

Application of the Complete Data Fusion to the ozone profiles measured by the Copernicus Atmospheric Sentinel missions: a feasibility study

Nicola Zoppetti¹, Simone Ceccherini¹, Bruno Carli¹, Samuele Del Bianco¹, Marco Gai¹, Cecilia Tirelli¹,
5 Flavio Barbara¹, Rossana Dragani², Antti Arola³, Jukka Kujanpää⁴, Jacob C.A. van Peet^{5,6}, Ronald van
der A⁵ and Ugo Cortesi¹

¹ Istituto di Fisica Applicata “Nello Carrara” del Consiglio Nazionale delle Ricerche, Via Madonna del Piano 10, 50019 Sesto Fiorentino, Italy

² European Centre for Medium-Range Weather Forecasts, Shinfield Park, Reading, RG2 9AX, UK

10 ³ Finnish Meteorological Institute, Atmospheric Research Centre of Eastern Finland, P.O.Box 1627, 70211 Kuopio, Finland

⁴ Finnish Meteorological Institute, Space and Earth Observation Centre, P.O. Box 503, FI-00101 Helsinki, Finland

⁵ Royal Netherlands Meteorological Institute, Utrechtseweg 297, 3731 GA De Bilt, The Netherlands

⁶ Vrije Universiteit Amsterdam, Department of Earth Sciences, Amsterdam, The Netherlands

15 *Correspondence to:* Nicola Zoppetti (N.Zoppetti@ifac.cnr.it)

Abstract. The new platforms for Earth observation from space are characterized by measurements made with great spatial and temporal resolution. While this abundance of information makes it possible to detect and study localized phenomena, on the other hand it may be difficult to manage this large amount of data in the study of global and large-scale phenomena.

A particularly significant example is the use by assimilation systems of Level 2 products that represent gas profiles in the
20 atmosphere. The models on which assimilation systems are based are discretized on spatial grids with horizontal dimensions of the order of tens of kilometres in which tens or hundreds of measurements may fall in the near future.

A simple procedure to overcome this problem is to extract a subset of the original measurements but this involves a loss of information; another option is the use of simple averages of the profiles but also this approach has some limitations that will be discussed in the paper. A more refined solution is to resort to the so-called fusion algorithms, capable of compressing the
25 size of the dataset while limiting the information loss. A novel data fusion method, the Complete Data Fusion, was recently developed to merge a-posteriori a set of retrieved products in a single product. In the present paper, the Complete Data Fusion method is applied to ozone profile measurements simulated in the thermal infrared and ultraviolet bands, in a realistic scenario, according to the specifications of the Sentinel 4 and 5 missions of the Copernicus programme. Then the fused products are compared with the input profiles; comparisons show that the output products of data fusion have in general smaller total errors
30 and higher information contents. The comparisons of the fused with the fusing products are presented both at single fusion grid-box scale and with a statistical analysis and the grid box size impact was also evaluated, showing that the Complete Data Fusion method can be used with different grid-box sizes even if this possibility is strictly connected to the natural variability of the considered atmospheric molecule.

1. Introduction

35 In the context of the Copernicus programme (<https://www.copernicus.eu>) of the European Commission, the European Space Agency is responsible for the Space Component consisting of a novel set of Earth Observation (EO) satellite missions for environmental monitoring applications: the Sentinels (<https://sentinel.esa.int/web/sentinel/missions>). Each mission focuses on a specific aspect of EO. In particular, the geostationary mission Sentinel-4 (S4) and the two Low Earth Orbit missions (Sentinel-5p and Sentinel 5 (S5)), referred to as the atmospheric Sentinels, are dedicated to monitoring air quality, stratospheric
40 ozone, ultraviolet surface radiation and climate.

The atmospheric Sentinels will provide an enormous amount of data with unprecedented accuracy and spatio-temporal resolution. In this scenario, a central challenge is to enable a generic data user (for example, an assimilation system) to exploit such a large amount of data.

45 A variety of approaches can serve the purpose to convey in a single product the information associated to remote sensing observations of the vertical distribution of a given atmospheric target from multiple independent sources. Despite the fact that methodologies to combine coincident measurements from vertical sounders developed to a relatively lesser extent compared to similar classes of algorithms applicable to imaging systems, they are of great and increasing importance to respond to the need for full exploitation of data from new satellites, such as the Copernicus Sentinels and contributing missions. Strategies for combined use of multiple atmospheric profile datasets encompass from a posteriori data fusion techniques to synergistic inversion processes (Aires et al., 2012 and references therein; Natraj et al., 2011; Cuesta et al., 2013; Cortesi et al., 2016; Sato et al., 2018), and, in broader terms, might include assimilation systems with their unique capability of gap filling by merging model and experimental data (Lahoz and Schneider, 2014).

55 These three different approaches differ in the accepted inputs and in the involved models. In the synergistic inversion the inputs consist in the radiance observations (Level 1 products) of all the involved measurements and the output profiles are obtained by a simultaneous retrieval of these observations. A posteriori fusion techniques consist in sophisticated averaging processes in which the inputs are profiles (L2 products) retrieved from the single measurements. The assimilation techniques, in their more general implementations, can accept as inputs both radiances and profiles and use the information of the measurements as inputs of an atmospheric model. Each of these strategies clearly implies different advantages and drawbacks, ultimately assessing the cost-to-benefit ratio that drives the selection of the option of choice for the specific case under investigation.

60 In particular Data fusion algorithms, such as the Complete Data Fusion (CDF) (Ceccherini et al., 2015), can be particularly well suited to reduce the data volume that users need to access and handle while retaining the information content of the whole level 2 (L2) products.

The CDF input is any number of L2 profiles retrieved with the optimal estimation technique and characterized by their a-priori information, covariance matrix (CM) and averaging kernel (AK) matrix. The output of the CDF is a single product (also characterized by an a-priori, a CM and AK matrices) which collects all the available information content.

This work is based on the simulated data produced in the context of the Advanced Ultraviolet Radiation and Ozone Retrieval for Applications project (AURORA, Cortesi et al., 2018), funded by the European Commission in the framework of the Horizon 2020 programme. The project regards the sequential application of fusion and assimilation algorithms to ozone profiles simulated according to the specifications of the atmospheric Sentinels.

70 The use of synthetic data allows evaluating the performances of the algorithm also in terms of differences between the products of interest and a reference truth, represented by the atmospheric scenario used in the procedure to simulate the L2 products. On the other hand, the absence of systematic errors in the simulated measurements limits the study to ideal measurement conditions. However, the CDF algorithm intrinsically provides a mechanism to include different kinds of errors into the analysis. For instance, Ceccherini et al. (2018) discussed how interpolation and coincidence errors can be accounted for and Ceccherini et al. (2019) explicitly introduces the treatment of systematic errors.

75 This work is divided in two parts. In the first part, we describe the datasets and methodologies (the L2 simulation procedure and the CDF) used in the present paper and discuss the use of the profiles average as fusion technique. In the second part, the quality of the fused products obtained from L2 profiles that are not perfectly co-located in space and in time is analysed. To account for the geo-temporal differences in the L2 profiles, a coincidence error is added to the fused product error budget. The fused and standard L2 products are compared and assessed in terms of their information content, highlighting the better data quality provided by the fusion. Finally, we also show that the CDF can be applied with different coincidence grid-box sizes, allowing for different compression factors of the Level 2 input data volume.

The application of CDF to L2 products simulated with the characteristics expected from the atmospheric Sentinel 4 and 5 allows to establish the possible benefits in case of real Sentinel data.

85 2. Material and methods

2.1. Atmospheric scenario and ozone climatology

Two basic external sources have been used to generate the database of the standard L2 ozone products used in this work: the ozone climatology and the atmospheric scenario.

The ozone climatology was used as a priori information for both the simulation of L2 products and the calculation of the CDF.

90 The atmospheric scenario represents the true state of the atmosphere and is used for both the simulation of L2 products and the quality assessment of the fused ones.

The ozone climatology was derived from McPeters and Labow (McPeters and Labow, 2012) and directly provides the a priori profile x_a used either in the simulation expressions (Eq (1)), as well as in the fusion (see Eq. (6)). The CM of the a priori profile S_a is obtained setting the diagonal terms equal to the square of the standard deviation of McPeters and Labow climatology

95 where this standard deviation is larger than 20% of the a priori profile and to the square of 20% of the a priori profile otherwise.

The off-diagonal elements are calculated using a correlation length of 6 km. The correlation length is used to reduce oscillations in the simulated profiles and the value of 6 km is typically used for nadir ozone profile retrieval (Liu et al., 2010, Miles et al., 2015). The a priori CM is used in the simulation of the L2 products, and in particular in the expression of the AK matrix (Eq. (2)) and in the CMs of Eq. (4) and (5) of the next paragraph. The a priori CM S_a plays an important role also in the CDF

100 equations (see Eq. (6)).

The atmospheric scenario is taken from the Modern Era-Retrospective analysis for Research and Applications version 2 (MERRA2) reanalysis (Gelaro et al., 2017). The MERRA2 data are provided by the Global Modelling and Assimilation Office (GMAO) at NASA Goddard Space Flight Center. This reanalysis covers the recent time of remotely sensed data, from 1979 through the present. The atmospheric scenario is the source of true profile x_t used in Eq. (1) to synthesize the simulated L2

105 products and represents the main reference for the comparison of the quality of L2 and fused products.

2.2. L2 Product Simulation

The simulation algorithm has been originally formalized in the context of the AURORA project, aiming at an efficient computational process. The L2 retrieved state is simulated on a fixed vertical grid with a 3 km step, by the linear approximation given in Eq. (1):

110

$$\hat{\mathbf{x}} = \mathbf{A}\mathbf{x}_t + (\mathbf{I} - \mathbf{A})\mathbf{x}_a + \boldsymbol{\delta} \quad (1)$$

where \mathbf{x}_t is the true state of the atmosphere represented by the atmospheric scenarios, \mathbf{x}_a is the a priori estimate of the state vector provided by the ozone climatology, $\boldsymbol{\delta}$ is the uncertainty in the retrieved value due to measurement noise, and $\mathbf{A} = \partial\hat{\mathbf{x}}/\partial\mathbf{x}_t$ is the AK matrix (Rodgers, 2000) calculated according to Eq. (2):

115

$$\mathbf{A} = (\mathbf{K}^T\mathbf{S}_y^{-1}\mathbf{K} + \mathbf{S}_a^{-1})^{-1}\mathbf{K}^T\mathbf{S}_y^{-1}\mathbf{K} \quad (2)$$

In Eq. (2), \mathbf{K} is the Jacobian matrix of the forward model, the superscript T represents the transpose operator, \mathbf{S}_y is the CM of the observations and \mathbf{S}_a is the CM of the a priori profile. The retrieval error $\boldsymbol{\delta}$ is calculated applying the gain matrix \mathbf{G} (Rodgers,

2000) to an error $\boldsymbol{\varepsilon}$ on the observations randomly taken from a Gaussian distribution with average equal to zero and CM given by \mathbf{S}_y :

$$\boldsymbol{\delta} = \mathbf{G}\boldsymbol{\varepsilon} = (\mathbf{K}^T \mathbf{S}_y^{-1} \mathbf{K} + \mathbf{S}_a^{-1})^{-1} \mathbf{K}^T \mathbf{S}_y^{-1} \boldsymbol{\varepsilon} \quad (3)$$

The CM \mathbf{S} associated to the retrieval error $\boldsymbol{\delta}$ (introduced in Eq. (3)) is given by Eq. (4) (Rodgers, 2000):

$$\mathbf{S} = \langle \boldsymbol{\delta} \boldsymbol{\delta}^T \rangle = (\mathbf{K}^T \mathbf{S}_y^{-1} \mathbf{K} + \mathbf{S}_a^{-1})^{-1} \mathbf{K}^T \mathbf{S}_y^{-1} \mathbf{K} (\mathbf{K}^T \mathbf{S}_y^{-1} \mathbf{K} + \mathbf{S}_a^{-1})^{-1} \quad (4)$$

The CM \mathbf{S}_{total} associated to the total error $\boldsymbol{\delta}_{total}$ (that is the difference between the simulated and the true profiles, equal to the random $\boldsymbol{\delta}$ plus the so-called smoothing error caused by the limited vertical resolution of the measurement; see Eq. (7)), is given by Eq. (5) (Rodgers, 2000):

$$\mathbf{S}_{total} = \langle \boldsymbol{\delta}_{total} \boldsymbol{\delta}_{total}^T \rangle = (\mathbf{K}^T \mathbf{S}_y^{-1} \mathbf{K} + \mathbf{S}_a^{-1})^{-1} \quad (5)$$

It should be noted that through the term $\boldsymbol{\delta}$ it is possible to simulate additional error components with respect to the random one considered in this study and this fact adds flexibility to the simulation method.

In this study, the above formulation was used to simulate ozone profiles in the two spectral bands (UV and TIR) for both S4 and S5, after considering the instrument specifications and accounting for the differences in the two spectral bands. In particular, if a fixed geo-location is considered, starting from the same true profile and the same a priori information, the L2 products of the different instruments are obtained by the choice of the suitable Jacobian matrix \mathbf{K} and of the CM \mathbf{S}_y that have been synthesized using the technical requirements of the considered platforms and their foreseen performances.

It is important to remind here that ozone profiles measured in the UV region will be retrieved from spectral radiances acquired by the UVNS/Sentinel-5 spectrometer on-board Meteorological Operational satellite - Second Generation (MetOp-SG) and by the UVN/Sentinel-4 spectrometer on-board Meteosat Third Generation Sounder (MTG-S). For ozone and other targets observed in the TIR, the atmospheric Sentinel missions will be using the operational products of IASI-NG on MetOp-SG and of IRS on MTG.

In the framework of the AURORA project, simulated ozone products from the above-mentioned instruments operating in the UV and TIR regions were generated by using the most up-to-date information available. A detailed description of the instrumental and observational features, as well as of the characteristics of different L2 products goes beyond the scope of this article. All the relevant information was reported in the Technical Note on L2 Data Simulations (AURORA 2017) and can also be found in (Cortesi et al., 2018) and in table 1 of (Tirelli et al 2020), where the spectral bands and the available products specifications used in for the simulation of ozone products in the context of the AURORA project are summarized.

For the sake of shortness, in the next sections of the paper, we will refer to UVNS/MetOp-SG as S5-UV1, to UVN/MTG as S4-UV1, to IASI-NG/MetOp-SG as S5-TIR and to IR/MTG as S4-TIR.

2.3. The CDF method

In this Section, we briefly recall the formulas of the CDF method (Ceccherini et al., 2015). We assume to have N independent simultaneous measurements of the vertical profile of an atmospheric species that can be referred to the same geo-location. Performing the retrieval of the N measurements, we obtain N vectors \boldsymbol{x}_i ($i=1, 2, \dots, N$) providing independent estimates of the

profile, here assumed to be represented on a common vertical grid. Using as inputs these N measurements, the CDF produces as output a single product characterized by a profile \mathbf{x}_f , an AK matrix \mathbf{A}_f and a CM matrix \mathbf{S}_f with the procedure summarized by Eq.s (6). These three quantities are function of the input products, \mathbf{A}_i , \mathbf{S}_i , hereafter referred to as fusing products, and depend on the a priori information $(\mathbf{x}_a, \mathbf{S}_a)$ used as a constraint for the fused product.

$\boldsymbol{\alpha}_i = \hat{\mathbf{x}}_i - (\mathbf{I} - \mathbf{A}_i)\mathbf{x}_{a_i} = \mathbf{A}_i\mathbf{x}_t + \boldsymbol{\delta}_i + \mathbf{A}_i\boldsymbol{\delta}_{coinc,i}$	(6a)
$\tilde{\mathbf{S}}_i = \mathbf{S}_i + \mathbf{A}_i\mathbf{S}_{coinc,i}\mathbf{A}_i^T$	(6b)
$\mathbf{x}_f = \left(\sum_{i=1}^N \mathbf{A}_i^T \tilde{\mathbf{S}}_i^{-1} \mathbf{A}_i + \mathbf{S}_a^{-1} \right)^{-1} \left(\sum_{i=1}^N \mathbf{A}_i^T \tilde{\mathbf{S}}_i^{-1} \boldsymbol{\alpha}_i + \mathbf{S}_a^{-1} \mathbf{x}_a \right)$	(6c)
$\mathbf{A}_f = \left(\sum_{i=1}^N \mathbf{A}_i^T \tilde{\mathbf{S}}_i^{-1} \mathbf{A}_i + \mathbf{S}_a^{-1} \right)^{-1} \sum_{i=1}^N \mathbf{A}_i^T \tilde{\mathbf{S}}_i^{-1} \mathbf{A}_i$	(6d)
$\mathbf{S}_f = \left(\sum_{i=1}^N \mathbf{A}_i^T \tilde{\mathbf{S}}_i^{-1} \mathbf{A}_i + \mathbf{S}_a^{-1} \right)^{-1} \sum_{i=1}^N \mathbf{A}_i^T \tilde{\mathbf{S}}_i^{-1} \mathbf{A}_i \left(\sum_{i=1}^N \mathbf{A}_i^T \tilde{\mathbf{S}}_i^{-1} \mathbf{A}_i + \mathbf{S}_a^{-1} \right)^{-1}$	(6e)
$\mathbf{S}_{f\ total} = \left(\sum_{i=1}^N \mathbf{A}_i^T \tilde{\mathbf{S}}_i^{-1} \mathbf{A}_i + \mathbf{S}_a^{-1} \right)^{-1}$	(6f)

For what concern the profile and the error the CDF can be thought as a “smart average” in which the a priori information is removed from the L2 profiles and CMs before they are put together in the average. The total error of the single product whose a priori information has been removed is higher than the original one and this error increase is only partially compensated by the effect of the average. This is the reason why, even if the total error of the fused product is generally lower than the one of the single L2 fusing product, it is in general higher than the error of the average. The behaviour of the AK matrix is less intuitive and will be thoroughly analysed in the presentation of the results.

A coincidence error characterized by a CM \mathbf{S}_{coinc} is added if the input products are not coincident in time and space. When the CDF is applied to not perfectly coincident products, the diagonal elements of \mathbf{S}_{coinc} are calculated as the square of the 5% of the a priori profile \mathbf{x}_a ; this value has been chosen considering the size of the coincidence grid cells used in this study. The off-diagonal elements of \mathbf{S}_{coinc} are obtained applying an exponential decay with a correlation length of 6 km (Ceccherini et al. 2018). In (Ceccherini et al. 2019) the dynamical choice of \mathbf{S}_{coinc} is presented and in particular the a priori error (coincident with the climatological variability) is used as reference for the diagonal elements and a fixed exponential decay is applied too but the multiplicative factor is calculated by imposing that the cost function of the retrieval is equal to its expectation value. That study, which is based on simulated products similar to the ones of this work, shows that even if the coincidence error is strictly needed for the correct behaviour of the CDF product, this is not strongly dependant by its exact amount until it is smaller with respect to the errors of the individual L2 products.

The formulas of Eqs.(6) refer to the case of measurements made on the same vertical grid. In general, also an interpolation error may be needed considering that the retrievals of the products to be fused can be furnished on different vertical grids. In (Ceccherini et al. 2018) the general expressions of CDF in the case of the fusion of products characterized by different vertical grids are presented and discussed together with the expression of the interpolation error that depends on the involved grids and on the AK matrices of the fusing products. However, since the interpolation error does not apply to the present study (the L2 products have been simulated on the same vertical grid) it has not been considered in Eqs. (6) and in the following discussion.

185 **2.4. Arithmetical average and biases**

Before proceeding, it is necessary to clarify why the arithmetic average of the profiles cannot be considered as a good option to represent a set of products retrieved with optimal estimation techniques.

To do this, we consider N coincident L2 measurements ($i=1,\dots,N$) referring to the same true profile, the same AK matrix and the same CM but having different (noise) errors δ_i randomly generated according to Eq. (3). The total error expression for the
 190 i -th measurement is given in Eq. (7) that can be easily derived from Eq. (1).

$\delta_{i,\text{total}} = \hat{x}_i - x_t = (\mathbf{I} - \mathbf{A}_i)(x_a - x_t) + \delta_i$	(7)
---	-----

Considering that the individual measurements are co-located in space and time, thus they refer to the same truth, the same a priori profile and the same AK matrix \mathbf{A} , the mean total error is equal to:

195

$\langle \delta_{i,\text{total}} \rangle = \langle \hat{x}_i \rangle - x_t = \dots = (\mathbf{I} - \mathbf{A})(x_a - x_t) + \frac{1}{N} \sum_{i=1}^N \delta_i$	(8)
--	-----

It follows that the averaging process reduces the random component of the total error, but does not reduce the bias, due to the a priori information and equal to the term $(\mathbf{I} - \mathbf{A})(x_a - x_t)$ of Eq. (8), which therefore becomes a dominant component. The existence of this bias is one of the reasons why the arithmetic mean cannot be considered as a reference algorithm to collect
 200 the information of several products into one. Further reasons concern the choice of a suitable AK matrix to be assigned to the average (see also von Clarmann. and Glatthor 2019) and the management of possible coincidence and interpolation errors. An explicative comparison of the application of CDF and standard averages in the case of 1000 coincident L2 products is reported in the supplementary material.

3. Results and discussion

205 **3.1. Fusion in realistic spatial and temporal resolution conditions: the L2 Datasets**

To analyse the behaviour of CDF in realistic spatial and temporal resolution conditions four sets of measurements were considered. These measurements correspond to the cloud free observations that were possible between 9:00am and 10:00am on the 1st April 2012. Table 1 lists the L2 product types, namely S4-TIR, S4-UV1, S5-TIR and S5-UV1, used in the remaining of the article. The L2 datasets have been generated according to the equations (1)-(5) described in the paragraph 2.2. The
 210 details of the simulation process can be explored in the technical note (AURORA 2017) considering that here we simulated all the pixels corresponding to a clear sky line of sight in the atmospheric scenario without applying any additional selection criteria. In fact, in the AURORA project 4 months of data have been considered, but a subset of the clear sky pixels has been simulated to reduce the computational cost of the simulations; for this study all the clear sky pixels in the considered hour of data have been simulated, without additional filter, choosing the orbits so that S4-S5 coincidences occur; the spatial distribution
 215 of the simulated products is indirectly represented in the left panel of Figure 5.

3.2. Single grid-box analysis (0.5°x0.625°)

We consider first the case of a single grid-box (Figure 1). In the selected grid-box, 118 measurements were available (55 of S4-TIR, 55 of S4-UV1, 8 of S5 TIR, no S5 UV1). The cell has the size of 0.5 degrees in latitude and 0.625 degrees in longitude,
 220 centred on the Egina Island in the Aegean Sea. The cell size has been chosen to be comparable with the assimilation grid used in the AURORA project. We assign the geo-location of the fused product to be the barycentre of the horizontal coordinates of

the L2 measurements in the grid-box. In this particular case, since the horizontal distribution of the 118 L2 profiles is quite homogeneous, the barycentre is practically placed at the centre of the grid-cell.

225 Figure 2 shows with green lines the absolute (left panel) and relative (right panel) differences between each L2 profile and the corresponding true profile, with a red line the difference between the fused profile and the mean truth (computed as average of the 118 true profiles), with a black dash-dotted line the average of the estimated standard deviation of total error of the individual L2 measurements σ_{total} , and with a red dash-dotted line the estimated standard deviation of the total error of the fused profile σ_{ftotal} . The last two quantities have been calculated as the square root of the diagonals of the \mathbf{S}_{total} and \mathbf{S}_{ftotal} CMs given by Eqs. (5) and (6) respectively. Figure 2 shows that the fused product is in better agreement with its truth than the individual profiles with their own, and presents a smaller estimated total error than the individual L2 products. In particular the right panel allows to see in the detail the performances of CDF in the tropospheric region.

The representation of a retrieved profile is always a compromise between the amplitude of the errors and the vertical resolution. The latter can be quantified by the AKs, which ideally would be equal to the identity matrix in the case of a profile that has a vertical resolution equal to that defined by the sampling grid. Diagonal elements with values smaller than 1 correspond to a loss of vertical resolution. In the left panel of Figure 3, the diagonal elements of the AKs of the L2 products, are compared with the one of the fused product where we have also computed the number of Degrees Of Freedom (DOFs), given by the sum of the diagonal elements of the AK matrix (Rodgers, 2000), for both L2 and fused products, and reported the values in the text box of the left panel; as it can be noted the number of DOFs of the fused product is about twice the number of DOFs of the best L2 one. In the right panel of Figure 3, the comparison of the vertical resolution profiles of L2 and FUS products is shown, where the vertical resolution is calculated starting from AK matrices according the Full Width Half Maximum (FWHM) approach (Rodgers, 2000) and in particular with the algorithm defined in (Ridolfi and Sgheri, 2009).

From the comparison of the left and the right panel of Figure 3 it can be noted that the increase of the AK matrix diagonal values of FUS product, and consequently the increase of the number of DOFs, implies an improved vertical resolution only in a subset of the vertical levels. To better understand the effect of the fusion on the AK matrices, it is useful to analyse the behaviour of their individual rows. In Figure 4, two rows are represented, one that refers to the troposphere (left panel, 6 km) one to the middle stratosphere (right panel, 39 km), where the reference altitude is the one corresponding to the diagonal value of the row. The value of the vertical resolution at the considered altitude is reported in the legend (the minimum of vertical resolution at the considered vertical level for each type of L2 product) while the diagonal value of each row is evidenced in the graphs with cross (L2) and dot markers (FUS). At lower altitudes (left panel), as suggested by one of the reviewers, the DOFs increase can be attributed to three distinct phenomena. The first is the constriction of the main FUS AK lobe and the consequent improvement (of more than 30%) of the vertical resolution with respect to L2 products. The second phenomenon is linked to the fact that, while for the FUS product the maximum value of the AK row corresponds to its diagonal element, for the L2 products these maxima are shifted with respect to the reference altitude of the rows. The last phenomenon is a stronger contribute of the measures with respect to the a priori in the FUS product, where the latter effect can be evidenced considering the sum of all the elements of the rows that assume 0.913 as maximum value for the L2 products and 0.956 for the FUS product. In the particular case all these three effects go in such a direction that can be considered as benefits of CDF application. The results at higher altitudes (39 km, right panel) are primarily influenced by the shape of the AK rows that exhibit large secondary lobes that degrade the vertical resolution.

3.3. Statistical analysis for a large domain

260 While the analysis of the previous paragraph focuses on a particular grid-box, here an analysis of the CDF behaviour is presented, referring to all the 1939 fusion grid-boxes in which more than one of the 79781 L2 simulated products, considered in Table 1 is placed. The fused products can be classified depending on the types of L2 measurements falling inside the coincidence grid cell. Since S4-TIR and S4-UV1 products are in perfect coincidence and S5-UV1 products have a horizontal

spacing larger than the cell size, only six fused product types (FUS type), listed in Table 2, effectively occur. In this table, the FUS type and its description are reported together with the following complementary data:

- Ncells: the number of grid-boxes characterized by the considered FUS type.
- <NL2>: the mean number of individual L2 fusing profiles per grid-box.
- Max NL2: the maximum number of individual L2 fusing products per grid-box.

The left-hand side panel of Figure 5 shows the geographical distribution of the FUS products. Different colours have been used to classify the fused data according to their provenance type. The irregular geographical coverage is due to the realistic distribution of the cloud free measurements. The histogram in the right-hand side panel of Figure 5 shows the number of cells that contain a given number of measurements, divided in different colours depending on the FUS type. The FUS cells, in which only S5 platform L2 products fall, are characterized by a small number of L2 measurements, while when S4 products are present, many L2 measurements can be present.

With the selected grid-box size and the multitude of different products that are present in each cell, the question is which product can be used in alternative to the fusion process in those operations in which a single product is requested in each grid-box. Since the averaging process is affected by a large bias error, a viable alternative is the use of the best fusing product present in the cell and we want to compare the CDF result with this product. This comparison is the so-called Synergy Factor (SF), introduced by Aires et al. (2012). Although Aires introduces SF only for errors (Eq. (11)), we extend his definition also for other quantities because they constitute a useful tool to synthetically represent the performances of fusion algorithms.

The *SF DOF*, defined by Eq. (9), is a pure number that can be calculated for every FUS pixel by the ratio of the number of DOFs of the FUS product and the maximum number of DOFs of the L2 measurements that have been fused. In this equation the index *l* enumerates the vertical levels and the index *i* enumerates the L2 products fused in each grid-box.

$SFDOF = \frac{\sum_l \mathbf{A}_{f,l}}{\max_{i \in L2} \sum_l \mathbf{A}_{i,l}}$	(9)
---	-----

When *SF DOF* is larger than 1.0, the FUS product carries more information than the individual L2 measurements. Figure 6 shows that the *SF DOF* computed for all the fused products (and plotted as a function of the number of L2 profiles in each grid-box) is always larger than 1.0. This means that the information content of the fused product is always larger than that of the standard L2 retrievals. It is also worth noticing that *SF DOF* increases approximately linearly with the logarithm of the number of fusing products, although the proportionality depends on the FUS type. The two different clusters of red symbols (S4:TIR+UV1) are caused by the different latitude bands in which these products are distributed (see also left panel of Figure 5). It is important to underline that the improvement in vertical resolution, which cannot be obtained with the arithmetic averaging, is the most demanding requirement (in terms of observation time and instrument sensitivity) in remote sensing observations and, considering the significant gain obtained relative to the single product selection, is the most important feature of fused products.

While *SF DOF* is a scalar quantity, both *SF AK* and *SF ERR*, defined by Eqs. (10) and (11), are vertical profiles of pure numbers. *SF AK* represents an expansion on the vertical dimension of *SF DOF* and, in particular, is calculated, level by level, as the ratio between the diagonal elements of the AK matrix of the FUS product and the maximum of the corresponding elements of the AK matrices of the fusing L2 measurements.

A value of *SF AK* larger than 1.0 at a specific vertical level (indicated by the index *l*) means that, at that level, the diagonal value of the AK matrix of the FUS product has a larger value than that of all the individual products. As we have seen commenting Figure 3 and Figure 4, the increase of the AK diagonal values at a specific level can happen for different reasons but all of them can be considered as an improvement in the product quality.

$$SFAK_l = \frac{A_{f,l}}{\max_{i \in L2} A_{i,l}} \quad (10)$$

305

The *SF ERR* (Eq. (11)) at a given level is calculated as the ratio between the minimum total error of the L2 measurements that have been fused and the total error of the FUS product. A value of *SF ERR* larger than 1.0 means that at a specific level the error of the FUS product is smaller than that of all the individual products.

$$SFERR_l = \frac{\min_{i \in L2} \sigma_{total,i,l}}{\sigma_{total,l}} \quad (11)$$

310

The SFs defined by Eqs. (10) and (11) provide a conservative comparison because the fused product is compared with the L2 product that at that level has the largest diagonal value in its AK matrix and with the one that has the smallest total error at the same level (generally, these are two distinct L2 products).

Figure 7 shows the *SF AK* (left panel) and *SF ERR* (right panel) profiles for the 1939 FUS products considered in Table 2.

315

We have used different colours to denote the provenance of the L2 data contributing to the fused products and different symbol size to infer the number of L2 fusing measurements (the larger the symbol size, the larger the number of L2 fusing profiles).

As mentioned above, the merit of the fused product in terms of SF is higher than that of the L2 retrievals if the *SF AK* and *SF ERR* are greater than 1. The significant improvement obtained with the fused products is confirmed by Figure 6. It can also be

320

noticed that considering symbols of the same colour the symbol size (*N*) tends to increase moving horizontally in the graph (same vertical level) from left to right (*SF* increasing) denoting that, for each FUS type, *SF* increases with *N*. This is not in

contradiction with the fact that symbols with different colours (FUS types) and different sizes (*N*) can share the same position (*SF*, vertical level) on the graph. Some *SF AK* values, both in the troposphere and in middle upper atmosphere are smaller

325

than one; in the troposphere this happens in 20 cells out of 1939 while in the middle upper atmosphere this happens in almost 500 cells for two possible and sometimes simultaneous circumstances. The first one happens when the introduction of the

coincidence error provokes a sensible degradation of the quality of the FUS AKM. The second circumstance happens for example when one of the L2 products is characterized by a vertical resolution that is much better than all the other fusing products and in particular the peaks of their AKM rows tends to not coincide with the nominal vertical level of the row itself.

3.4. Statistical analysis on a coarse horizontal resolution

330

We have seen that, starting from 79781 L2 measurements (Table 1), when a coincidence grid-box with size 0.5°x0.625° is used, the number of fused profiles is 1939 (Table 2), with a reduction of the data volume of more than a factor 40.

Table 3 provides a summary of the number of fused profiles and the provenance of the L2 profiles that contribute to them for a fusion grid resolution of 1°x1°. In this case, the total number of FUS products is 775 with a reduction of data volume of more than a factor 100.

335

The Synergy factors *SF DOF*, *SF AK* and *SF ERR* have been considered also in this case and the figures (similar to Figure 6 and Figure 7) are reported in the supplementary material. In summary, the greater number of fusing observations in each fusion

cell produces a further improvement for both the vertical resolution and the total error proving that the CDF method can be used with a wide range of grid-box size and data compression and the quality of the products generally improves with larger

cells. An upper limit to the grid-box size is caused by the coincidence error amplitude, which increases with the geographical

340

variability degrading the quality of the fused product; the study of this aspect will be of crucial importance if the CDF will be applied to species with greater spatial-temporal variability than ozone or in any case to very large spatial-temporal domains.

4. Conclusions

This paper presents a sensitivity study of the Complete Data Fusion technique, applied to L2 measurements simulated with the characteristics of the measurements, which will be acquired in the context of the atmospheric Sentinel missions. This analysis allows to evaluate the performances of the CDF algorithm in ideal conditions (i.e., with no systematic errors added) and to quantify the possible benefits of the application of CDF to real Sentinel data.

In particular, we show the application of CDF to a single cell with size of 0.5 degrees in latitude and 0.625 degrees in longitude in which more than 100 L2 products are fused. Results show that the fused product is characterized by higher information content, smaller errors and smaller residuals (i.e., anomaly from the true profile) compared to individual L2 products. The information content being, with its improvement of the vertical resolution, the most important achievement.

This analysis is then extended to a larger domain consisting in 79781 L2 products subdivided in 1939 grid boxes with 0.5°x0.625° size. In this case, the comparison of L2 products and CDF output are carried on in terms of synergy factors. This analysis shows that the CDF can be applied to a wide range of situations and that the benefits of the fusion strongly depend on the number of the measurements that are fused together and from their characteristics. It is also shown that CDF can be run customizing grid resolutions, e.g. to match the resolution requirements of the process that will ingest the products, with full exploitation of all the available measurements.

As the fused products are traced back to a regular, fixed horizontal grid and, as shown here, are not affected by the bias introduced by the a priori information, they can be considered as a new type of Level 3 products with improved quality (reduced bias) and the same characteristics (AK included) with respect to L2 products, even if further analysis are needed especially for what concern the coincidence error to be applied fusing data on large spatial-temporal domains.

Data availability

The data of the simulations presented in the paper are available from the authors upon request.

MERRA-2 data (atmospheric scenario) are available at MDISC (<https://disc.gsfc.nasa.gov>), managed by the NASA Goddard Earth Sciences (GES) Data and Information Services Center (DISC).

The ML climatology (McPeters and Labow, 2012) is available online from the Goddard anonymous ftp account: <ftp://toms.gsfc.nasa.gov>.

Author contributions (according to CRediT <https://casrai.org/credit/>)

N. Zoppetti: Conceptualization, Methodology, Software, Writing – Original Draft, Writing – Review & Editing, Investigation, Data curation, Visualization S. Ceccherini: Conceptualization, Methodology, Investigation, Writing – Review & Editing B. Carli: Conceptualization, Methodology Writing – Review & Editing, Supervision S. Del Bianco: Investigation, Data curation, Project Administration M. Gai: Investigation, Data curation C. Tirelli: Investigation, Data curation, Project Administration F. Barbara: Resources R. Dragani: Investigation, Data curation, Writing – Review & Editing A. Arola: Investigation, Data curation J. Kujanpää: Investigation, Data curation R. Van Der A: Investigation, Data curation U. Cortesi: Project Administration, Supervision, Writing – review & editing.

Competing interests.

The authors declare that they have no conflict of interest.

Acknowledgments

The results presented in this paper arise from research activities conducted in the framework of the AURORA project (<http://www.aurora-copernicus.eu/>) supported by the Horizon 2020 research and innovation programme of the European Union (Call: H2020-EO-2015; Topic: EO-2-2015) under Grant Agreement N. 687428.

Financial support.

This research has been supported by the European Commission, H2020 (AURORA, grant no. 687428).

References

- Aires, F., Aznay, O., Prigent, C., Paul, M. and Bernardo, F.: Synergistic multi-wavelength remote sensing versus a posterior combination of retrieved products: Application for the retrieval of atmospheric profiles using MetOp-A. *J GEOPHYS RES*, Vol. 117, D18304, <https://doi.org/10.1029/2011JD017188>, 2012.
- AURORA consortium, (Advanced Ultraviolet Radiation and Ozone Retrieval For Applications, grant no. 687428): Technical Note On L2 Data Simulations [D3.4], 35 pp., available for download at <https://cordis.europa.eu/project/id/687428/results>, 2017.
- Ceccherini, S., Carli, B., and Raspollini, P.: Equivalence of data fusion and simultaneous retrieval, 2015. *Opt. Express*, 23, 8476-8488, <https://doi.org/10.1364/OE.23.008476>, 2015.
- Ceccherini, S., Carli, B., Tirelli, C., Zoppetti, N., Del Bianco, S., Cortesi, U., Kujanpää, J., and Dragani, R.: Importance of interpolation and coincidence errors in data fusion. *Atmos. Meas. Tech.*, 11, 1009–1017, <https://doi.org/10.5194/amt-11-1009-2018>, 2018.
- Ceccherini S., Zoppetti N., Carli B., Cortesi U., Del Bianco S., and Tirelli C.: The cost function of the data fusion process and its application. *Atmos. Meas. Tech.*, 12, 2967–2977, <https://doi.org/10.5194/amt-12-2967-2019>, 2019.
- Cortesi, U., S. Del Bianco, S. Ceccherini, M. Gai, B.M. Dinelli, E. Castelli, H. Oelhaf, W. Woiwode, M. Höpfner, D. Gerber, Synergy between middle infrared and millimeter-wave limb sounding of atmospheric temperature and minor constituents, *Atmos. Meas. Tech.*, 9, 2267-2289, <https://doi.org/10.5194/amt-9-2267-2016>, 2016.
- Cortesi, U., Ceccherini, S., Del Bianco, S., Gai, M., Tirelli, C., Zoppetti, N., Barbara, F., Bonazountas, M., Argyridis, A., Bós, A., Loenen, E., Arola, A., Kujanpää, J., Lipponen, A., Nyamsi, W.W., van der A, R., van Peet, J., Tuinder, O., Farruggia, V., Masini, A., Simeone, E., Dragani, R., Keppens, A., Lambert, J.-C., van Roozendaal, M., Lerot, C., Yu, H., and Verberne, K.: Advanced Ultraviolet Radiation and Ozone Retrieval for Applications (AURORA): A Project Overview, *Atmosphere*, 9, 454, <https://doi.org/10.3390/atmos9110454>, 2018.
- Cuesta, J., M. Eremenko, X. Liu, G. Dufour, Z. Cai, M. Höpfner, T. von Clarmann, P. Sellitto, G. Foret, B. Gaubert, M. Beekmann, J. Orphal, K. Chance, R. Spurr and J. M. Flaud: Satellite observation of lowermost tropospheric ozone by multispectral synergism of IASI thermal infrared and GOME-2 ultraviolet measurements over Europe, *Atmos. Chem. Phys.*, 13 (19), pp.9675-9693, 2013.

Gelaro, R., McCarty, W., Max J. Suárez, M. J., Todling, R., Molod, A., Takacs, L., Randles, C. A., Darmenov, A., Bosilovich, M. G., Reichle, R., Wargan, K., Coy, L., Cullather, R., Draper, C., Akella, S., Buchard, V., Conaty, A., da Silva, A. M., Gu, W., Kim, G. K., Koster, R., Lucchesi, R., Merkova, D., Nielsen, J. E., Partyka, G., Pawson, S., Putman, W., Rienecker, M., Schubert, S. D., Sienkiewicz, M., and Zhao, B. The Modern-Era Retrospective Analysis for Research and Applications, Version 2 (MERRA-2), *J. Climate*, 30, 5419-5454, <https://doi.org/10.1175/JCLI-D-16-0758.1>, 2017.

Lahoz, W.A. and Schneider, P., Data assimilation: making sense of Earth Observation, *Frontiers in Environmental Science*, 2, <https://doi.org/10.3389/fenvs.2014.00016>, 2014.

Liu, X., Bhartia, P. K., Chance, K., Spurr, R. J. D., and Kurosu, T. P.: Ozone profile retrievals from the Ozone Monitoring Instrument, *Atmos. Chem. Phys.*, 10, 2521-2537, <https://doi.org/10.5194/acp-10-2521-2010>, 2010.

McPeters, R.D., and Labow, G.J.: Climatology 2011: An MLS and sonde derived ozone climatology for satellite retrieval algorithms, *J. Geophys. Res.*, 117, D10303, <https://doi.org/10.1029/2011JD017006>, 2012.

Miles, G. M., Siddans, R., Kerridge, B. J., Latter, B. G., and Richards, N. A. D.: Tropospheric ozone and ozone profiles retrieved from GOME-2 and their validation, *Atmos. Meas. Tech.*, 8, 385-398, <https://doi.org/10.5194/amt-8-385-2>, 2015.

Natraj, V., Liu, X., Kulawik, S., Chance, K., Chatfield, R., Edwards, D. P., Eldering, A., Francis, G., Kurosu, T., Pickering, K., Spurr, R., and Worden, H.: Multi-spectral sensitivity studies for the retrieval of tropospheric and lowermost tropospheric ozone from simulated clear-sky GEO-CAPE measurements, *Atmos. Environ.*, 45, 7151–7165, 2011

Ridolfi, M. and Sgheri, L. 2009: A self-adapting and altitude-dependent regularization method for atmospheric profile retrievals, *Atmos. Chem. Phys.*, 9, 1883–1897, 2009 <https://doi.org/10.5194/acp-9-1883-2009>

Rodgers, C.D.: *Inverse Methods for Atmospheric Sounding: Theory and Practice*. Vol. 2 of Series on Atmospheric, Oceanic and Planetary Physics. World Scientific: Singapore, 2000.

Sato, T. O., Sato, T. M., Sagawa, H., Noguchi, K., Saitoh, N., Irie, H., Kita, K., Mahani, M. E., Zetsu, K., Imasu, R., Hayashida, S., and Kasai, Y.: Vertical profile of tropospheric ozone derived from synergetic retrieval using three different wavelength ranges, UV, IR, and microwave: sensitivity study for satellite observation, *Atmos. Meas. Tech.*, 11, 1653–1668, <https://doi.org/10.5194/amt-11-1653-2018>, 2018.

Tirelli, C., Ceccherini, S.; Zoppetti, N., Del Bianco, S., Gai, M., Barbara, F., Cortesi, U., Kujanpää, J., Huan, Y., Dragani, R. Data fusion analysis of Sentinel-4 and Sentinel-5 simulated ozone data. *J. Atmos. Ocean. Technol.*, 37 (4), 573–587, <https://doi.org/10.1175/JTECH-D-19-0063.1>, 2020.

von Clarmann, T. and Glatthor, N.: The application of mean averaging kernels to mean trace gas distributions, *Atmos. Meas. Tech.*, 12, 5155–5160, <https://doi.org/10.5194/amt-12-5155-2019>, 2019.

L2 Type	Platform	Band	Number of simulated measurements	Minimal distance between measurements across x along track [km]
S4-TIR	S4	TIR	35594	5.7 x 7.4
S4-UV1	S4	UV1	35594	
S5-TIR	S5	TIR	8023	12.2 x 12.3
S5-UV1	S5	UV1	570	46.2 x 46.7
TOTAL			79781	

Table 1: Characteristics of the simulated measurements. For S4 platform across-track is South-North direction and along-track is East-West direction.

FUS Type	Description	<i>N_{cells}</i>	<i><NL₂></i>	<i>max NL₂</i>
S4:TIR+UV1	Two or more S4 pixels, no S5 pixels.	908	29.3	160
S4:TIR+UV1_S5:TIR+UV1	Two or more S4 pixels, one or more S5_TIR pixel, one or more S5_UV1 pixel.	245	114.7	163
S4:TIR+UV1_S5:TIR	Two or more S4 pixels, one or more S5_TIR pixel, no S5_UV1 pixels.	299	69.4	165
S4:TIR+UV1_S5:UV1	Two or more S4 pixels, one or more S5_UV1 pixel, no S5_TIR pixels.	2	20	37
S5:TIR+UV1	No S4 pixels, one or more S5_TIR pixels, one or more S5_UV1 pixels.	247	11.1	24
S5:TIR	No S4 pixels, two or more S5_TIR pixels, no S5_UV1 pixels.	238	6.2	14
TOTAL		1939	41.1	165

Table 2: types and characteristics of fused product when a coincidence grid cell size of 0.5°x0.625° is used. N_{cells} is the number of grid-boxes characterized by the considered FUS type; <NL₂> is the mean number of individual L2 fusing profiles per grid-box and Max NL₂ is the maximum number of individual L2 fusing products per grid-box.

460

FUS Type	Description	<i>N_{cells}</i>	<i><NL₂></i>	<i>max NL₂</i>
S4:TIR+UV1	Two or more S4 pixels, no S5 pixels.	354	73.1	420
S4:TIR+UV1_S5:TIR+UV1	Two or more S4 pixels, one or more S5_TIR pixel, one or more S5_UV1 pixel.	140	289.4	504
S4:TIR+UV1_S5:TIR	Two or more S4 pixels, one or more S5_TIR pixel, no S5_UV1 pixels.	79	115.4	442
S4:TIR+UV1_S5:UV1	Two or more S4 pixels, one or more S5_UV1 pixel, no S5_TIR pixels.	0	0	0
S5:TIR+UV1	No S4 pixels, one or more S5_TIR pixels, one or more S5_UV1 pixels.	142	26.2	71
S5:TIR	No S4 pixels, two or more S5_TIR pixels, no S5_UV1 pixels.	60	8.9	26
TOTAL		775	102.9	504

Table 3: Like in Table2 but with a grid-box size of 1°x1°. N_{cells} is the number of grid-boxes characterized by the considered FUS type; <NL₂> is the the mean number of individual L2 fusing profiles per grid-box and Max NL₂ is the maximum number of individual L2 fusing products per grid-box.

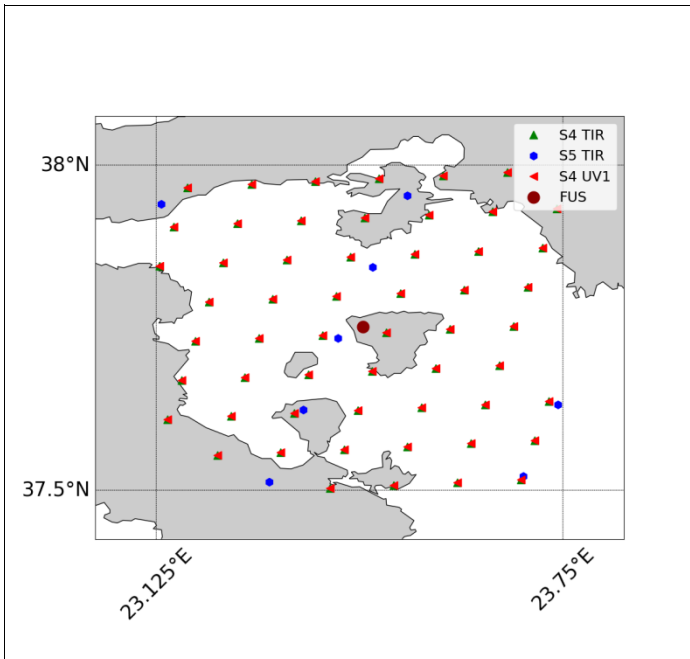


Figure 1: geographical distribution of the simulated L2 measurements and geo-location of the fused product. The black dash-dotted lines represent the borders of the $0.5^\circ \times 0.625^\circ$ grid cells.

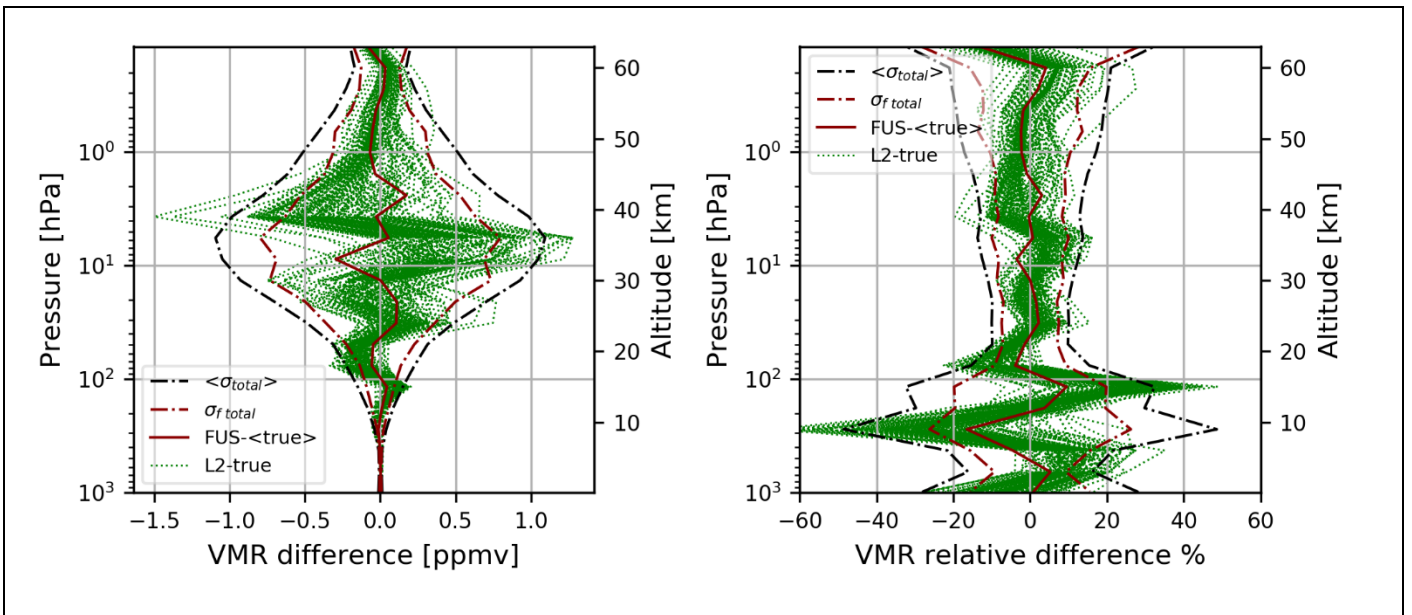


Figure 2 (Left panel): absolute differences between L2 profiles and their true profiles (green lines), absolute difference between the fused profile and the average of the true profiles (dark red continuous line), the average of σ_{total} of L2 simulations (black dash-dotted lines), $\sigma_{f total}$ (dark red dash-dotted lines). (Right panel): relative percentage differences between L2 profiles and their true profiles (green lines), relative percentage difference between the fused profile and the average of the true profiles (dark red continuous line), the average of σ_{total} of L2 simulations normalized wrt the true profile and expressed in percentage (black dash-dotted lines), $\sigma_{f total}$ normalized wrt the true profile and expressed in percentage (dark red dash-dotted lines).

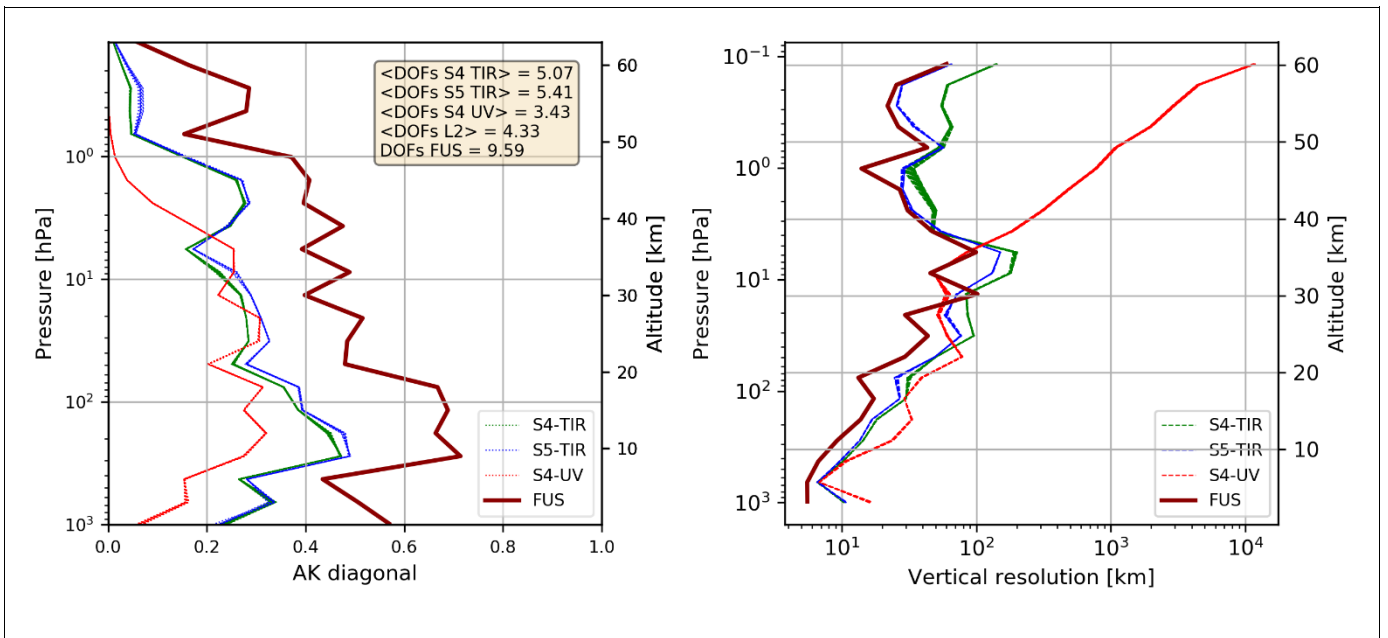


Figure.3 (Left panel): AKs diagonal of: S4-TIR products (red lines), S5-TIR products (blue lines) S4-UV products (red lines) and FUS product (dark red line). In the text box the average number of DOFs for each type of L2 product, the average number of DOFs for all L2 products and the number of DOFs of the FUS product are reported. (Right panel): Vertical resolution (FWHM) profiles of: S4-TIR products (red lines), S5-TIR products (blue lines) S4-UV products (red lines) and FUS product (dark red line). In each panel, while solid dark red line is a single one, red and green lines are both 55 overlapped lines and blue lines are 8 overlapped lines (one for each L2 product).

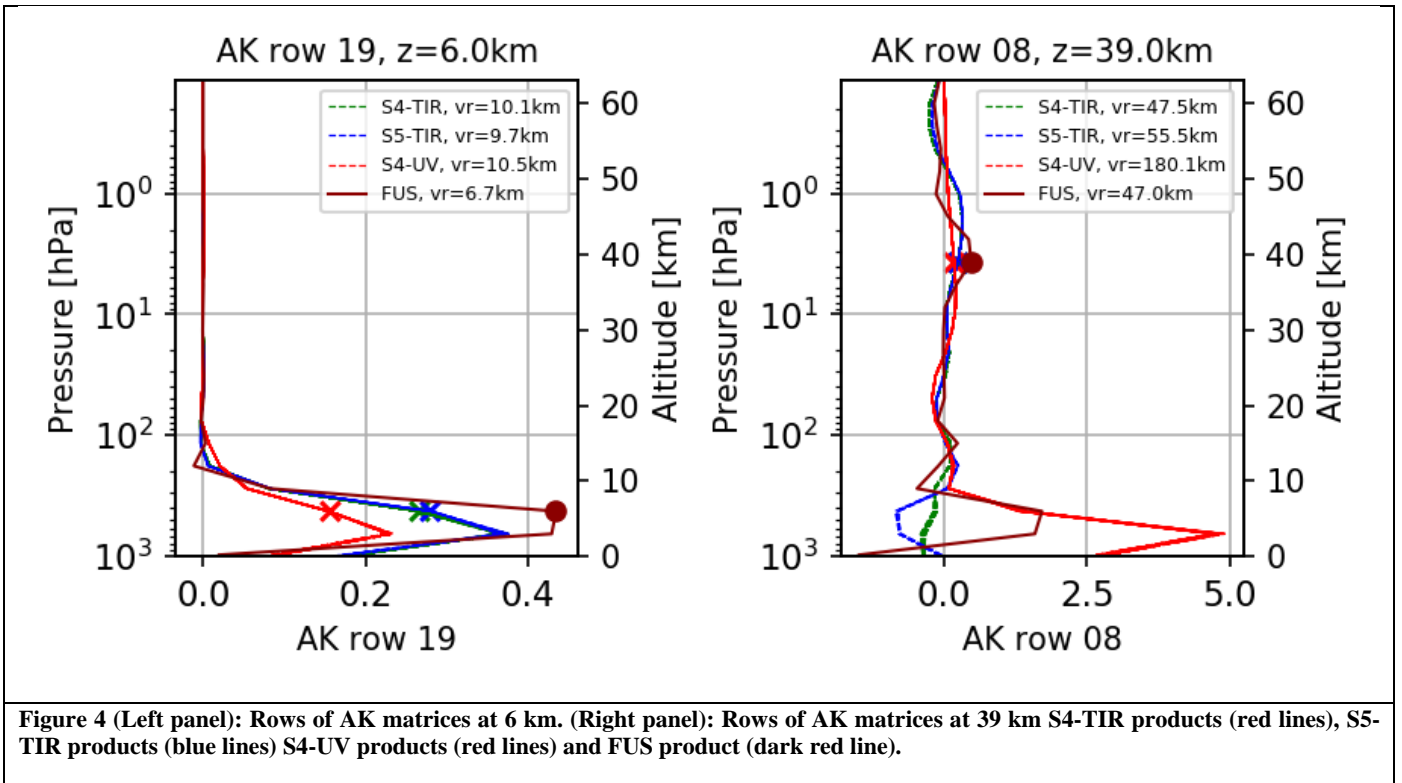


Figure 4 (Left panel): Rows of AK matrices at 6 km. (Right panel): Rows of AK matrices at 39 km S4-TIR products (red lines), S5-TIR products (blue lines) S4-UV products (red lines) and FUS product (dark red line).

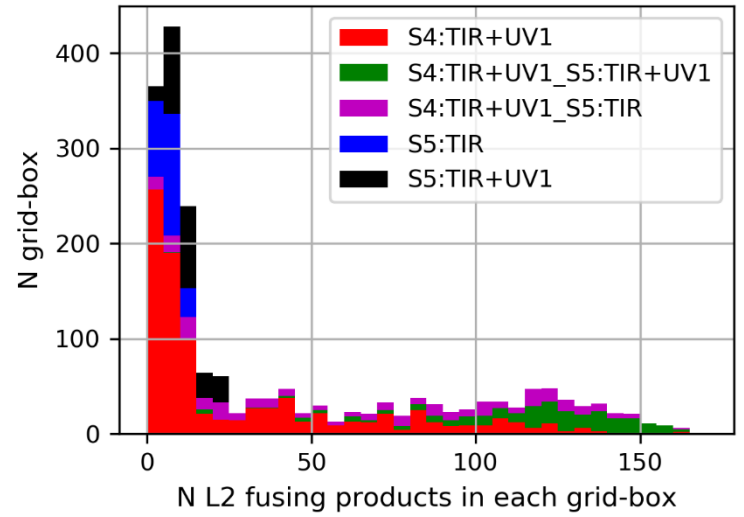
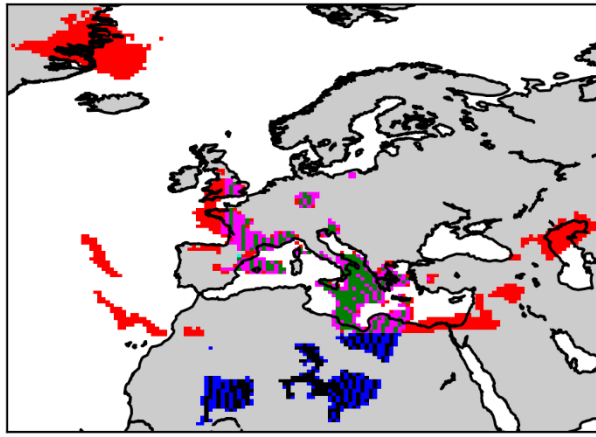


Figure 5. Left panel: geographical distribution of FUS products differentiated by FUS type where the effect of the lower resolution of S5-UV1 respect to the other L2 products is the cause of the periodic FUS type transitions in the Mediterranean area. Right panel: histogram of the number of cells with a given number of L2 measurements differentiated by FUS type.

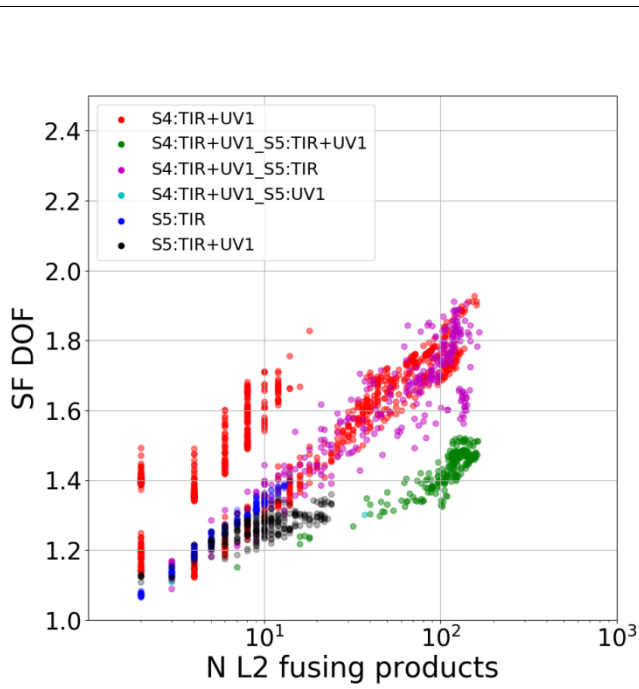


Figure 6: scatter plot of SF DOF as a function of the number of L2 measurements fused in each coincidence grid cell; different colours represent different FUS types.

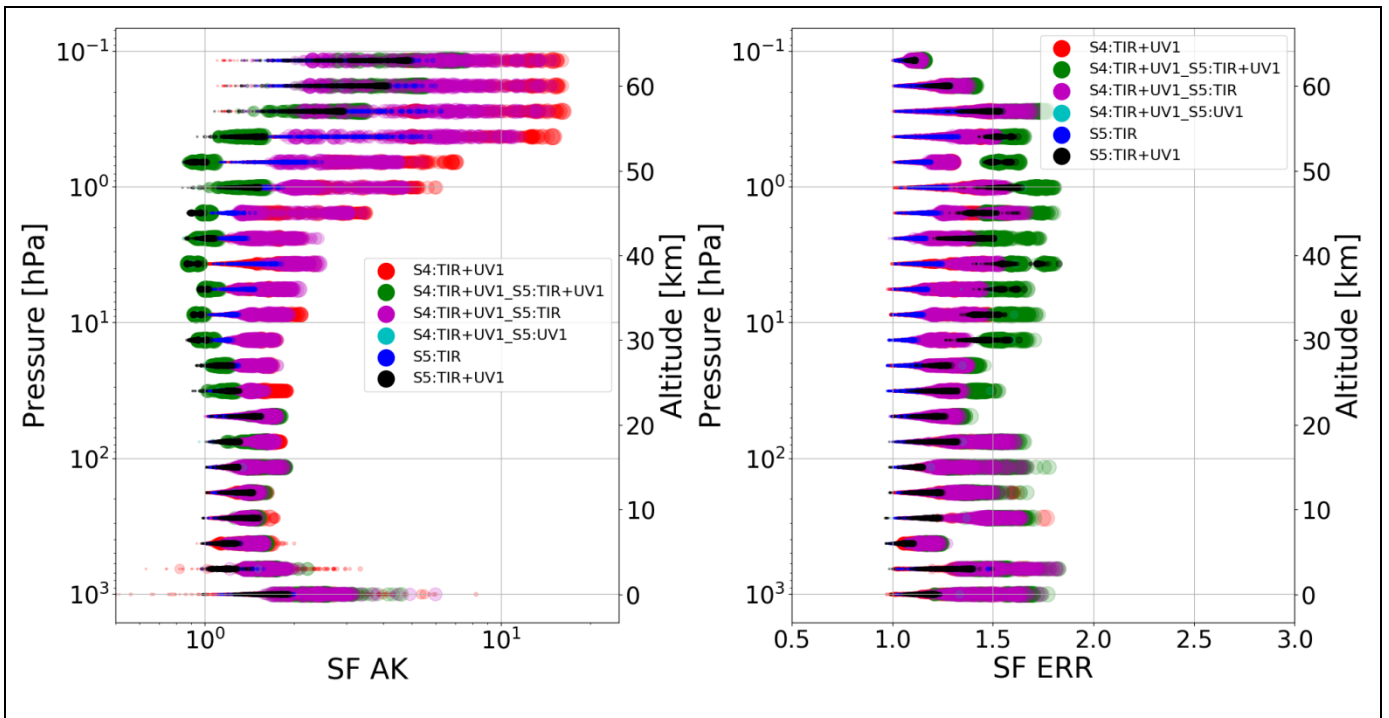


Figure 7 (Left panel): SF AK versus vertical level. (Right panel): SF ERR versus vertical level. In both panels, different colours of the symbols represent the FUS type, different sizes of the symbols represent the number of measurements that have been fused. The maximum symbol size shown in the legend corresponds to N=160.