**Response to Referee #1**

**The authors would like to thank the referee for her/his review. Below are our responses to the comments brought up by the referee. Referee's comments and our replies are marked in blue and in black, respectively. In italic are the changes made in the manuscript.**

The manuscript describes a PCA based method on real time detection and characterization of atmospheric events. For this, they are applying measured data from three IASI satellites. The manuscript is nicely written and offers interesting application of the PCA on detection of extreme events. However, some clarifications are needed, as described below.

Major comment
The methodology description needs to be improved. Even though majority of the methodology is described in previous study, it would be important to provide here necessary details on the method for replicating the analyses with similar data. e.g. PCA could be described more clearly and it is not clear how to you get GMI and GMA from the PC's. This makes it more difficult to follow the results from the case studies.

Section 3 has been extensively modified and rewritten to improve the PCA methodology description, See the revised Section 3 at the end of this document in Appendix A.

The explanation to obtain the GMI/GMA was clarified in Section 4.1 as follows:

*Each granule contains ~2700 radiance spectra, from which the corresponding IFOV-residuals are computed based on the IASI-PCA method. For each granule, the largest positive and negative residual value for each spectral channel is recorded in two arrays, called hereafter "Granule Maxima" (GMA) and "Granule Minima" (GMI). GMI and GMA are defined as pseudo-residuals of dimension 8461 (the number of radiance channels) and represent the spectral envelope of the statistics of residuals over the granule. Physically, the GMI (GMA) pseudo-residual is associated with reconstruction errors of spectral absorption (emission) lines. Since the method is based on the granule extrema (GMI and GMA), the method is therefore called: IASI-PCA-GE, with GE standing for Granule-Extrema. It is important to note that these pseudo-residuals associated with a granule are different from the individual IFOV-residual associated with each IFOV.*

Specific comments
Abstract: Point out the focus and true novelty of this manuscript in the abstract. Now it sounds more like the introduction

As suggested by the Referee, we added a paragraph pointing out the focus and true novelty of the manuscript in the abstract:

*The method is running continuously, delivering email alerts on a routine basis using the near real time IASI L1C radiance data. It is planned to be used as an online tool for the early and automatic detection of extreme events, which was not done before.*

line 117, Antonelli 2004 is not in the list of references. In addition, with Atkinson 2008 and 2010 they are not the original or the best references of methods for defining the optimal number of components

We thank the referee. Antonelli (2004) was added to the list of references. In addition, this sentence and corresponding references have been consolidated and moved to Section 3.3, see our answer to the corresponding comment below.

Line 173: Define IASI-PCA-GE more clearly
We added the following sentence in Section 4.1:

*Since the method is based on the granule extrema (GMI and GMA), in the following the method is called: IASI-PCA-GE, with GE standing for Granule-Extrema.*

Section 3.3.: As the explained variance is not really increasing after ~25 components, using 150 PC sounds a bit of overfitting. How did you define the number? How many PC would e.g. Scree test or Kaiser criterion suggest?

Section 3.3 has been consolidated to better explain and justify the choice of 150 PC. Additional references have been also added (Hultberg, 2009, Atkinson, 2009). Note however that :
- We chose to provide only references related to the use of PCA on IASI, as there is already a large experience with the PCA on these specific instruments, thanks to the work performed at EUMETSAT for defining and operationnaly implementing the Principal Component Compression for the IASI L1D products, and to the scientific work performed for testing and analysing this processing and the outliers. As already explained in these references, different approaches can be used for the choice of the truncation. One important aspect is that the performance of the reconstruction depends (in a complex and correlated manner) on several choices : the training dataset, the normalisation matrix, the truncation threshold. For these choices, the specific experience gained on IASI is critical.
- At the end the "optimal" choice remains empirical and statistically-based. Sensitivity tests on the different parameters (including the choice of the truncation) is thus a key point. This is now mentioned in the revised manuscript. It is of particular importance in our work, as the objective here is not to perform the best reconstruction of all the measurements, but to detect outliers.

Atkinson, N. C., Ponsard, C., and Hultberg, T.: AAPP enhancements for the EARS-IASI service, Proc. EUMETSAT Meteorological Satellite Conf., Bath, UK, 21–25 September 2009, available at https://www-cdn-int.eumetsat.int/files/2020-04/pdf_conf_p55_s8_39_atkinson_p.pdf, 2009.

Hultberg, T.: IASI Principal Component Compression (IASI PCC) FAQ, March 2009, EUMETSAT technical note, available at https://www.eumetsat.int/media/8306, 2009.

See the revised Section 3.3 at the end of this document in Appendix A.

Lines 501-506: With this high number of observations in the training set, it is not probable that few outliers would affect drastically to the sensitivity of the method. As already the Atkinson papers pointed out, there has been suspicions that PCA might not be the best method for this type of analysis. Have you considered other possible factorization methods like EFA, NMF or PMF discussed e.g. in Isokääntä et al. 2020 (https://doi.org/10.5194/amt-13-2995-2020)?
In addition, have you considered accounting for the geophysical parameter possibly acting as confounding factors in your analysis?

We agree that it was not correct to write that few outliers would affect drastically to the sensitivity of the method. To improve and clarify the discussion, the corresponding sentences have been modified in the revised manuscript. In particular, the main argument explaining the unconclusive results on CO has been added, in agreement with your comment:

*Also, unconclusive results were obtained for CO because its variability is already well captured by a truncated reconstruction due to the high variability of this species, from background conditions (50 ppb) to highly polluted areas (4000 ppb).*

However, other possible methods have not been tested in this work, as the choice of testing PCA analysis of IASI measurements for extreme event detection is at the origin of the presented work. Finally, we considered accounting for the geophysical parameters acting as confounding factors, as it is illustrated for the $HNO_3$ species in the discussion of the The Ubinas case in Section 5.1.1. a sentence was also added in the revised manuscript following your remark :

*Finally, as explained above concerning SO$_2$ and HNO$_3$, the spectral coincidence of some of the intense spectral features of these two species can affect the reconstruction of one when the other one is highly present. In the frame of this study, this is the only identified example of confounding situations (i.e., unusual perturbation in a limited number of channels impacts the reconstruction residual in other channels) leading to false detection.*

Conclusions: point out that the method can be used as online tool for detecting extreme events, as mentioned in the text earlier.

We thank the Referee for his suggestion. We added the following paragraph in Section 6:

*A first version of this method is currently running continuously, delivering email alerts on a routine basis using the near real time IASI L1C radiance data. Although the method is still being tested, it is planned to be used as an online tool for the early and systematic detection of extreme events.*

***Appendix A: Revised Section 3***

***3 The Principal Component Analysis Method***

***3.1 Basic concepts***

*The PCA method for high spectral resolution sounders, such as IASI, is described in Atkinson et al. (2008). This method is well suited to efficiently represent the amount of information contained in the 8641 IASI channels. It relies on the use of a dataset of thousands of spectra representing the full range of atmospheric conditions from which the principal components are calculated, the so-called "training database".*

*One considers an ensemble Y of n IASI radiance spectra $\boldsymbol{y}$ of dimension m (where m is the number of channels and n is the number of observations). Let denote $\boldsymbol{N}^{-1}\overline{\boldsymbol{y}}$ the mean and $\boldsymbol{S}_{\epsilon}$ $(m \times m)$ the covariance of the normalized ensemble of spectra $\boldsymbol{N}^{-1}Y$. $\boldsymbol{N}$ is the noise normalisation matrix and is defined as the square root of $\boldsymbol{S}_{\boldsymbol{y}}(m \times m)$ the instrument noise covariance matrix associated to the IASI spectra.*

*The PCA is based on the eigen decomposition of the matrix $\boldsymbol{S}_{\epsilon}$ :*

$$\boldsymbol{S}_{\epsilon} = \boldsymbol{E} \, \boldsymbol{\Lambda} \, \boldsymbol{E}^{T} \qquad (1)$$

*where $\boldsymbol{E}$ is the matrix m x m of eigenvectors and $\boldsymbol{\Lambda}$ the diagonal matrix of their associated eigenvalues. The representation of a measured spectrum $\boldsymbol{y}$ in the eigenspace $\boldsymbol{E}$ is obtained by:*

$$\boldsymbol{p} = \boldsymbol{E}^{T}\boldsymbol{N}^{-1}(\boldsymbol{y} - \overline{\boldsymbol{y}}) \qquad (2)$$

*$\boldsymbol{p}$ (dimension m) is the vector of the principal component scores.*

*The analysis consists in representing the multidimensional IASI spectra in a lower dimensional space, which accounts for most of the variance seen in the data. This space is spanned by a truncated set of the eigenvectors of the data covariance matrix. By noise-normalizing the spectra prior to the application of the PCA, the ability to fit the data is enhanced by avoiding giving too much weight to variance caused by noise. Giving m\* the number of most significant eigenvectors of $\boldsymbol{S}_{\epsilon}$, one can represent the spectrum in the eigenspace by a truncated vector of principal component scores, $\boldsymbol{p}$\* of rank m\* (m\* < m). $\boldsymbol{p}$\* is thus a compressed representation of $\boldsymbol{y}$. The reconstructed spectrum, $\widetilde{\boldsymbol{y}}$ (dimension m) is given by:*

$$\widetilde{\boldsymbol{y}} = \overline{\boldsymbol{y}} + \boldsymbol{N}\boldsymbol{E}^{*}\boldsymbol{p}^{*} \qquad (3)$$

*where $\boldsymbol{E}^{*}$ is the matrix of the m\* first eigenvectors or principal components. We define the noise normalized residual vector $\boldsymbol{r}$ (dimension m) of the reconstruction by:*

$$\boldsymbol{r} = \boldsymbol{N}^{-1}(\boldsymbol{y} - \widetilde{\boldsymbol{y}}) \qquad (4)$$

*By definition, if m\* is taken equal to m, $\widetilde{\boldsymbol{y}} = \boldsymbol{y}$ and the residual is the null vector. In nominal cases if the truncation rank is carefully chosen, $\boldsymbol{r}$ essentially contains noise. Several techniques exist to estimate m\* in order to keep the essential part of the atmospheric signal and to remove the eigenvectors containing mainly the measurement noise (e.g., Antonelli et al. (2004), Atkinson et al. (2010)).*

*In the following the noise normalized residual, which is calculated for each IASI IFOV, is called IFOV-residual.*

***3.2 Construction of the training database***

*The training set includes spectra observed over different types of atmospheric/surface conditions at different scan angles and for different pixel numbers to ensure that a truncated set of eigenvectors can be adequately used to represent any observed spectrum. Additionally, if the training set is too small, the specific outcome of the random noise will not be sufficiently uncorrelated and uniform, and will therefore have an influence on the computed eigenvectors and eigenvalues. Extensive experience on IASI spectra from EUMETSAT (Hultberg, 2009, https://www.eumetsat.int/media/8306) and additional experiments with different dataset sizes show that a number of about 70000 spectra is a reasonable lower limit. For this study, around 120000 IASI/Metop-A L1C spectra were selected during a full year (which was chosen as a nominal year for avoiding excessive occurrence of extreme events such fires and volcanoes) on the global scale. The database contains spectra associated with a good quality flag in order to only keep reliable data, acquired indifferently during the day and the night, over land and*

sea, and regardless of the cloud cover. For each month of the year 2013 spectra were selected every five days (1, 6, 11, 16, 21 and 26 of each month). To avoid over-representing high latitudes, because of the large swath of IASI (~2200 km) and frequent overpasses over this area with the polar orbiting satellites, the following method was applied:
- between 90 and 75° only one spectrum is selected
- between 75 and 60°, two spectra are selected
- between 60 and 45°, three spectra are selected
- between 45 and 30°, four spectra are selected
- between 30 and 15°, five spectra are selected
- between 15 and 0°, six spectra are selected
To reach a sufficient but reasonable number of IASI spectra/IFOVs (1.3 $10^6$ spectra per day, 4.7 $10^8$ per year), 120000 IFOVs for year 2013 were randomly chosen to represent all atmospheric/surface situations (air masses, land/sea, day/night, clear/cloudy) and acquisition conditions (IASI scan mirror position and pixel number).

### 3.3 Number of eigenvectors

Several techniques exist to estimate m* in order to keep the essential part of the atmospheric signal and to remove the eigenvectors containing mainly the measurement noise. Antonelli et al (2004) define a criterium based on the spectral RMS reconstruction residuals, finding the optimal truncation rank when this value approach the spectral RMS of the instrument noise. Other methods test directly the behavior of the reconstruction score $\sqrt{\frac{1}{m}\sum_{i=1}^{m}r_i^2}$ as a function of the truncation rank, by looking at the second derivative of the reconstruction score as a function of the truncation rank (e.g., Hultberg, 2009) or plot the principal component score ($p$) spatial correlation as a function of eigenvector rank (Atkinson et al., 2009). In this study, the estimation of m* is based on the analysis of the eigenvalues. The eigenvalues (sorted in descending order) quantify the variability explained by the corresponding eigenvectors, and the optimal number of eigenvectors needed to reproduce the signal in the raw radiances can be determined by analyzing their magnitude and behavior. In the present implementation of the PCA method we process the full IASI spectrum and use a simple method for selecting the truncation rank. The plot of the eigenvalues was examined and PCs were selected up to the point where the slope of the curve stabilized. This leads to choose the first 150 eigenvectors as done in Atkinson et al., 2010. Sensitivity tests has been performed to test the impact of using different values (from 120 to 250) on the reconstructed scores obtained on several atmospheric events (fires and volcanoes cases discussed in the next sections) and confirm this value.