifier’ is the electronic circuit at the end of the read-out register where the registration and pre-amplification of the photoelectric charges takes place.
- The ‘time delay integration’ mode (TDI mode) consists of gradually shifting the photoelectric charges from “left” to “right” — and eventually into the read-out register — to follow moving optical images over the light-sensitive area.
- ‘Gates’ are special lines within the light-sensitive area. If activated they act like summing registers, holding up charges, preventing them from moving along scan in spite of the TDI clocking. This causes a compression (summing) of the already accumulated TDI images into a single line, and the creation of a new “blank” empty space “in front of” this line.

Figure 1 summarizes these definitions.

### 2.3.2 Reference System

Figure 2 shows the Gaia Focal Plan Array composed of 106 CCDs and the different CCD variants (BAM/AF).

## 2.4 Data Source

For the analysis of the cosmic rays in Gaia CCDs we have used a number of different *engineering* datasets. Namely:

- Sky Mapper images downlinked via the Service Interface Function - from here on shortened to SM-SIF. These images are recorded fully in TDI mode and contain objects of interest such as planets, comets, dense stellar fields such as globular clusters or the Galactic center as well as, most frequently, very bright stars that drive the CCDs into saturation.

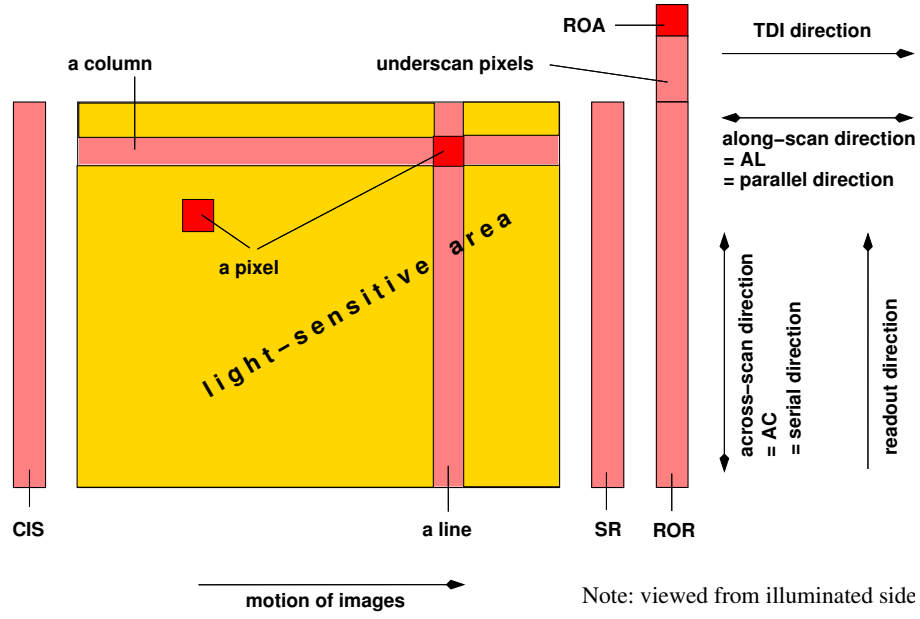


Figure 1: CCD-related notations, taken from [1].

- Nominal BAM observations, from here on shortened to BAM-OBS. These images consist of two readout windows roughly  $1 \times 1$  cm in size located on one chip, which continuously observe the interference pattern of the BAM in staring mode.
- SIF BAM observations, shortened to BAM-SIF. These images record the same AC region as the BAM-OBS source images, but span the width of the entire chip in AL. They thus contain off-pattern regions with a very weak background. This dataset has been previously analyzed for cosemics in [2].

### 3 Algorithm Description

A number of different event detection algorithms have been used for the generation of the catalog. Each algorithm is tailored to the Gaia detector type. Where possible, existing algorithms for cosmic ray detection in CCDs have been used. The following subsections briefly describe the types of algorithms used.

The extraction of the data product has been split into two steps, a feature extraction and a post-processing. The first step uses the input of the different data sources to produce for every observation a list of objects that have been recognized as cosmic ray tracks. The second step then analyzes the extracted tracks further, adding information such as their orientation on the FPA or their length, as well as the flux measured in each observation.

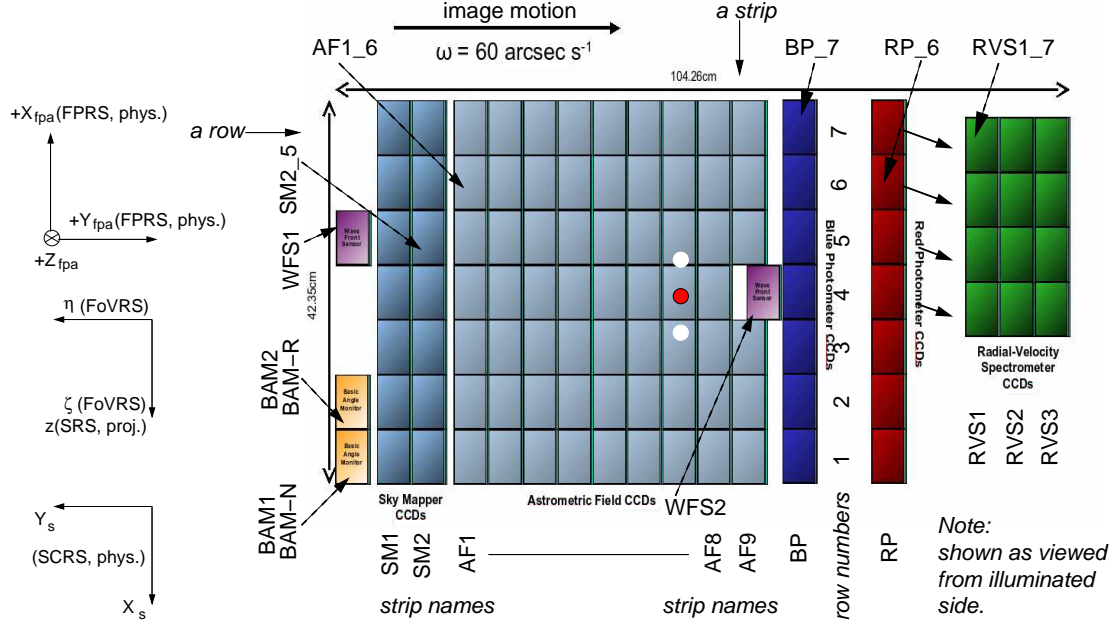


Figure 2: The Focal-Plane Reference System, taken from [1].

### 3.1 Feature Extraction

The goal of the feature extraction step is to collect all the objects recognized as cosmic ray tracks from the different datasets into a source-independent dataset for the post-processing. To that extent, we have defined an intermediary dataset named TrackObs (short for track observation), which collects necessary metadata on the source observation and contains a list of all identified objects.

Each TrackObs contains a list of metadata on the source observation, including:

- Source observation (SM-SIF, BAM-OBS, BAM-SIF)
- Acquisition time in OBMT
- Source image dimension (in AL and AC)
- Number of pixels not used for extracting cosemics (for the calculation of the exposed area)
- Sample binning, in AL and AC
- CCD gain
- Chip location on the FPA (Row and FOV)

The output of all of the following algorithms are three two-dimensional arrays of the same size as the input image: A boolean mask array set to True for all pixels identified as cosemics, a

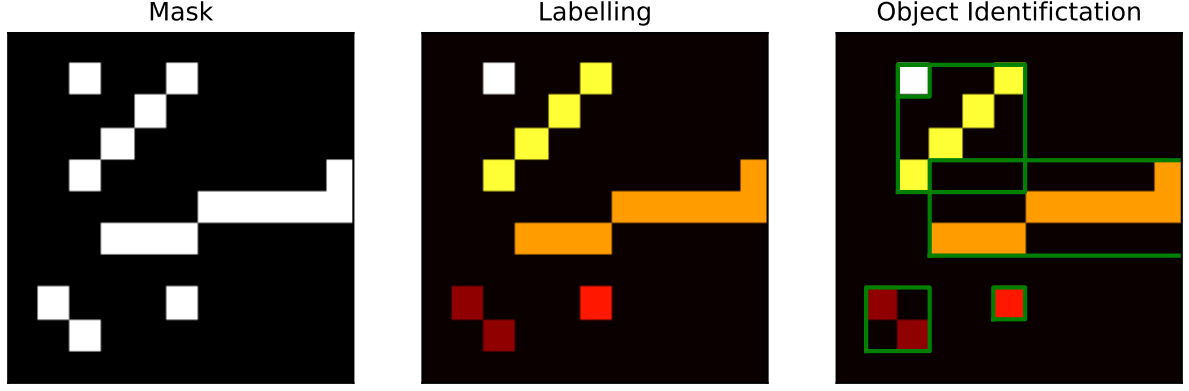


Figure 3: Identification of individual connected objects from a cosmic ray mask. The boolean mask (left panel) is used to calculate an integer array for all 8-connected pixels, which are labeled with a unique integer per object (middle panel). Objects are then identified by their labels as rectangular image sections around all pixels with a common label (green boxes in right panel).

signal array consisting of the original image minus the non-cosmic 'background', and an array of the uncertainty of this signal array. As the conversion of these arrays into cosmic ray tracks is the same in all of the following algorithms, we will explain the process here.

First, individual objects are identified from the mask. This process is shown in Figure 3. From the boolean mask, a labelling algorithm from the `scipy` Python distribution is used to assign the same integer value to all neighbouring pixels set to True – in this case, if they are horizontally, vertically or diagonally adjacent. The labelled array is then used by a `scipy` object finding algorithm to identify the bounding boxes of all pixels of the same label.

Each `TrackObs` then contains an array whose rows represent information on one of the objects identified in the above steps. For each track, it lists:

- The AL and AC lengths of the bounding box
- The location of the track, encoded as the lowest AL and AC value in its bounding box
- The sub-image of the track from the signal array, encoded in ADU. Note that this only includes those pixels inside the bounding box with the corresponding label - in the right panel of Fig. 3, the image of the  $4 \times 4$  diagonal track in the upper left does not include the single pixel track in its bounding box, which has a different label.
- The sum of signal values of the track in electrons, corresponding to the total measured ionizing energy.
- The uncertainty of the total ionizing energy, being the square root of the sum of the squares of the individual pixel uncertainties.

Each of the following algorithms produces such `TrackObs` objects, which are saved intermediately as FITS files.

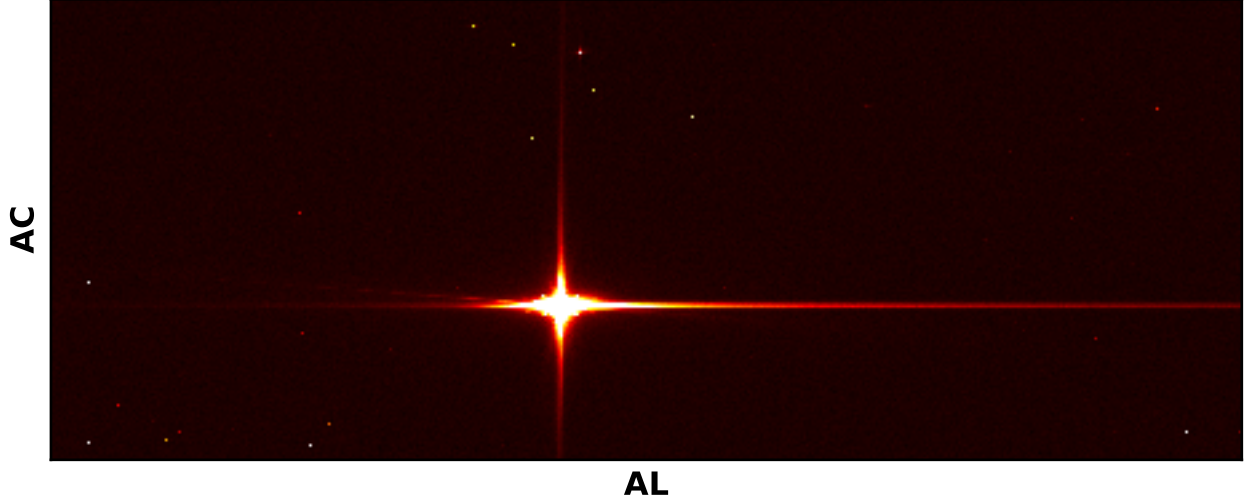


Figure 4: An example of a SM-SIF source image. Axis orientation is the standard Gaia convention, with the serial register to the right and readout to the upper right corner. The aspect ratio has not been corrected for pixel geometry and binning.

### 3.1.1 SM-SIF: Laplacian Edge Detection

Figure 4 demonstrates a typical image in the SM-SIF dataset. The image contains a bright star object with a PSF that contains long diffraction spikes and a saturated core. Smaller stars and what appear to be cosmic ray tracks can be detected by bare eye.

To automate the detection of cosmic ray tracks in these images, they must be successfully separated from the imaged astronomical objects. Such algorithms, usually used to remove unwanted cosmic rays, have been developed for earlier ground- and space-based observatories and can be used in to our advantage. For this dataset, we decided to utilize the L.A. COSMIC algorithm described in [4]. Briefly summarized, this algorithm utilizes a variation of Laplacian edge detection to discriminate cosmic ray tracks from stars by their sharp edges, which separates them from the smoother PSFs of stars.

In more detail, the algorithm first convolves a subsampled version of the original image with a derivative filter, which assigns high values to sharp edges. The resulting image is compared with a noise image, calculated by a  $5 \times 5$  median filter of the image and including the CCD readout noise. Pixels with a signal-to-noise-ratio above a high threshold are then selected as first candidates for noise pixels, after which a lower threshold is applied to neighbouring pixels. The identified pixels are then masked and the original image is cleaned, replacing the cosmic ray pixels with the median or mean of the surrounding, unmasked pixels. This procedure, starting from the convolution, is applied iteratively until either no more pixels are identified or the maximum number of iterations has been reached.

This analysis utilizes a Python implementation of the L.A. COSMIC algorithm called ASTRO-SCRAPPY [5] for cosmic ray identification.

First attempts at simply applying ASTRO-SCRAPPY to images like the one in Fig. 4 revealed that the bright stars in most of the SM-SIF images were still falsely detected as cosmics. The reason for this is their saturation of the detector – see especially the middle panel of Figure 5. The bright star, as it saturates the detector, causes a strong blooming effect, saturating even



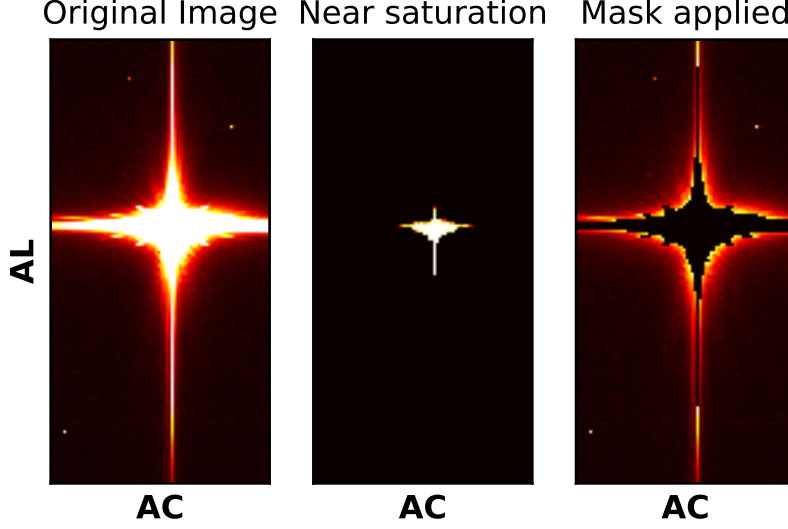


Figure 5: Bright stars in SM-SIF and their masking. *Left:* The appearance of the bright star from Fig. 4. *Middle:* The left panel after applying a high threshold. The bright pixels are fully saturated and show a sharp charge transfer pattern. *Right:* The left panel after applying bright star filtering - the core and a large part of the diffraction spikes are masked.

the pixels read in after those exposed to the bright stars (the image read out direction is to the top). As the panel shows, the resulting pattern has a very sharp gradient, which is picked up by the Laplacian filter and mis-identified as a cosmic.

To eliminate these bright star cores and mask their extended field, in which cosmics will disappear under the poisson noise of the star, we apply a simple algorithm before extracting cosmics. Its steps are:

1. Object finding: Construct a mask set to 1 for every pixel above a low threshold and apply object labelling to this mask. This step essentially catches all stars and cosmics.
2. Saturated object identification: Of all objects, select only those that contain at least one pixel at a value of 65535 ADU, the maximum of a 16-bit integer.
3. Masking: Add all the objects identified in the previous steps to a mask (set to 1 for invalid pixels). This mask can be forwarded to `ASTRO-SCRAPPY`, which supports masked regions and does not use them for cosmic identification.

The results of this algorithm can be seen in the right panel of Fig. 5, where all pixels belonging to the bright star are set to zero. There still remain some parts of the bright star's PSF - these, however, are very close to the bias level and decay smoothly, thus causing no issues.

Using the bright star removal, we then constructed the following algorithm to extract cosmics from SM-SIF observations:

1. Read in the image and construct the bright star mask.
2. Subtract the AL-dependent bias and apply the gain-scale, as L.A. COSMIC requires an image in units of electrons to calculate the Poisson noise.

3. Apply L.A. COSMIC, which returns the cosmic ray mask and a cleaned image. The values of the cleaned image are calculated by taking the mean of the surrounding, non masked pixels in a  $5 \times 5$  pixel area, i.e. for a pixel located at  $(AL, AC)$ , the cleaned value is

$$\text{Clean}(AL, AC) = \left( \sum_{\substack{\text{mask=False} \\ AL-2 \leq i \leq AL+2 \\ AC-2 \leq j \leq AC+2}} \text{Source}(i, j) \right) / \left( \sum_{\substack{\text{mask=False} \\ AL-2 \leq i \leq AL+2 \\ AC-2 \leq j \leq AC+2}} 1 \right) \quad (1)$$

4. Calculate the signal array as  $\text{Signal}(AL, AC) = \text{Source}(AL, AC) - \text{Clean}(AL, AC)$
5. Calculate the uncertainty array as the square root of the variance of the signal array, calculated as

$$\text{Var}(AL, AC) = \text{Source}(AL, AC) + \sigma_{\text{rn}}^2 + \sum_{\substack{\text{mask=False} \\ AL-2 \leq i \leq AL+2 \\ AC-2 \leq j \leq AC+2}} \frac{\text{Source}(i, j) + \sigma_{\text{rn}}^2}{N_{\text{mean}}(AL, AC)}, \quad (2)$$

with  $\sigma_{\text{rn}}^2$  denoting the pixel readnoise, and  $N_{\text{mean}}(AL, AC)$  denoting the number of pixels used to calculate the mean for cleaning this pixel. As one can see from the equation, the noise of the source image has been assumed to be Poissonian noise and readnoise.

6. Extract tracks from the mask, signal and uncertainty images and save them as a TrackObs, as described above.

The parameters used for the L.A. COSMIC algorithm were first selected manually, by applying the extraction algorithm to individual images and viewing the output. Selection criteria were that no stars should be misidentified as cosmic ray tracks, and that cosmic ray tracks should be complete, meaning that a single, long track should result in exactly one identified object.

To verify and further tune the extraction algorithm with regards to the detection efficiency and measurement accuracy, it was necessary to have known input data with which to compare the measured output. To that extent, we utilized a previously developed CCD cosmic ray events simulator called T.A.R.S. (Tools for Astronomical Radiation Simulations) developed at ESTEC [6]. The tool simulates the energy deposition of a cosmic ray in a CCD detector along with the charge collection to produce simulated CCD images. The first implementation of the tool was built to simulate Gaia images and was thus suited for this activity.

To test the extraction algorithm, we first prepared an SM observation by cleaning it with L.A. COSMIC, removing any existing cosmics that would be detected by our algorithm. For the source image, we selected a large image without any saturated stars, as these cosmics could not be detected there in any case and we were more interested in 'normal' stars.

We then generated cosmic ray tracks on the image, saving the locations and deposited energies for every simulated cosmic. Afterwards, we ran the extraction algorithm from above and analysed the results by associating each simulated track with the tracks recovered in its image location. Figure 6 shows two figures of merit for the extraction algorithm: The detection efficiency and the error of the extracted energies with respect to the calculated uncertainty.

The detection efficiency (left panel of Fig. 6) shows the number of detected tracks per input track - a value of 1 being a correct detection, a value of 0 being a non-detection and a value greater than 1 implying that parts of the track could not be identified, leading to multiple detections. The overall number of correctly detected events is greater than 97%.

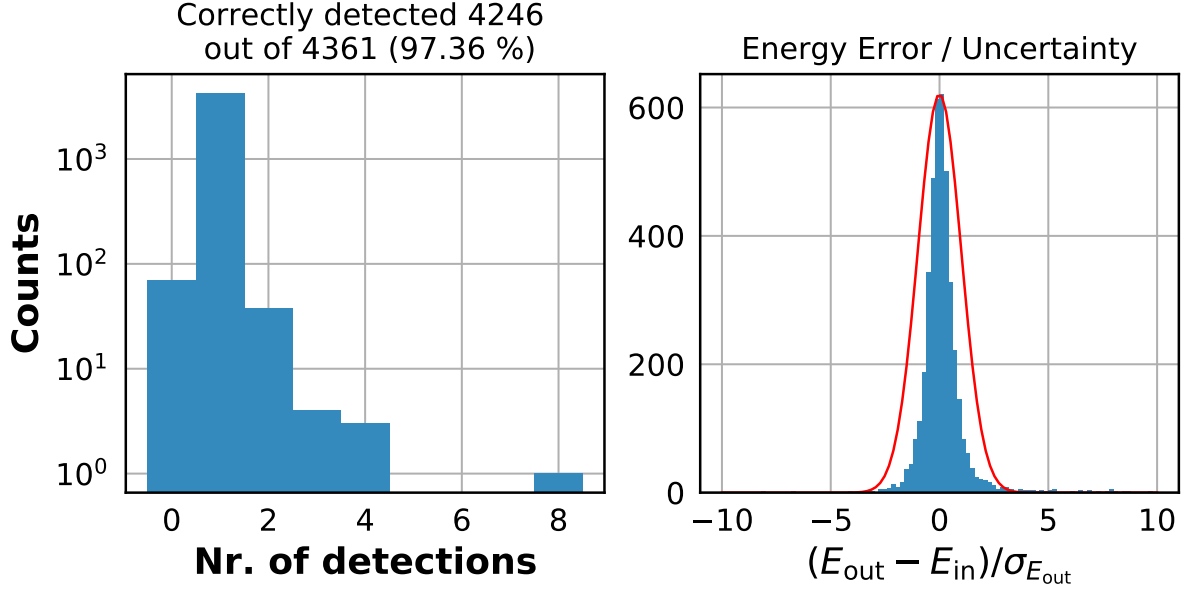


Figure 6: Figures of merit for the SM-SIF extraction. *Left:* Histogram of the number of associated track detections per input cosmic. A value of 0 is a non-detection, a value of 1 is a correct detection, a value of  $> 1$  means a single cosmic has resulted in multiple detections. *Right:* Histogram of the difference between measured and input track energy per cosmic, normalized by the calculated uncertainty. A Gaussian with a standard deviation of 1 is plotted in red.

The error in extracted energies is shown in the right panel of Fig. 6, and plots the difference between the measured track energy and the ionizing energy dumped by the simulated cosmic on the detector. For each track, this value is normalized by the uncertainty of the measured energy. A comparison of the resulting histogram with a Gaussian with a standard deviation of 1 shows that the energy measurements are correct within their estimated errors.

### 3.1.2 BAM-OBS: Boxcar filtering / Stacking

Figure 7 shows a series of observations from FOV 1 of the BAM. The three images were recorded in series (left to right) with an exposure time of about 23 seconds each. It is visibly apparent that the BAM patterns do not significantly differ between single observations, yet all three images contain pixels deviating from the pattern that only exist in one given image. These spurious pixels are a clear sign of cosmic rays.

Figure 8 further demonstrates that the interference pattern does not vary aside from Poissonian noise between the three images. It also shows why a simple application of the L.A. COSMIC algorithm to this dataset will immediately fail: The interference patterns themselves are very sharp, with a peak-to-trough distance of about 2 to 3 pixels. The algorithm is totally blinded by the pattern and can not separate it from the cosemics.

For our track extraction algorithm, we exploited the fact that the interference pattern is relatively static. In astronomy, multiple exposures of the same object (assuming constant luminosity) are used to reject the spurious high pixel values caused by cosmic rays – a similar algorithm is used to reject cosemics in BamObservations when analyzing the pattern. In our case, we can simply do the reverse, and accept pixels that vary significantly stronger than the inherent

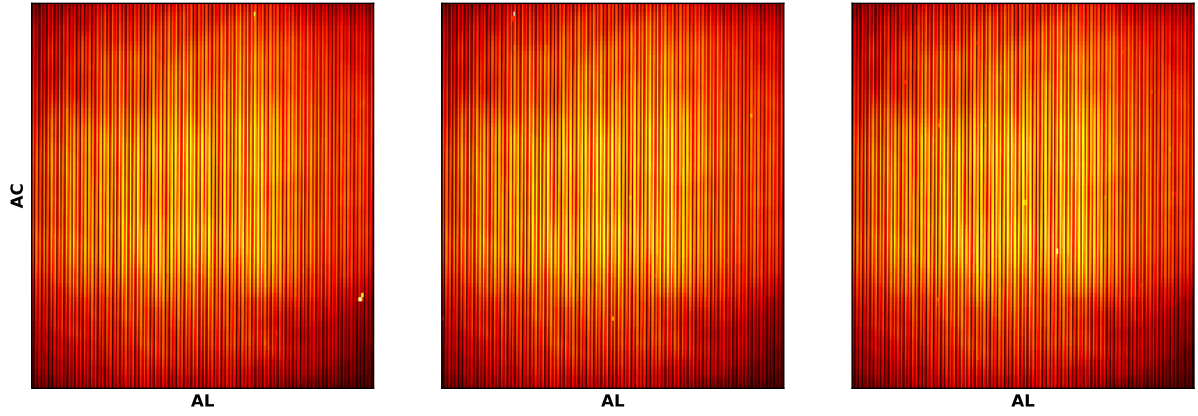


Figure 7: A series of three BamObservations from FOV 1. The dimensions have been rescaled to the physical dimension of the CCD – each image pixel represents a  $10 \times 120 \mu m$  (AL  $\times$  AC) sample. Individual pixels that deviate from the patterns can be recognized as cosmic ray tracks.

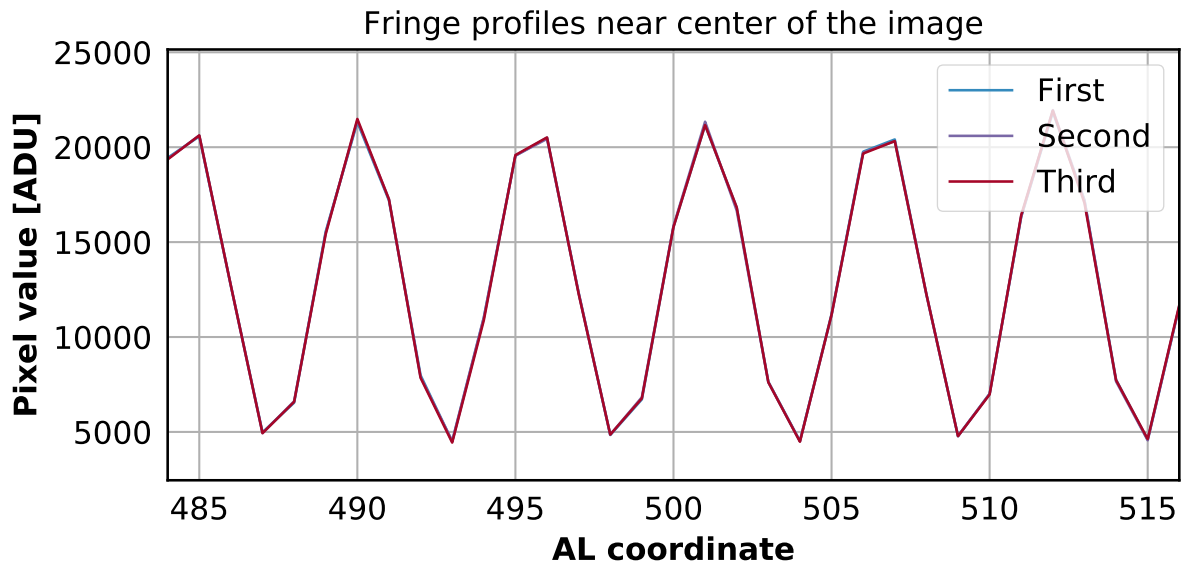


Figure 8: The ADU values of the three patterns in Fig. 7 plotted along the center of the images in AC across a narrow AL region in the center. The interference pattern is very sharp, with a period of about 5 pixels in AL. The fringes of all three images overlap very closely.

counting and readout noise as the tracks of cosmic rays.

The concrete algorithm assumes as an input a time-ordered series of BamObservations from a single FOV, and operates on a range of observations, which we refer to as a BoxCar. To create a good estimate of the pattern across a period of time where it does not vary, we chose to use a total of 7 images at once - a central image for cosmic ray extraction and the three images recorded immediately after and before. The extraction algorithm works as follows.

1. Construct the BoxCar: Read the first 7 BamObservations – applying the electron conversion gain – into a three dimensional array with the indices  $t$ ,  $AL$  and  $AC$  - with  $t \in [1, \dots, 7]$  ( $t = 4$  being the central observation), and  $AL$  and  $AC$  denoting the image dimensions.
2. Sigma-clipping: Compute the median and standard deviation across time and mark all the pixels whose difference to the median exceeds a chosen threshold - these are likely cosemics, and will not be used to calculate the pattern.
3. Calculate a two-dimensional image array of the pattern, by taking the mean across time without using the sigma-clipped pixels, e.g.

$$\text{Pattern}(AL, AC) = \left( \sum_{t \text{ not clipped}} \text{BoxCar}(t, AL, AC) \right) / \left( \sum_{t \text{ not clipped}} 1 \right). \quad (3)$$

4. Calculate the signal array by subtracting the Pattern from the central array in the BoxCar, e.g.

$$\text{Signal}(AL, AC) = \text{BoxCar}(t = 4, AL, AC) - \text{Pattern}(AL, AC). \quad (4)$$

5. Calculate the uncertainty array as the square root of the variance of the signal array, calculated as

$$\text{Var}(AL, AC) = \text{BoxCar}(t = 4, AL, AC) + \sigma_{\text{rn}}^2 + \sum_{t \text{ not clipped}} \frac{\text{Boxcar}(t, AL, AC) + \sigma_{\text{rn}}^2}{N_{\text{mean}}(AL, AC)}, \quad (5)$$

with  $\sigma_{\text{rn}}^2$  denoting the pixel readnoise, and  $N_{\text{mean}}(AL, AC)$  denoting the number of pixels in the BoxCar that have not been rejected by sigma clipping at this  $(AL, AC)$  pair. As one can see from the equation, the noise of the BamObservations has been assumed to be Poissonian noise and readnoise.

6. Create the cosmic mask: Divide the signal array by the uncertainty array and create the boolean mask array, set to True where the signal/uncertainty ratio is greater than a threshold  $f$ . Due to the strong photon background, some cosmic pixels are missed in this step. We attempt to mitigate this by investigating all the pixels that neighbour the previously masked pixels, found via binary dilation, and applying for them a signal/uncertainty threshold of  $r_{\text{Neighbour}} \cdot f$ , adding the pixels that satisfy this condition to the mask as well.
7. Apply a track connection algorithm: Even after the previous step, visually apparent long cosmic ray tracks still remain separated - refer to the left panel of Figure 9 for an illustration. Tracks that are oriented in the AL direction see a strongly variable photon background, seen as fluctuations in the uncertainty array. These pixels are less likely to be detected by the thresholding of before. Especially towards the center of the image, where the pattern amplitude is very high, this leads to a single, long cosmic being detected as

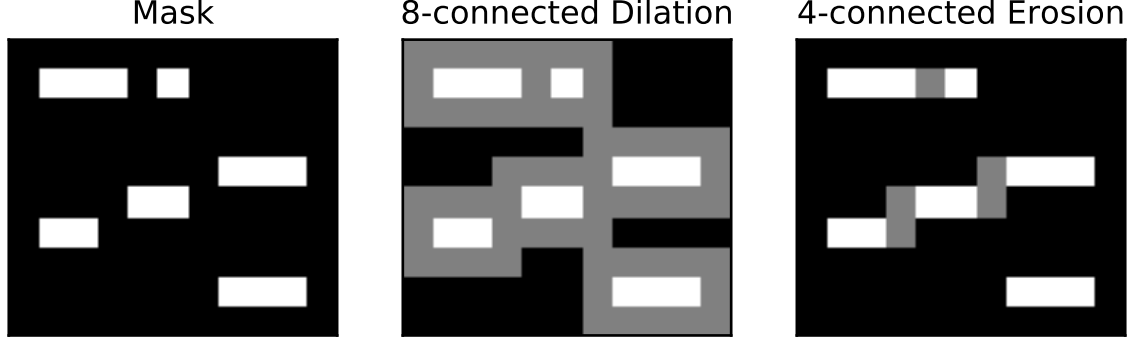


Figure 9: Illustration of the track connection algorithm used for this dataset. *Left:* An example of a cosmic ray mask where ray tracks oriented in the AL direction have not been fully identified. *Middle:* Application of 8-connected binary dilation to the mask - the added pixels are shown in gray. *Right:* Additional application of 4-connected binary erosion, shown in the gray pixels, which are examined by the connection algorithm.

multiple short cosemics in the troughs of the interference pattern, which eventually leads to an overestimation of the cosmic ray rate.

In an attempt to correct for this track separation, we implemented a track connection algorithm: We first apply to the mask an 8-connected binary dilation followed by a 4-connected binary erosion (this is a variant on binary closing, which uses the same structuring elements for the dilation and erosion instead), illustrated in the middle and right panels of Fig. 9. As the Figure shows, this marks several pixels of interest that may connect the separated objects.

We then examine each of the newly merged tracks, which are identified by the standard labelling algorithm. For each of the tracks, we first calculate the mean of the previously identified signal pixels of this track,  $m_{\text{Object}}$ . Assuming that a cosmic ray creating a long track should have a constant energy deposition per path length, we examine each of the candidate pixels (at the location  $(AL, AC)$ ), checking for the condition

$$|Signal(AL, AC) - m_{\text{Object}}| > cfac * \sqrt{Var(AL, AC)}, \quad (6)$$

with  $cfac$  being an algorithm parameter. The above equation checks whether the signal value of the candidate pixel is sufficiently close to the value of the other pixels in this object. If a pixel satisfies the condition, it is added to the pixel mask.

8. Extract tracks from the mask, signal and uncertainty images and save them as a TrackObs, as in SM-SIF and described above.
9. Update the BoxCar: Read in the next BamObservation and replace the oldest pattern in the BoxCar with it – the BoxCar is utilized as a FIFO, in this case. We then update the central pixel index and repeat starting step 2 until all BamObservations have been processed.

As we did for the SM-SIF extraction, we also tested the BAM-OBS extraction algorithm using the T.A.R.S. package. In the case of BAM-OBS, we first created a series of fake BAM patterns by calculating a running median over a set of 100 successive BamObservations. Concretely, we took for every observation the four preceding and four next observation and calculated the

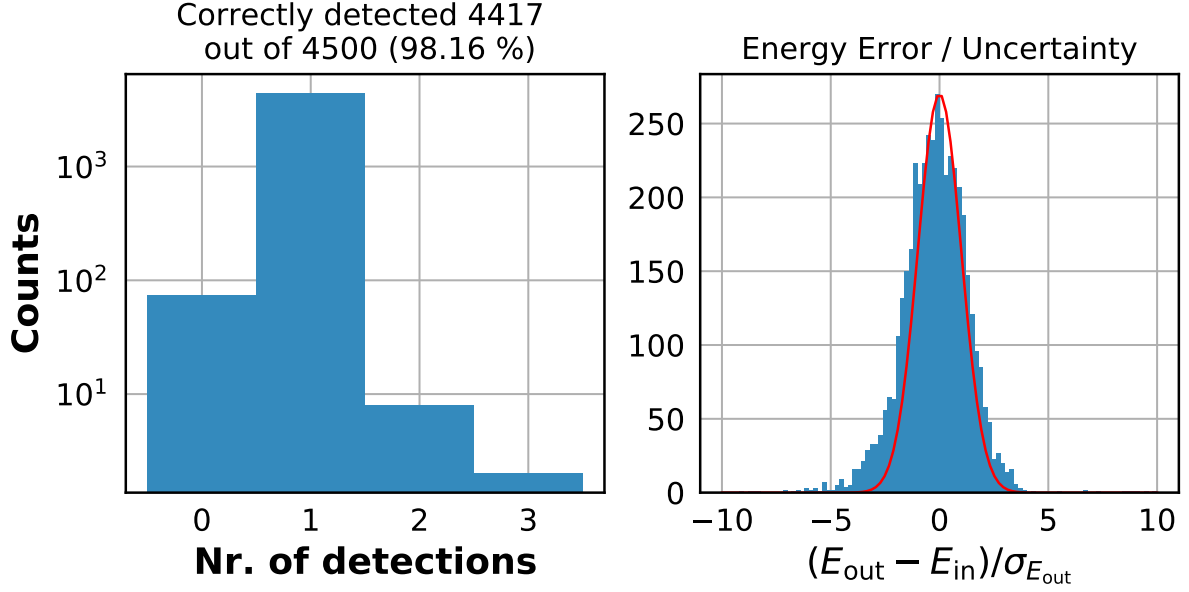


Figure 10: Figures of merit for the BAM-OBS extraction, as in Fig. 6.

median image of these observations. Afterwards, we added Poissonian counting noise to each image.

We then added simulated cosmic ray tracks to all images and ran the extraction algorithm over the list of observations. As in the SM-SIF verification, we associated each input track with the detected tracks in its image location. As in the previous verification, Figure 10 shows the detection efficiency and normed energy error for the BAM-OBS extraction, once again returning consistent values for both. Note that these plots show results for FOV 2, whose interference pattern has a higher intensity and thus noise.

To justify the use of the connection algorithm (step 7 above), Fig. 11 shows the figures of merit for the same input data as Fig. 10 when the connection algorithm is not used. While the measured energies don't differ significantly, the number of input tracks recognised as multiple output tracks is lowered by the connection algorithm.

### 3.1.3 BAM-SIF: Thresholding

Figure 12 depicts an example of a BAM-SIF image used in this extraction. These images both contain the interference pattern in the BAM-OBS images, as well as off-pattern regions showing the TDI footprint of the pattern, which is only AC-dependent. These images are much less frequently sampled than the BAM-OBS images – usually a handful of images twice per week – which is why they are not very useful for the purpose of particle monitors. However, the off-pattern regions offer a good opportunity to extract cosmons from the BAM chip with a much lower and simpler background than that of BAM-OBS, allowing us to examine the effect the pattern has on the above extraction.

As this dataset is very attractive for cosmic ray extraction, it has been used before for cosmic ray studies, see [2]. This algorithm builds on the previous extraction, merging it into the TrackObs-based approach utilized here.

For the extraction of cosmons, we utilize the two regions outlined in blue and green boxes in

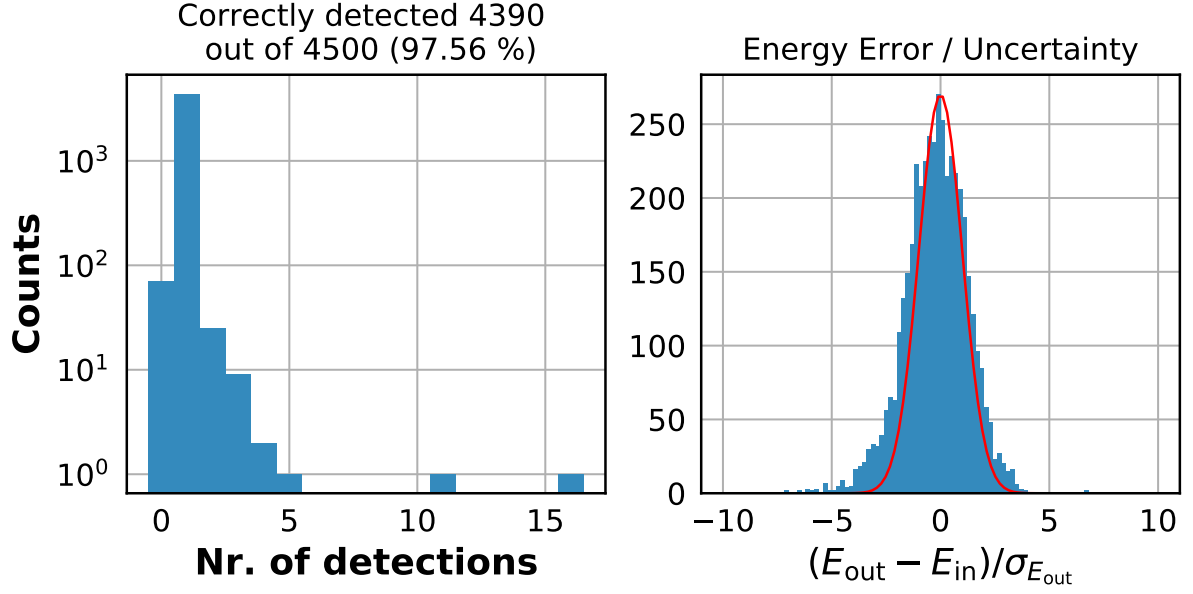


Figure 11: Figures of merit for the BAM-OBS extraction without the track connection algorithm. Compared to Fig. 10, individual cosmics are more likely to be identified as multiple tracks.

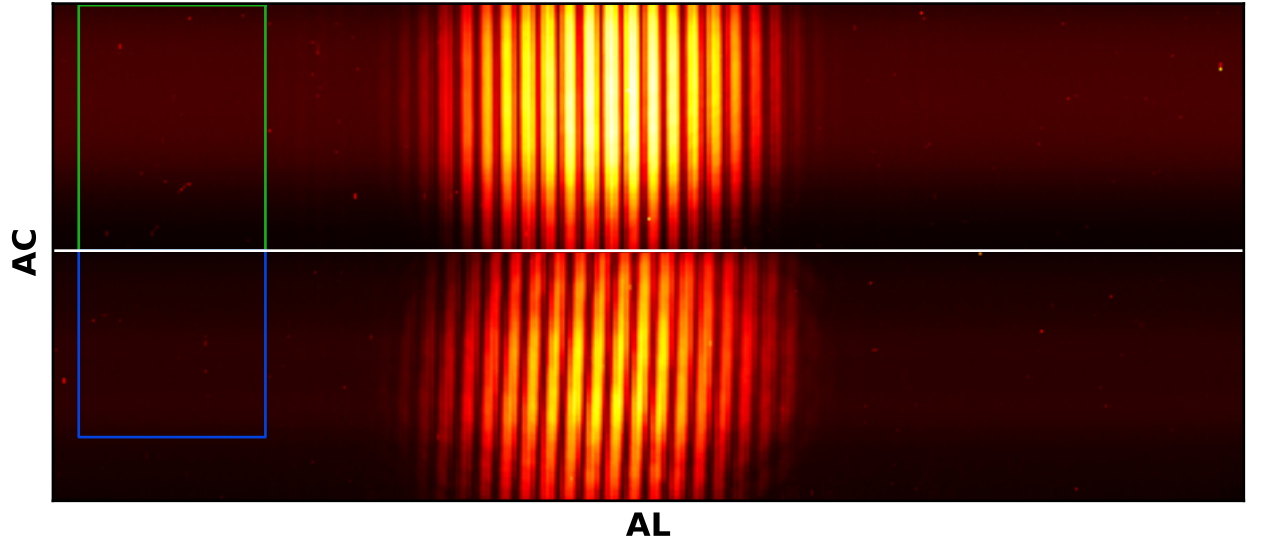


Figure 12: An example of a BAM-SIF source image, which has been rescaled to the physical CCD dimensions. The top half of the image shows the AC region of FOV 2, the bottom half that of FOV 1. The off-pattern regions in green and blue boxes are used for cosmic extraction.



Fig. 12. A smaller region in AC has been selected for the FOV 1 region - the ignored region has shown a high stray light contamination, which causes mis-identifications.

The extraction algorithm for this dataset is then (for each region separately):

1. Read in the image and apply the gain scale.
2. Determine the AC-dependent TDI background: Take the 99 columns with the highest AL value (to the left in Fig. 12) and remove the outliers – for each column within, discard the highest and lowest 25 % of pixel values to remove cosmics. Then, calculate the mean and standard deviation of the remaining image pixels for each line. The output are two one-dimensional arrays of the AC-dependent TDI background and its standard deviation. The removal of the TDI background is demonstrated in Figure 13.

3. Calculate the signal and variance arrays as

$$Signal(AL, AC) = Source(AL, AC) - Background(AC) \quad (7)$$

and

$$Var(AL, AC) = Source(AL, AC) + \sigma_{rn}^2 + \sigma_{Background}^2(AC), \quad (8)$$

with  $\sigma_{rn}$  being the readnoise and once again assuming Poisson noise in the source array. The uncertainty array is then the square root of the variance array.

4. Construct the mask array in exactly the same way as for BAM-OBS, including the track connection. While this should not be necessary here, we want to have as similar as possible algorithms for BAM-OBS and BAM-SIF, to see the effect of the pattern.
5. Extract tracks from the mask, signal and uncertainty images and save them as a TrackObs, as in the algorithms above.

All steps until the calculation of the variance are the same as in [2], which only returned the signal images.

As in both previous algorithms, we once again verified the BAM-SIF extraction against simulated cosmics. In this case, a fake TDI background was obtained by removing outliers from a BAM-SIF observation. A background image is then randomly sampled line by line from the extracted background. Afterwards, simulated cosmic ray tracks are added on top of the produced image and extracted via the BAM-SIF algorithm. The figures of merit obtained in the same way as before are plotted in Figure 14, showing a very good detection efficiency and energy errors within the bounds expected from the calculated uncertainties.

### 3.1.4 Extraction Pipelines

Having specified the extraction of cosmic ray tracks from the images themselves, we now describe the extraction pipelines used to systematically retrieve the image data files from the Gaia Main Database in order to produce and save TrackObs from them.

While the pipeline is implemented in Apache Spark, we will describe the pipeline independently of the utilized software.

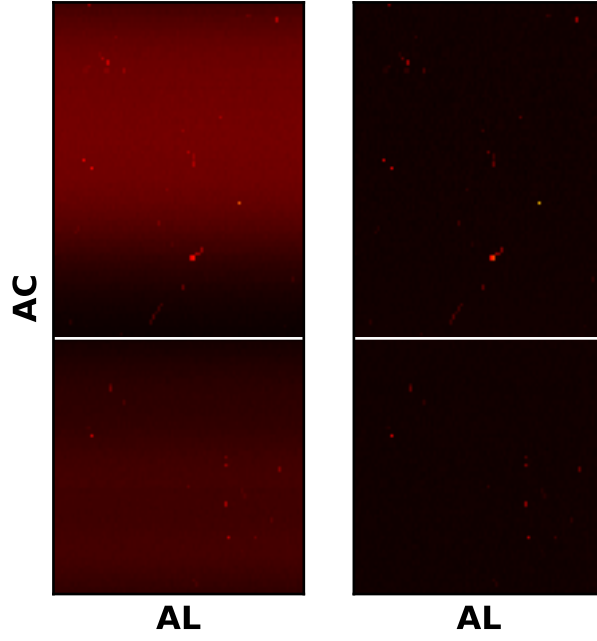


Figure 13: *Left:* The two extraction regions from Fig. 12. *Right:* The same image with the TDI background removed, leaving only cosmic ray tracks.

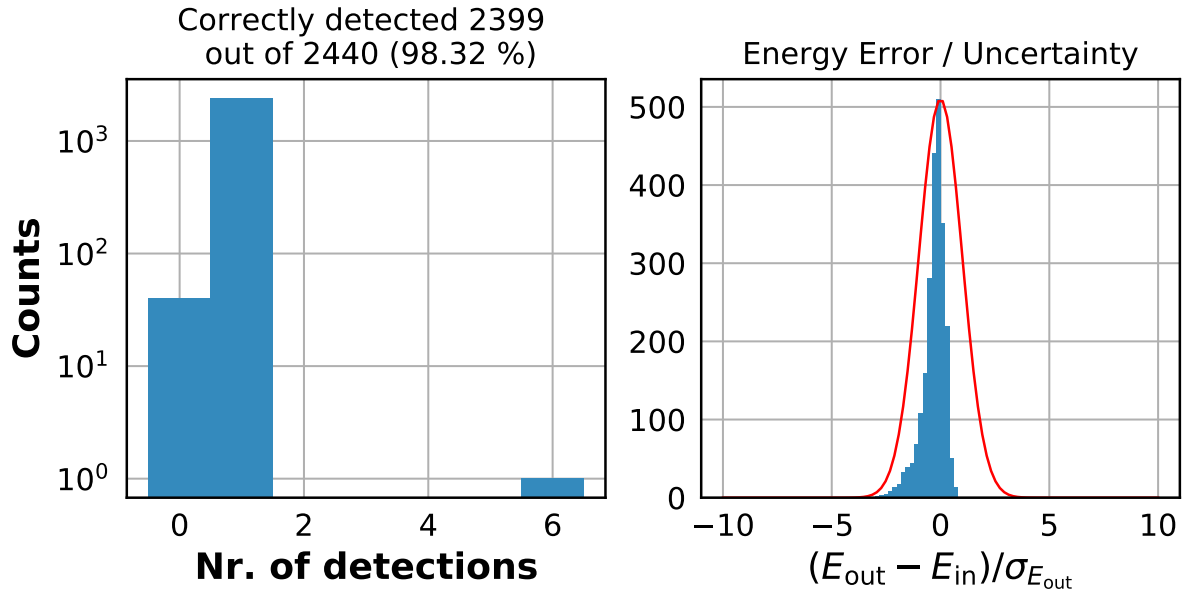


Figure 14: Figures of merit for the BAM-SIF extraction, as in Fig. 6.

## SM-SIF and BAM-SIF

Input data for the SM-SIF and BAM-SIF are both saved in in FITS files, each containing the data of a single observation. All input files are stored in a SIF directory, in individual folders sorted by date (format: yyyy-mm-dd). As the analysis of one observation is fully independent, this extraction is very easy to parallelize.

For the SM-SIF and BAM-SIF dataset respectively, we use the extraction algorithms listed above to obtain TrackObs from images. The steps of the extraction pipeline around the extraction are:

1. Group all files by their date
2. For each date:
  - 2.1. Create a buffer for TrackObs
  - 2.2. For each file:
    - 2.2.1. Extract a TrackObs from the image using algorithms from Sec. 3.1.1/3.1.3
    - 2.2.2. Add the TrackObs to the buffer
  - 2.3. Write all the TrackObs in the buffer to a FITS file with the directory structure  
yyyy/mm/dd/[source]\_OBMT\_START\_[obmt of first TrackObs].fits,  
with [source] being SM-SIF or BAM-SIF, respectively.

Step 2 and all substeps within are completely independent for each day and can as such very easily be parallelized.

## BAM-OBS

The input data for the BAM-OBS dataset is saved in compressed gbin files, each containing normally around 200 BamObservation images. The BAM-OBS extraction pipeline – see the flowchart in Fig. 15 – is designed to take into account that the algorithm requires the utilization of a BoxCar moving over a continuous list of BamObservations. Parallelization is achieved by sorting the gbin files by the acquisition time of their BamObservations and dividing this sequence into individual chunks, to be handled separately. The extracted TrackObs are saved once again in FITS files, in the same directory structure as in the SIF extraction above.

Note that for the processing of each file group, a small overlap to the next file group needs to be considered in order to process all the BamObservations: For two sequential groups of BamObservations A and B, the last observation in A needs to use the first three observations in B to fill the BoxCar for its extraction, and analogously for extracting the first observations in B.

## 3.2 Post-Processing

For the final data product described in Section 4.1, we derive additional data from the raw extracted tracks in a TrackObs. The algorithms for deriving these are outlined in this section.

### 3.2.1 Observed Flux

One of the first – and most simple – data products that can be derived from a TrackObs is the flux measured per observation. The calculation of a flux is in this case very simple: Given an exposed CCD area  $A_{\text{exp}}$  and exposure time  $T_{\text{exp}}$  the flux is simply calculated as

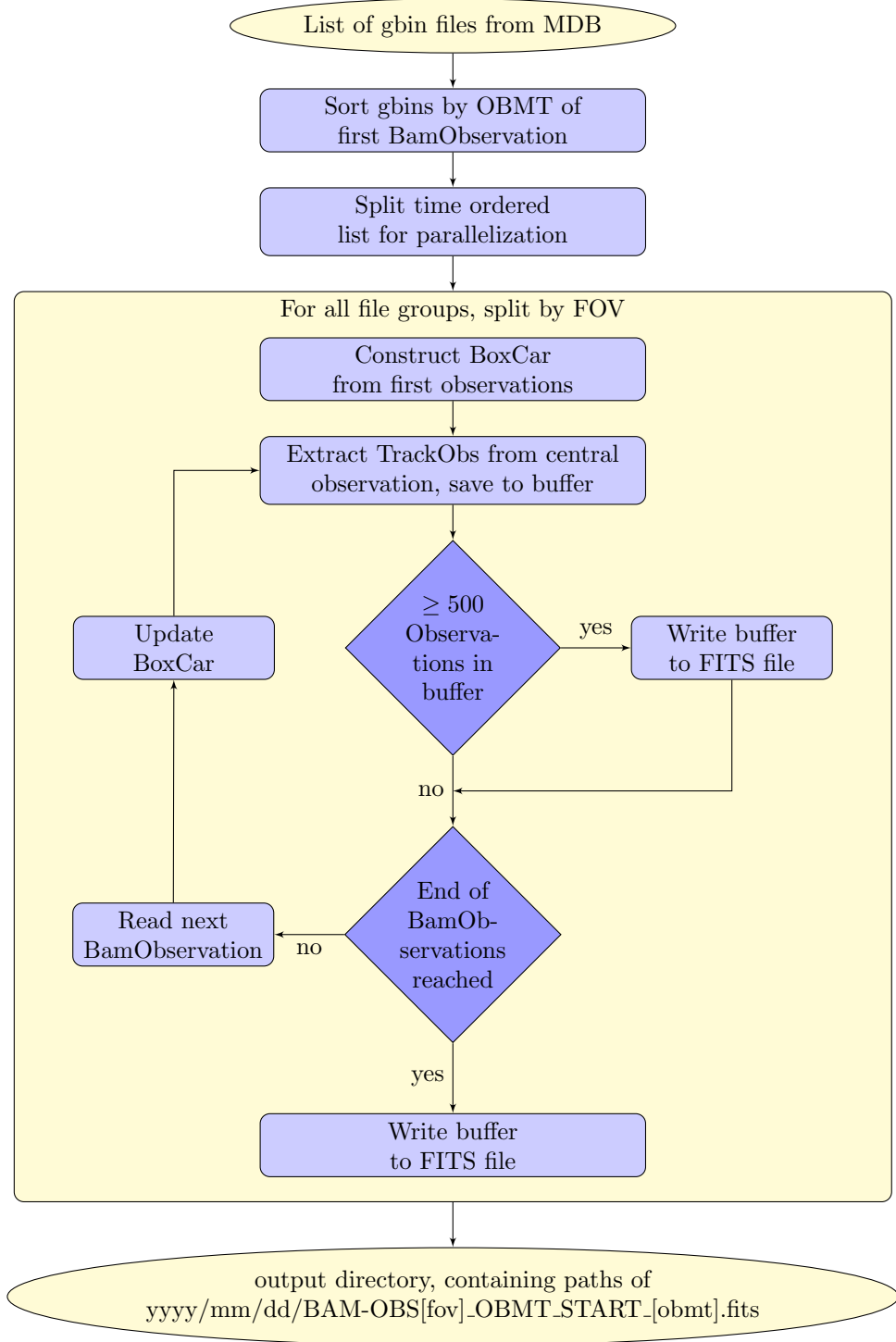


Figure 15: Flowchart for the extraction pipeline of the BAM-OBS dataset. The `[obmt]` in the last element refers to the acquisition time of the first `BamObservation` in the output file. The operations on each file group can be parallelized.

$$F = \frac{\text{Counts}}{A_{\text{exp}} \cdot T_{\text{exp}}} \quad (9)$$

In the following, we determine these values for SM-SIF and BAM observations, whose computation are not trivial due to the use of TDI (and in the latter case, staring) observations and the size of the used AL readout windows relative to the physical CCD size.

For our calculations: A Gaia CCD is segmented into 4500 pixels in AL direction and 1980 pixels in AC direction - of the latter, 14 pixels are in the pre-scan region and not used for observations. Pixels are  $10 \times 30 \mu\text{m}$  (AL  $\times$  AC) and TDI transfer time is  $T_{\text{TDI}} = 0.9828 \text{ ms/pixel}$ , so the effective CCD exposure time is 4.4226 s.

### SM-SIF

The readout of the SM-SIF images are a continuous readout of the signal acquired from TDI (see also Fig. 16). To avoid saturation from bright stars, Gate 12 is permanently activated, leading to a readout of only the last 2900 CCD lines in AL. Pixels are then binned to a factor of 2 in both AL and AC.

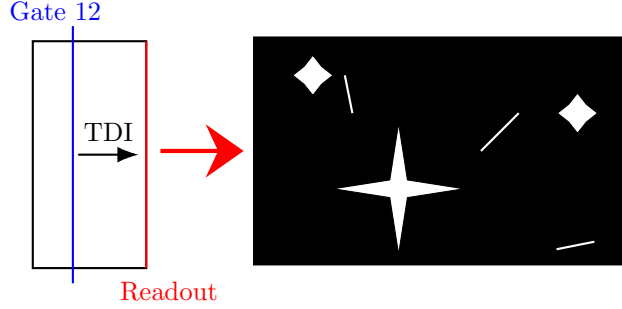


Figure 16: Scheme of an SM acquisition. The CCD area of each chip after Gate 12 is continuously read out via TDI (left). The resulting image (right) has the same dimension in AC but a dimension in AL dependent on the duration of the readout.

For clarity in notation, I will, from now on, use the following terms:

- A **CCD pixel** refers to a physical pixel on the actual CCD - i.e. one physical unit of  $10 \times 30 \mu\text{m}$
- An **image pixel** refers to a pixel in the image recorded after the TDI scan, including binning. While it still carries spatial information in AC, its AL information is a result of the integration over the sampled region in AL - it encodes time and location. Thus, an image pixel does not correspond to a single physical pixel, but rather the history of an entire AL line during TDI.

The resulting image, then, consists of 983 image pixels in AC and a variable number  $N_{\text{AL}}$  of image pixels in AL - depending on the time over which data has been recorded (which is  $T_{\text{rec}} = T_{\text{TDI}} \cdot 2N_{\text{AL}}$ ).

From a physical pixel perspective, each CCD pixel has been exposed for the time  $T_{\text{rec}}$ , with the information of individual exposures spread over different image pixels in AL. The exposed area corresponds to that of  $2900 \times 1966$  physical pixels.

From an image pixel perspective, each pixel has been exposed for exactly  $T_{\text{TDI}} \cdot 2900$ , regardless of  $N_{\text{AL}}$ . The physical area of each image pixel corresponds to that of 4 physical pixels, due to binning.

As one can tell from the above, the physical perspective has a lower number of pixels, but a higher exposure time per pixel (assuming  $2N_{\text{AL}} > 2900$ ). The image perspective has more pixels, but lower exposure times. Irregardless of the perspective, the product  $A_{\text{exp}} \cdot T_{\text{exp}}$  – which is of interest for computing fluxes – is always the same, being

$$\text{Physical : } A_{\text{exp}} \cdot T_{\text{exp}} = (2900 \cdot 1966 \cdot 10 \mu\text{m} \cdot 30 \mu\text{m}) \cdot (T_{\text{TDI}} \cdot 2N_{\text{AL}}) \quad (10)$$

$$\text{Image : } A_{\text{exp}} \cdot T_{\text{exp}} = (N_{\text{AL}} \cdot 893 \cdot 4 \cdot 10 \mu\text{m} \cdot 30 \mu\text{m}) \cdot (T_{\text{TDI}} \cdot 2900) \quad (11)$$

$$= 33.62 \left( \frac{N_{\text{AL}}}{1000} \right) \text{ cm}^2 \text{ s} \quad (12)$$

Going forward, we will use the image perspective, as it is consistent with prior studies of the BAM image and more useful for the next section. It is also, as we'll see later, easier to treat masked pixels. Nevertheless, for the calculation of fluxes, both perspectives are equivalent.

### BAM-OBS and BAM-SIF

For the BamObservations, pixels are binned by a factor of 4 in AC and not binned at all in AL.

The determination of an exposure time in both pictures for the BamObservations is complicated by the fact that it is not only a TDI observation, but also contains a period of staring. Normally, the BAM CCD is operated in TDI mode, in the same way as the SM CCD. To record the interference pattern, an exception is made – for an integration time of 19 s, the shifting of charges is stopped, and the chip operates in a kind of staring mode. In nominal BamObservations, data is then only acquired from sub-regions of the CCD that have been illuminated by the BAM pattern. BAM-SIF images use the entire CCD in AL direction along a restricted AC range.

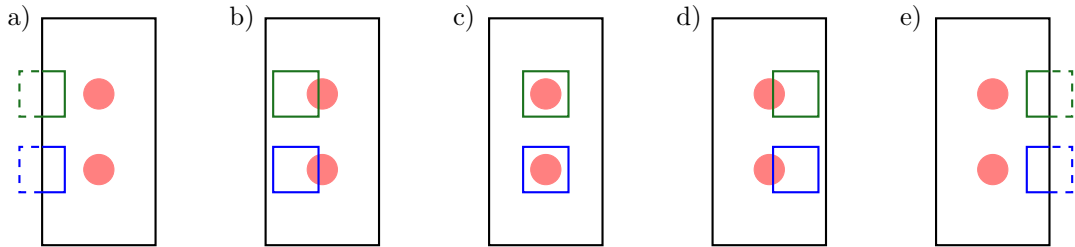


Figure 17: Readout scheme for a BAM pattern in BAM-OBS. BAM fringe patterns are outlined as red circles. a) Create the readout windows. b) Move the readout windows via TDI. c) Stare observation – charges are not shifted. d) Move readout windows via TDI. e) Readout, recording only contents of windows.

The observation principle and the effect on imaging particle tracks is shown in Fig. 17. Effectively, data is collected from two windows, outlined in green and blue. These two windows have an unbinned length of 1000 in AL and 320 in AC each. They first move across the chip, stay in place during the staring, and then move again to the readout register.

In BAM-SIF observations, we can make a similar estimation of the readout window sizes, although the image dimensions are variable. We will from now on use two values  $N_{\text{AL}}$  and  $N_{\text{AC}}$  to refer to the binned image dimensions of the result image. The used unbinned lengths are then  $N_{\text{AL}}$  and  $4 N_{\text{AC}}$ , respectively – for reference, a nominal BamObservation has  $N_{\text{AL}} = 1000$  and  $N_{\text{AC}} = 80$

The exposure can once again be viewed in two pictures, the physical pixel and image pixel based ones. In the former, it is easy to see that the entire region both readout windows cross are exposed – the exposed area is  $4500 \times 4 N_{\text{AC}}$  pixels ( $\text{AL} \times \text{AC}$ ) per pattern. The exposure time, however, is more complex: First, the entire exposed region is exposed for  $N_{\text{AL}} \times T_{\text{TDI}}$ , due to the movement of the readout windows across the chip. Additionally, the region used for staring is observed for a further 19 s. The exposure time is, as such, AL-dependent.

In the image picture, the calculation is much easier. Viewing each readout window as a physical unit, the exposed area per window is simply  $N_{\text{AL}} \times 4 N_{\text{AC}}$  pixels. The exposure time is simply the TDI transfer time across the whole chip plus the staring time.

The product of exposed area and exposure time, which is of interest for flux determination, is in both pictures (per readout window)

$$\text{Physical : } A_{\text{exp}} \cdot T_{\text{exp}} = (4500 \cdot 10 \mu\text{m} \cdot 4 N_{\text{AC}} \cdot 30 \mu\text{m}) \cdot (N_{\text{AL}} \cdot T_{\text{TDI}}) \quad (13)$$

$$+ (N_{\text{AL}} \cdot 10 \mu\text{m} \cdot 4 N_{\text{AC}} \cdot 30 \mu\text{m}) \cdot (19 \text{ s}) \quad (14)$$

$$\text{Image : } A_{\text{exp}} \cdot T_{\text{exp}} = (N_{\text{AL}} \cdot 10 \mu\text{m} \cdot 4 N_{\text{AC}} \cdot 30 \mu\text{m}) \cdot (4500 T_{\text{TDI}} + 19 \text{ s}) \quad (15)$$

$$= 22.49 \text{ cm}^2 \text{ s} \quad (\text{for BAM-OBS}). \quad (16)$$

As the computation is more simple in the image picture, we will from now on use this picture.

### Masked Pixels

In both SM and BAM images, some image pixels can not be used for the determination of cosmics. In the SM, these pixels are mostly caused by the bright star masking algorithm. For the BAM, charge overflow effects can be caused by a high energy cosmic coincident with an interference fringe, fully saturating the pixel and causing errors in the readout of several following pixels.

Correcting the exposed area for these masked pixels is very easy in the image picture, as we simply discard these pixels, meaning one can calculate the exposed area using

$$A_{\text{exp}} = (N_{\text{tot}} - N_{\text{mask}}) \cdot b_{\text{AL}} \cdot b_{\text{AC}} \cdot 10 \cdot 30 \mu\text{m}^2 \quad (17)$$

With  $N_{\text{tot}}$  being the total number of image pixels,  $N_{\text{mask}}$  being the number of masked pixels and  $b_{\text{AL}}$  and  $b_{\text{AC}}$  being the binning factor in AL and AC, respectively.

### 3.2.2 Track Geometries

For particle tracks that pass through multiple pixels, it is possible to estimate the entrance angle of a cosmic into a CCD with respect to the FPA and the length of the path of the particle through the chip. For a selected subset of tracks, this was attempted in this dataset using either simple estimates or a fitting technique.

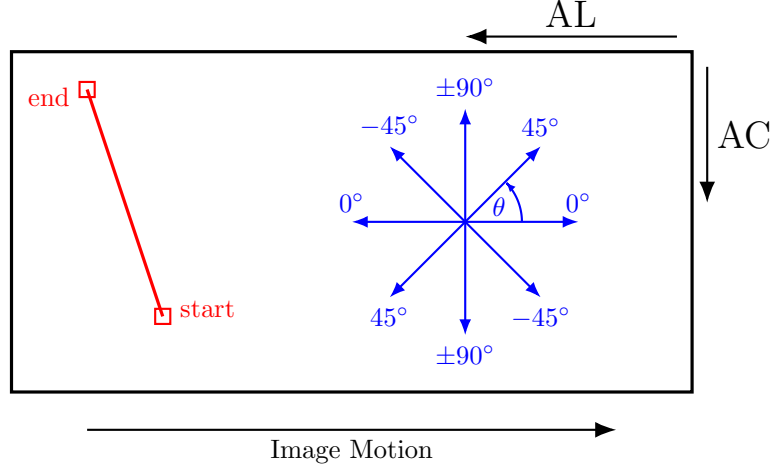


Figure 18: Definition of the track angle  $\theta$  in the standard Gaia CCD reference system. The red line on the left shows an example track, whose start is always defined as the end-point at the lower AL-coordinate. The blue compass rose shows the angle corresponding to different track orientations.

Figure 18 shows the definition of the angle  $\theta$  used in this analysis. It is obtained via the slope of a line describing the cosmic ray track in an (AL,AC) coordinate system, e.g.

$$\theta = \tan^{-1} \left( \frac{AC_{\text{end}} - AC_{\text{start}}}{AL_{\text{end}} - AL_{\text{start}}} \right) [^\circ], \quad (18)$$

where the starting point of a track is always defined to be its edge point at the lower AL value. We have made this choice because we do not determine at which of the edge points of a track the cosmic entered and exited the CCD. As such,  $\theta$  is defined in a range of  $[-90^\circ, 90^\circ]$ , with  $0^\circ$  corresponding to a track completely along the AL axis,  $\pm 90^\circ$  being oriented in the AC direction.

For a track that is one-dimensional, meaning that all its signal pixels are in one line or column, the angle and track length are constrained and can be estimated with a simple algorithm, outlined in Figure 19. The particle trajectory can be between two extremes: either perfectly aligned to the AL/AC direction or running between the corners of the signal pixels at the edge of the track.

In formulas: A track along an AL line has an angle  $\theta$  of

$$\theta = 0^\circ \pm \tan^{-1} \left( \frac{h_{\text{AC}}}{w_{\text{AL}} \times N_{\text{AL}}} \right) \quad (19)$$

and a length  $l$  of

$$w_{\text{AL}} \times (N_{\text{AL}} - 2) \leq l \leq \sqrt{(w_{\text{AL}} \times N_{\text{AL}})^2 + (h_{\text{AC}})^2}, \quad (20)$$

with

- $N_{\text{AL}}$  being the number of binned pixels in AL
- $w_{\text{AL}}$  being the binned pixel width (i.e. in AL direction)



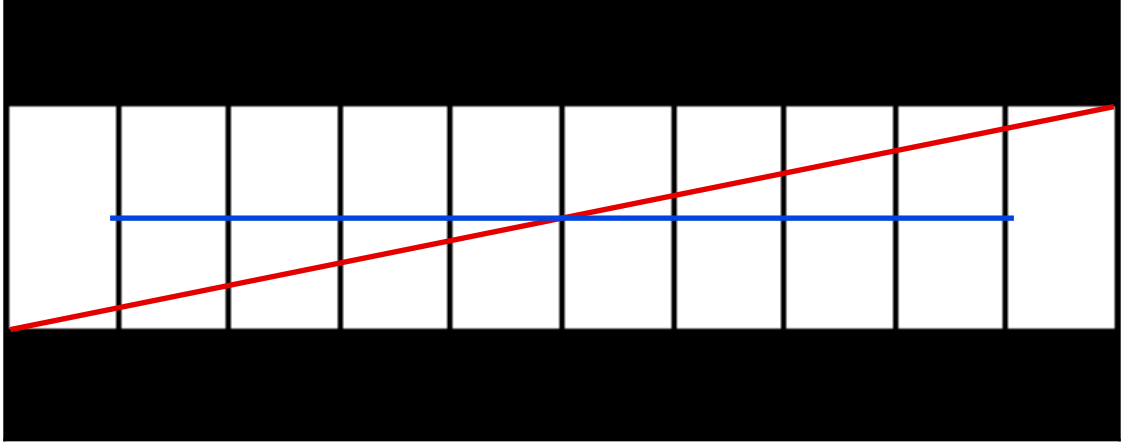


Figure 19: Possible track lengths and angles in a one-dimensional track in AL. The longest, steepest track in red passes through the corners of the most distant pixels, while the shortest, flattest track in blue starts and ends near the nearest edges of the outermost pixels.

- $h_{AC}$  being the binned pixel height.

This can be calculated analogously for a track along an AC column. For such tracks, we define our output values for the track angle and length as the mean of Equations 19 and 20 respectively. The uncertainty is then half the absolute value of the difference between the two extreme values of these equations.

For tracks that are not aligned with the AL and AC direction, we use a fitting routine to determine  $\theta$ , which we in turn use to determine  $l$ . The fitting routine is shown in Figure 20.

For each track image, we first determine the length of the bounding box in binned samples along the AL and AC direction. For each pixel coordinate in the longer axis, we determine the center of mass of the track image values along the shorter axis (e.g. if the image has more pixels in AL, we compute the center of mass along the AC axis for each line). The resulting pairs of values are saved in an  $(x, y)$  array, with the x-value denoting the location in AL, the y-value denoting the location in AC multiplied by the factor  $f_{\text{pix}} = 3 \times \frac{\text{AC binning}}{\text{AL binning}}$ .  $x$  and  $y$  are, as such, in the unit of binned AL pixels.

The rationale for using the center of mass is that can better estimate the trajectory of a cosmic if it causes a signal in, for instance, two adjacent pixels in a column – see the Figure 20.

For each  $(x, y)$  pair, we set an uncertainty ( $\sigma_x = 0.5, \sigma_y = 0.5f_{\text{pix}}$ ), representing the uncertainty of the detection location due to pixelisation. We then perform a fit of a simple fit to a line equation,

$$y(x) = m \cdot x + y_0, \quad (21)$$

using the `scipy` implementation of the Orthogonal Distance Regression algorithm [7]. This algorithm takes into account the uncertainties of both  $x$  and  $y$ , providing best fit values and estimations on the standard deviations  $\sigma_m$  and  $\sigma_{y_0}$  of the slope and offset.

To determine the angle theta and its uncertainty, we utilize the inverse tangent and gaussian error propagation, yielding

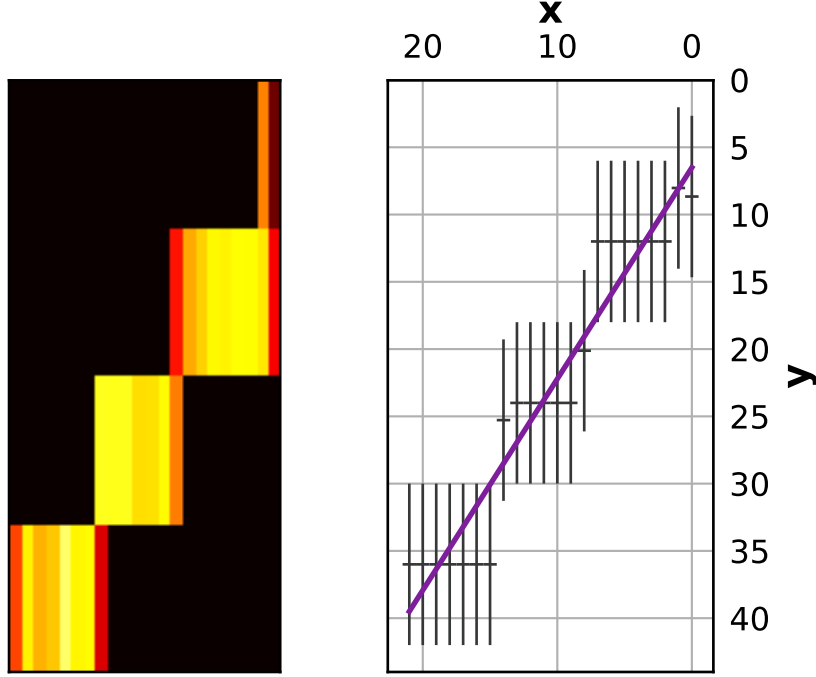


Figure 20: Fitting routine on two-dimensional tracks on the example of a BAM-OBS track. *Left:* Track image, with the pixels scaled to their physical dimension. *Right:* Best fit. The  $(x, y)$  pairs have been computed by calculating the center of mass along AC for each AL-line, with the error bars representing the pixel size. The purple line shows the best line fit as determined by orthogonal distance regression.

$$\theta = \tan^{-1}(m), \quad \sigma_\theta = \frac{1}{m^2 + 1} \sigma_m \quad (22)$$

to determine the length of the cosmic, we define the distances  $d_x$  and  $d_y$  as the difference of  $x$  and  $y$  between the first and last  $(x, y)$  pairs – if the image has more pixels in AL ( $N_{\text{AL}} \geq N_{\text{AC}}$ ), this would be the pairs at the lowest and highest AL, respectively.

To calculate the length, we separate into two cases:

1. If  $N_{\text{AL}} \geq N_{\text{AC}}$ , use  $d_y = m \cdot d_x$ . Then

$$l = \sqrt{d_x^2 + d_y^2} = \sqrt{m^2 + 1} \cdot |d_x|, \quad (23)$$

$$\sigma_l^2 = (m^2 + 1) \cdot 2\sigma_x^2 + \frac{m^2 d_x^2}{m^2 + 1} \cdot \sigma_m^2, \quad (24)$$

using Gaussian error propagation for  $\sigma_l$ .

2. If  $N_{\text{AL}} < N_{\text{AC}}$ , use  $d_x = m^{-1} \cdot d_y$ . Then

$$l = \sqrt{d_x^2 + d_y^2} = \sqrt{m^{-2} + 1} \cdot |d_y|, \quad (25)$$

$$\sigma_l^2 = (m^{-2} + 1) \cdot 2\sigma_y^2 + \frac{d_y^2}{m^4 + m^6} \cdot \sigma_m^2, \quad (26)$$

using Gaussian error propagation for  $\sigma_l$ .

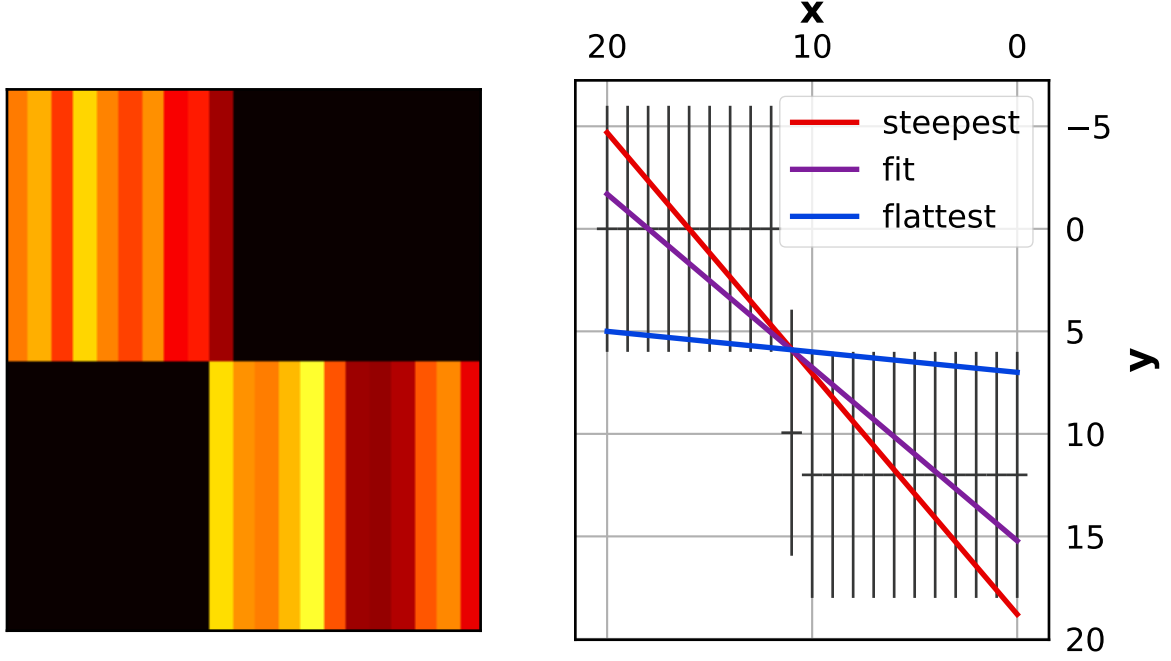


Figure 21: Issues with fitting tracks whose length in AC or AL is 2. *Left*: Track image, with the pixels scaled to their physical dimension. *Right*: Best line fit and the steepest and flattest plausible associated tracks.

the separation into two cases was chosen to avoid numerical instabilities, since  $d_y$  and  $m$  are very small in case 1, while  $d_x$  is very small in case 2, and  $m$  very large.

This fitting routine was then applied to a subset of all detected tracks in a TrackObs, namely those that

1. are two dimensional ( $N_{AL} > 1$  and  $N_{AC} > 1$ )
2. contain at least 5 pixels with a signal value

in addition, tracks with  $N_{AL} = 2$  or  $N_{AC} = 2$  must also be significantly longer in their other dimension. This is due to the fact that, as Figure 21 outlines, there is a very high uncertainty in the distance crossed along the dimension whose length is two pixels.

To test these analysis routines, we once again used the T.A.R.S. simulator to provide input cosmic rays with known quantities - in this case, the true angle  $\theta$  and projected path length  $l$ . We used the simulator to generate cosmic ray subimages akin to those saved into a track observation. This was done for two cases - one being an SM-like thin chip with a binning of 2 in both AL in AC, the other being a BAM-like thick chip with a binning of 4 in AC and no binning in AL.

In this way, we could predominantly examine the effects of binning and the resulting pixel geometry on the retrieval of the track geometries.

In the case of the SM-like chip, Figure 22 plots the input angle of the source cosmic versus the measured angle for a long simulation run. Within their accuracies, the measured angles represent the input angles.

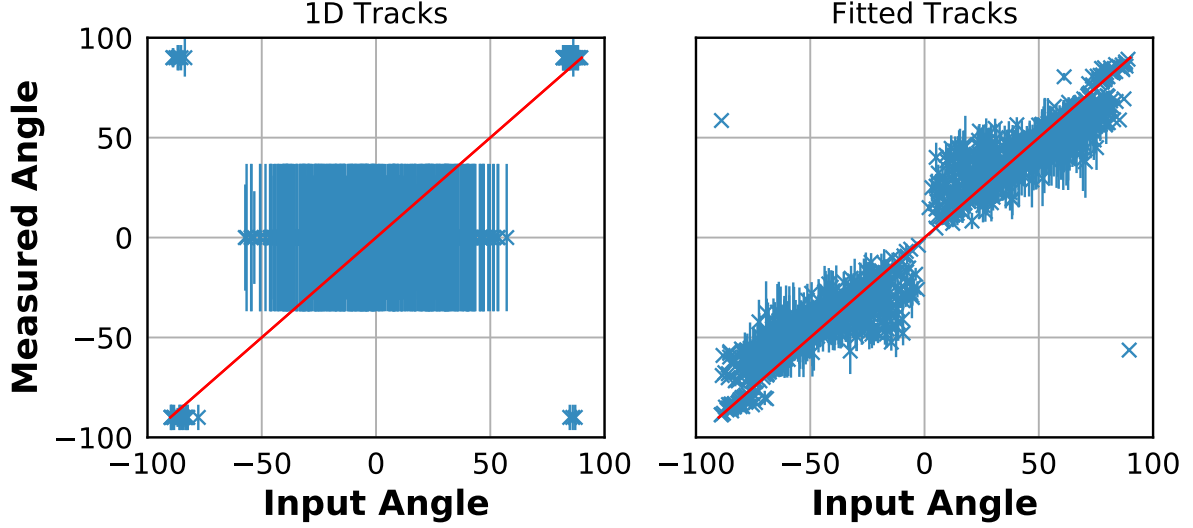


Figure 22: Input angle vs measured angle (in units of degrees) for SM-SIF observations. The error bars represent the uncertainty of the recovered angles, the red diagonal represents the optimal values. The results for one-dimensional and fitted tracks are in the left and right panel, respectively.

This is better illustrated in Fig. 23, which shows histograms of the difference between measured and input angles, normalized by the calculated uncertainties. Here we see that while the fitted cosmics fall well into a Gaussian, the one-dimensional tracks show a bi-modal distribution, where track angles are either systematically under- or over-estimated. Examining Fig. 22, these cosmics appear to be those with very large input angles that are nonetheless seen at  $\theta = 0^\circ$ .

Figure 24 shows similar histograms for the error in the measured track length. These are within the ranges expected by their uncertainty, with a slight systematic overestimation in both indicated by the offset of the maximum from 0.

The same analysis for the BAM-like chip shows different systematics. As Figure 25 shows, a population of the fitted tracks with  $|\theta| \leq 60^\circ$  is redistributed to  $|\theta| \approx \pm 60^\circ$ . This is most likely due to the pixel geometry making these tracks appear to have similar angles.

The error histograms of the measured angles (Fig 26) shows that while the one-dimensional tracks are all correct within the limits of their uncertainty, the fitted tracks contain some outliers with respect to the bell curve - these are the same cosmics that are mistakenly redistributed to higher angles.

The measured lengths, plotted for the BAM-like chip in Figure 27 are once again correctly estimated, with a trend towards an overestimation of the length.

### 3.2.3 Edge Tracks

Users may not want to use tracks that are located at the edges of a readout window, as the detection does not contain the full energy dumped by the cosmic on the CCD, leading to artificially low energies. We mark these tracks, applying a flag whenever one or more edges of the bounding box of a track is located on the edge of the input image.

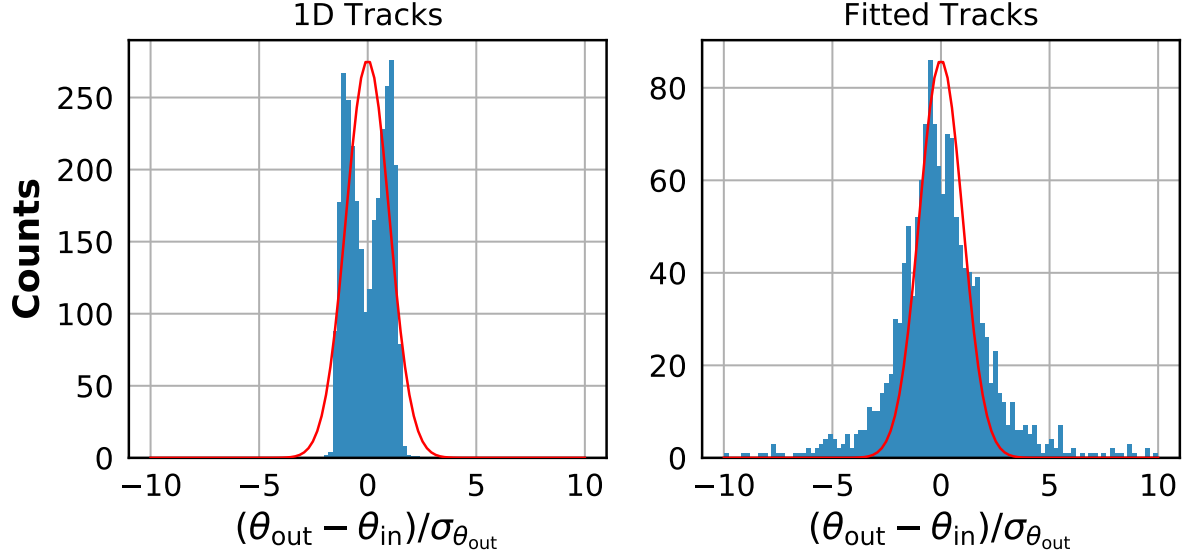


Figure 23: Histograms for the difference of the output angle to the input angle in SM-SIF observations, normalized by the calculated angle uncertainty. The left and right panel show the values for one-dimensional and fitted tracks, respectively. The red curves represent Gaussians with a standard deviation of 1.

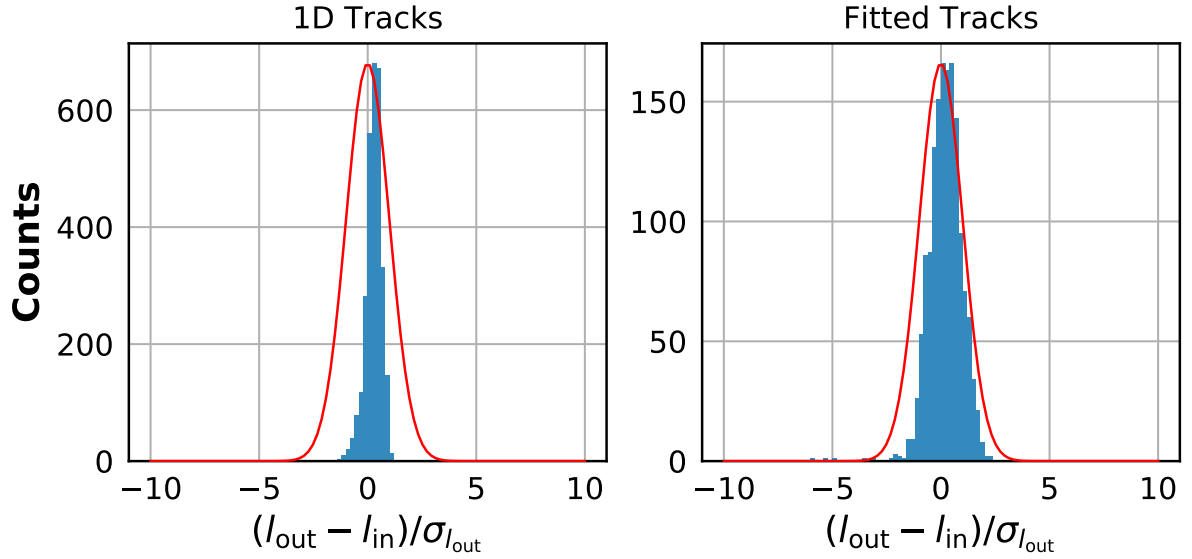


Figure 24: Histograms for the difference of the output track length to the input track length in SM-SIF observations, normalized by the calculated length uncertainty. The left and right panel show the values for one-dimensional and fitted tracks, respectively. The red curves represent Gaussians with a standard deviation of 1.

### 3.2.4 Post-Processing Pipeline

Using the algorithms outlined in this subsection, we process the TrackObs FITS files into FITS files containing the final output product of CosmicObservations, defined in Section 4.1. This

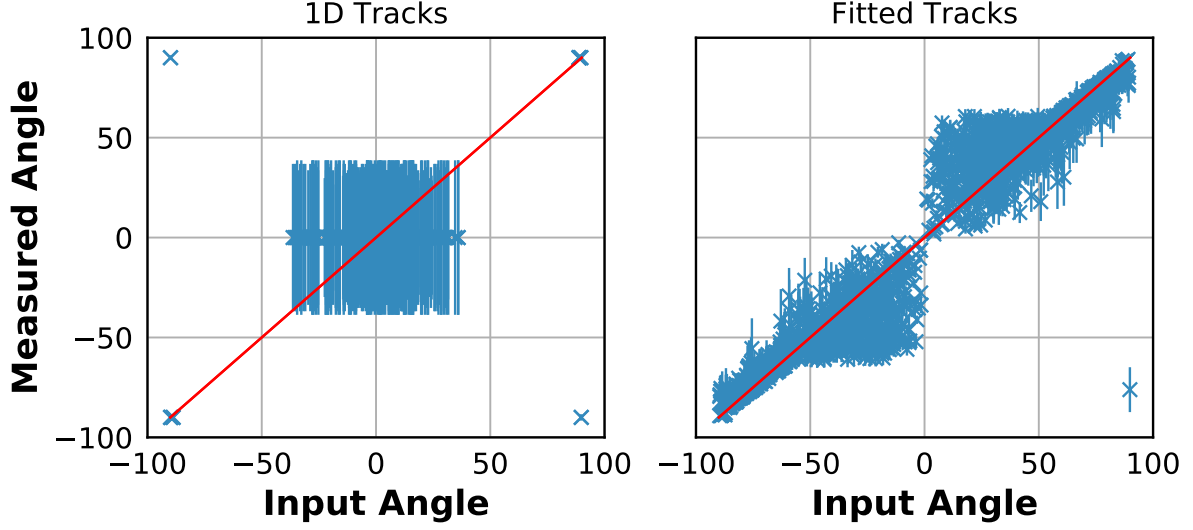


Figure 25: Input angle vs measured angle (in units of degrees) for BAM-OBS and BAM-SIF observations, as in Fig. 22.

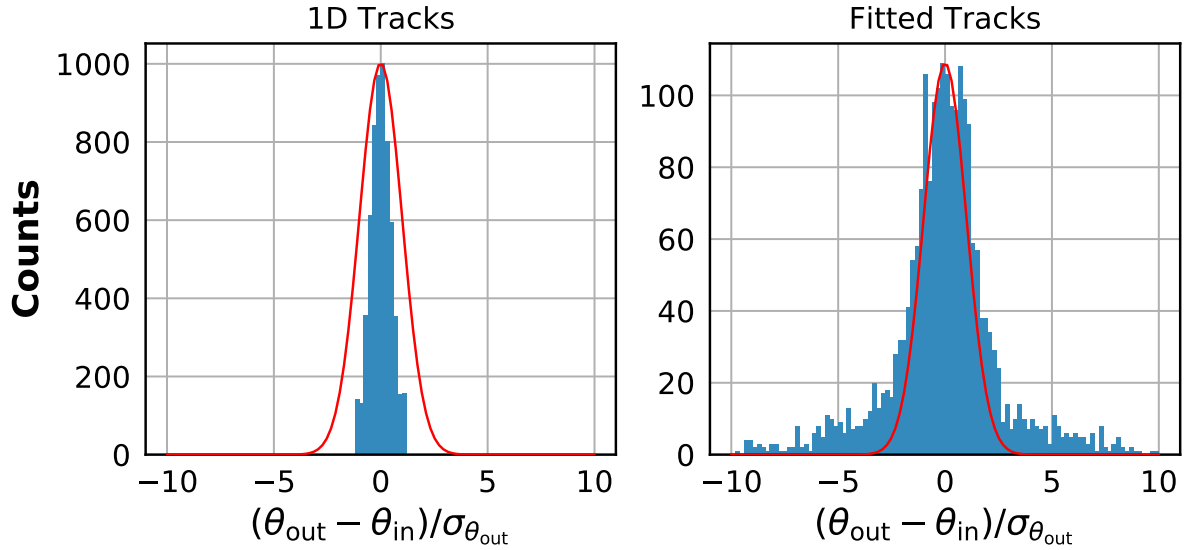


Figure 26: Histograms for the difference of the output angle to the input angle in BAM-OBS and BAM-SIF observations, as in Fig. 23.

dataset is essentially a modification of the TrackObs dataset, adding the calculated fluxes per observation and track geometries to individual tracks while removing the sub-image of the track to reduce file sizes.

Due to this definition of the CosmicObservation in relation to the TrackObs, there is a 1:1 correspondence of each element of the latter dataset to the output dataset. The operations of the post-processing pipeline are:

1. For each TrackObs FITS file:
  - 1.1. Create a buffer for CosmicObservations

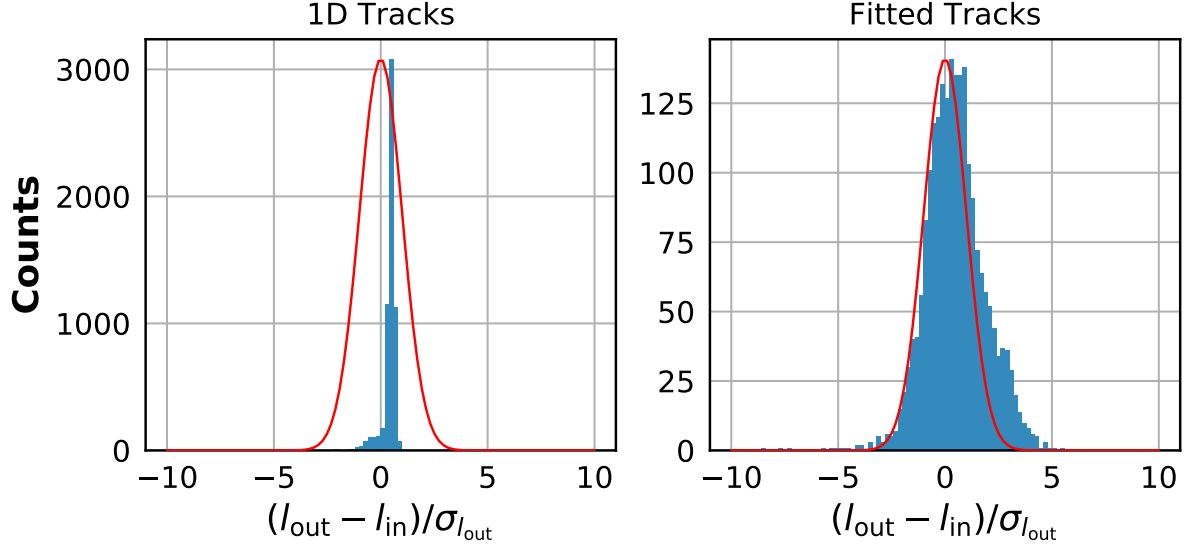


Figure 27: Histograms for the difference of the output track length to the input track length in BAM-OBS and BAM-SIF observations, as in Fig. 24.

- 1.2. For each TrackObs in the FITS file:
  - 1.2.1. Calculate the flux and its uncertainty (Sec. 3.2.1)
  - 1.2.2. Determine the track geometry, where applicable (Sec. 3.2.2)
  - 1.2.3. Determine which tracks are at the edge of the input image (Sec. 3.2.3)
  - 1.2.4. Create a CosmicObservation from the TrackObs and above data and save it to the buffer
- 1.3. Write all the CosmicObservations in the buffer to a FITS file, replicating the directory structure that the TrackObs file has been read from.

As this algorithm acts independently on each input file, it is very easily parallelized. The output directory is then a direct replica of the input TrackObs directory, including subfolders and filenames.

## 4 Data Model

This data model is composed of entities *CosmicObservation* which are described in full detail in the following subsection.

### 4.1 Main Dataset Description

**Entity Name:** *CosmicObservation*

**Entity Description:** A table of cosmic ray (prompt particle event) detections from any of the analyzed Gaia CCDs

**Entity Header:**

<i>Name</i>	<i>Description</i>	<i>Unit</i>	<i>Type</i>
SOURCE	Descriptor of the type of source observation	NA	String
CCD_ROW	CCD row of the source chip	NA	Short
FOV	FOV of the source chip	NA	Short
ACQTIME	Source acquisition time	OBMT	Integer
ACQDATE	Source acquisition time	UTC	String
SRC_AL	AL-dimension of the source image	pixels	Short
SRC_AC	AC-dimension of the source image	pixels	Short
MASKPIX	Number of masked pixels in the source image	pixels	Short
FLUX	Estimated particle flux	particles cm <sup>-2</sup> s <sup>-1</sup>	Float
FLUX_ERR	Poisson Uncertainty of estimated particle flux	particles cm <sup>-2</sup> s <sup>-1</sup>	Float

Table 2: CosmicObservation table header version v1.0

**Entity Columns:**

<i>Name</i>	<i>Description</i>	<i>Unit</i>	<i>Type</i>
LOC_AL	The start position of the event in the AL direction	pixels	Short
LOC_AC	The start position of the event in the AC direction	pixels	Short
DIM_AL	Event track dimension in AL	pixels	Short
DIM_AC	Event track dimension in AC	pixels	Short
TRACK_NPIX	Number of detected pixels	pixels	Short
TRACK_EN	Event track total energy	electrons	Integer
TRACK_EN_ERR	Uncertainty in the event track total energy	electrons	Integer
TRACK_TRUNCATED	Is the event truncated by start/end of CCD physical area	NA	Boolean
GEOMETRY_VALID	Is the track geometry (following rows LEN, THETA) applicable?	NA	Boolean
TRACK_LEN	Estimated incident particle track length	μm	Float
TRACK_LEN_ERR	Uncertainty in the estimated incident particle track length	μm	Float
TRACK_THETA	Estimated incident particle impact angle	deg	Float
TRACK_THETA_ERR	Uncertainty in the estimated incident particle impact angle	deg	Float

Table 3: CosmicObservation table version v1.0

**4.2 Flux Dataset Description**

We provide additional files listing just the measured flux per observation, sorted by the source observation type and chip. A file has been generated for each source type whose extraction has been defined above, as well as for the PPE counters (named SM-PPE).

**Entity Name:** *FluxTable*

**Entity Description:** A table of cosmic ray fluxes from any of the analyzed Gaia CCDs

**Entity Header:**



<i>Name</i>	<i>Description</i>	<i>Unit</i>	<i>Type</i>
SOURCE	Descriptor of the type of the source	NA	String
CCD_ROW	CCD row of the source chip	NA	Short
FOV	FOV of the source chip	NA	Short

Table 4: Flux table header version v1.0

#### Entity Columns:

<i>Name</i>	<i>Description</i>	<i>Unit</i>	<i>Type</i>
COUNTS	Number of particles in the corresponding observation	NA	Integer
FLUX	Estimated particle flux	particles cm <sup>-2</sup> s <sup>-1</sup>	Float
FLUX_ERR	Poisson Uncertainty of estimated particle flux	particles cm <sup>-2</sup> s <sup>-1</sup>	Float
TIME	Source acquisition time	MJD	Float
DATE	Source acquisition time	UTC	String

Table 5: Flux table version v1.0

### 4.3 Auxillary Datasets Description

The auxillary datasets are composed of the spin phase of the spacecraft as well as the location of the spacecraft in a co-rotating geocentric coordinate system.

#### 4.3.1 Spacecraft Spin Phase

The spin phase of the spacecraft yields information on the orientation of the FPA to the sun. It is saved in a FITS file with a single table extension, listing

- A time coordinate in units of MJD
- The spacecraft spin phase in units of degrees

The spin phase is defined such that an angle of 180°, the FPA is facing towards the sun. This is has been further analyzed in the context of cosmic rays in [3].

#### 4.3.2 Spacecraft Location

A low-resolution table of the spacecraft location as a function of time has been calculated and is provided in a fits file.

The coordinate system used is a geocentric ecliptic coordinate system. The X axis points away from the Sun along the Sun-Earth line, while the Y axis is perpendicular to the X axis, pointing in the direction of Earth’s orbital motion. The Z axis is defined to complete the right-handed coordinate system. This defines a co-rotating system to the Sun-Earth system. The definition of X and Y are in this case the opposite of the Geocentric Solar Ecliptic system.

The provided FITS file has a single table extension, listing

- A time coordinate in units of MJD
- The X, Y and Z coordinates of the spacecraft, in units of m.

## 5 Persistence Considerations

### 5.1 Persistence Method(s)

CosmicObservations, fluxes and auxillary data are currently stored in FITS files, with the CosmicObservations being in a date separated directory structure (following yyyy/mm/dd/). Each day-folder contains BAM-OBS files separated by the FOV, each file containing up to 500 CosmicObservations. If there are SM-SIF or BAM-SIF observations available for this day, they are saved in one file per source (i.e. at most one file for SM-SIF, one for BAM-SIF), each containing all the CosmicObservations for that source, regardless of CCD row and FOV.

The fluxes are saved in a single FITS file per source (BAM-OBS, SM-SIF, ...), each containing the tables of that dataset for each CCD row and FOV separately.

The auxillary datasets are saved in one FITS file each.

### 5.2 Size Estimations

Detection algorithms have been developed and written in Python and the default persistence method adopted is TBD. This results in the following estimations.

**Note:** the following estimations are applicable to catalog version : v1.0

The size of the output catalog of CosmicObservations per year can be estimated by extrapolating the size of the datasets for each fully extracted year so far, listed in Table 6.

Year	CosmicObservation data volume [GB]
2014	22
2015	23
2016	26
2017	29

Table 6: Data volume for CosmicObservations per year for catalog version v1.0

It can be clearly seen that the data volume per year is increasing with mission time. This has to do with the fact that the baseline level of galactic cosmic rays in every observation is increasing as the sun is moving toward solar minimum, less effectively shielding the inner solar system.

The data volume is entirely dominated by the BAM-OBS dataset, whose rate of acquisition is constant. While SM-SIF observations are observed to be more frequent as the mission progresses, this only has a minimal effect on the overall catalog size.

Extrapolating from table 6, we assume that the average added data volume per year is 30 GB, assuming that the increasing flux in galactic cosmic rays in further years is offset by the low flux during the early mission. This leads to an estimated catalog size of 30 GB/year, or 2.5 GB/month.

Table 7 lists the size of the datasets for the fluxes and auxillary data, from January 2014 to January 2018. The Flux data groups all datasets (SM-SIF, SM-PPE, BAM-SIF, BAM-OBS) and is dominated by the BAM-OBS and PPE data, whose sampling is constant and the only scaling factor for the catalog size. The spin phase and spacecraft location data sizes are dominated by the sampling at which they are provided.

Dataset	Data volume Jan 2014 - Jan 2018 [MB]
Flux	663
Spin Phase	69
Spacecraft Location	6.5

Table 7: Data volume for fluxes and auxillary datasets for catalog version v1.0

Extrapolating from table 7, we find that the flux dataset has an approximate size of 165 MB/year, the spin phase a size of 17 MB/year and the spacecraft location a size of 1.6 MB/year.

## References

- [1] U. Bastian *Reference Systems, Conventions And Notations for Gaia*.
- [2] R. Kohley *Cosmic rays on BAM CCDs data delivery note*.
- [3] A. Abreu *Radiation Environment and Cosmics from Gaia at L2*.
- [4] P. G. van Dokkum *Cosmic-Ray Rejection by Laplacian Edge Detection*.
- [5] C. McCully *Astro-SCRAPPY: The Speedy Cosmic Ray Annihilation Package in Python*, <https://github.com/astropy/astroscrappy>.
- [6] Lionel Garcia et al. *Tools for Astronomical Radiation Simulations*
- [7] P. T. Boggs and J. E. Rogers, *Orthogonal Distance Regression*, in “Statistical analysis of measurement error models and applications: proceedings of the AMS-IMS-SIAM joint summer research conference held June 10-16, 1989,” Contemporary Mathematics, vol. 112, pg. 186, 1990.