

## Replies to reviewers

Thanks a lot for taking the time to read the paper and giving us valuable comments. We have changed the manuscript according to the suggestions and have listed our replies and changes in blue below.

-----  
Reviewer 1:

The authors present the assimilation of SO<sub>2</sub> retrievals from Tropomi satellite observations in the global forecasting system used in CAMS for volcanic forecasting. As for other major centres, assimilating vertically-integrated information on SO<sub>2</sub> from space-borne sensors is a challenge which needs continuous improvement, as observational product and data assimilation settings can be refined or improved year after year. This paper is of interest to the community. I suggest it is accepted after modifications are made.

### Scope and title

The title is a bit misleading as the study presented in this manuscript is presenting assimilation experiments carried out in a different system than the near real time (NRT) CAMS system used for volcanic forecasting. Moreover, the present study mainly compares results obtained assimilating the new product proposed by the DLR including information on the SO<sub>2</sub> plume vertical extension, with several settings, to those obtained in the current operational setting with NRT Tropomi data disseminated by ESA. In addition, the study described in this manuscript only focusses on a particular eruptive event, the Raikoke 2019 eruption, which injects SO<sub>2</sub> plumes at very high altitudes. No other event is assessed in this study. Eruptive events release SO<sub>2</sub> plume at a large range of altitudes, depending on the volcano and the given episode. The present paper does not provide any guidance for other eruptive events. I suggest to change the title so as to reflect the content of the paper more closely, such as

"Evaluation of the assimilation of the S5P-Tropomi SO<sub>2</sub> layer height product in the CAMS global system in the case of the Raikoke 2019 eruption".

We have changed the title to:

*'Evaluation of the assimilation of S5P/Tropomi SO<sub>2</sub> layer height data in the CAMS global system for the Raikoke 2019 volcanic eruption.'*

### Assimilation settings for the observations

Section 3.2.1 (235)

The authors describe the baseline configuration and say "SO<sub>2</sub> observations are currently only assimilated ... when the observed SO<sub>2</sub> concentrations are considerably larger than the atmospheric background values". I suggest the authors clearly state that criterion, instead of vaguely referring to "considerably larger".

We already mentioned in section 2.1 : ' Furthermore, only TROPOMI SO<sub>2</sub> pixels with values greater than 5 DU are assimilated in the operational CAMS system to avoid assimilating SO<sub>2</sub> from outgassing volcanoes which are covered by SO<sub>2</sub> emissions in the CAMS model' . For GOME-2 we assimilate all the pixels flagged as volcanic, which is also stated in section 2.1. With the statement in line 235 we only wanted to illustrate that we can not use an NMC style method because the resulting background errors would peak at the surface where anthropogenic emissions lead to the largest SO<sub>2</sub> values. They would not give us background error statistics which would be useful for volcanic eruptions as that information is not in the model's background field. Reading the sentence again, the part 'considerably larger...' is not really needed and we have removed it, so that the sentence now simply reads:

*'SO<sub>2</sub> observations are currently only assimilated in the CAMS system in the event of volcanic eruptions.'*

I may have missed the description of the observation pre-processing in the paper. Can the authors state clearly how the mismatch between the observation resolution and the model resolution? Are data thinned? Is there a super-obbing step? What are the parameters of the pre-processing?

The TROPOMI data are super-obbed to the model resolution. We already mention this in Section 2.1: 'The TROPOMI SO<sub>2</sub> data are averaged to the model resolution (TL511, about 40km) before being used in the CAMS system.'

The GOME-2 data are used at the satellite resolution which is similar to the model resolution. We have added in Section 2.3:

*'The GOME-2 data are used at the satellite resolution which is similar to the resolution of the CAMS model used in this paper.'*

As the number of observations varies between NRT and LH SO<sub>2</sub> observations, a clear indication of the difference in the number of assimilated data should be clearly given.

We already show in Figure 8 a timeseries of the number of observations and have already this text in the paper: 'Figure 8 shows a timeseries of the number of observations that are actively assimilated in both experiments, i.e. the number of 1°x1° grid points with active observations, and illustrates that there are more active data in BLexp where NRT TROPOMI SO<sub>2</sub> data with values greater than 5 DU are assimilated (i.e. as done in the operational CAMS system) than in LHexp where only data with LH TCSO<sub>2</sub> greater than 20 DU are assimilated.'

No word is said on the observation errors, which are also important players in the game. The reader would benefit from a clear description on how the observation errors are handled.

We use the observation errors given by the data providers we have added a sentence in Section 'For the TROPOMI data (and also the other SO<sub>2</sub> products used in this paper) observation errors as given by the data providers are used.'

NRT Tropomi SO<sub>2</sub> observations are provided with averaging kernels. Are these averaging kernels used in the baseline configuration? Are SO<sub>2</sub>-LH observations provided with averaging kernels? If present, are the latter used in the assimilation? I suggest the authors clearly state all these "details".

The NRT Tropomi SO<sub>2</sub> observations are indeed provided with averaging kernels. However, for the volcanic SO<sub>2</sub> product the averaging kernels are simply 1 km box profiles that are used in the AMF calculation to represent typical volcanic SO<sub>2</sub> profiles and do not provide any real information about the current eruption. It therefore does not make sense to use these in the CAMS assimilation system. There are 3 different averaging kernels provide for each SO<sub>2</sub> column retrieval and the user can choose the product that best suits the situation. See TROPOMI ATBD for more information: <https://sentinel.esa.int/documents/247904/2476257/Sentinel-5P-ATBD-SO2-TROPOMI>.

We have added more information at the end of Section 2.1:

*The DOAS vertical column SO<sub>2</sub> retrieval requires knowledge of a prior SO<sub>2</sub> profile to convert the slant columns into vertical columns. Because this profile shape is generally not known at the time of the observation and it is also not known whether the observed SO<sub>2</sub> is of volcanic origin or from pollution (or both) the TROPOMI algorithm calculates four vertical columns for different hypothetical SO<sub>2</sub> profiles. One vertical column is provided for anthropogenic SO<sub>2</sub> with the prior SO<sub>2</sub> profile taken from the TM5 CTM and three for volcanic scenarios assuming the SO<sub>2</sub> is either located in the boundary layer, in the mid-troposphere (around 7 km) or in the stratosphere (around 15 km). These volcanic prior profiles are box profiles of 1 km thickness, located at the*

*corresponding altitudes. The NRT CAMS system uses the mid-troposphere product. TROPOMI SO2 data are provided with averaging kernels based on the prior hypothetical SO2 profiles (i.e. the 1 km box profiles centred around the assumed SO2 altitude for the volcanic columns). However, as these do not provide any real information about the altitude of the volcanic plume they are not used in the CAMS system. More information about the NRT TROPOMI SO2 retrieval can be found in the TROPOMI ATBD. For the TROPOMI data (and also the other SO2 products used in this paper) observation errors as given by the data providers are used.*

### **Minor comments**

line 397: data are gridded for comparison. What is the time step for this gridding: daily or hourly?

The calculation of the analysis or first-guess fields is done at the time and location of the observations in the observation operator of the model. Later, all data (obs or analysis/forecast) in a 12-hour analysis window are interpolated onto a 1x1 degree grid. We have added in section 4.1 (where we first mention the gridding):

*'All the satellite data available during a 12-hour assimilation window were gridded onto a 1°x1° degree grid....'*

Figures showing timeseries are numerous and sometimes hardly legible (eg. 12, 13).

We have improved several of the figures, including Fig 12 and 13.

Figures showing maps are sometimes a bit small (eg. 5, 9)

We think the quality of Figures 5 and 9 is good enough for publication. The main point of the figures is to give an overview of the evolution of the SO2 plume and they are big enough for that.

Do the authors think showing evaluation for D+5 forecasts is relevant for such a study which shows the high sensitivity to the assimilation settings?

As the CAMS forecast system provides 5-day forecasts we think it is relevant to show them.

Reviewer 2 (Nina Iren Kristiansen) :

### **General comments**

This paper presents the CAMS assimilation of volcanic SO2 satellite data, and in particular improvements made to the system by the use of layer height information retrieved from satellite, which show to improve the SO2 forecasts. The paper is interesting and presents both improvements to and current challenges with the system. The topic of the paper is highly relevant as it addresses a method which can be used to fuse models and observations and targets a particular application to volcanic clouds. The paper is well written and highly suited for publication; however, I would like the below comments to first be addressed.

I miss some discussion around the applied/assumed thickness of the SO2 plume and if/how this might affect the results. See specific comment on L260.

I am concerned that the model simulations do not directly consider the vertical averaging kernel information from the SO2 retrieval. See specific comment on L262.

I miss some further details on the TROPOMI SO2 retrieval. Other TROPOMI SO2 total column retrievals are interlinked with assumptions on the SO2 plume altitude and often different products are available based on different a priori plume altitudes (e.g., the Copernicus SP5 products). L122 mention prior SO2 vertical profile shapes but this is not mentioned again or discussed any further for the DLR TROPOMI retrieval (only for IASI on L186). Please elaborate further on which prior profiles are used in the DLR TROPOMI

retrievals (both NRT and LH) and if these vary, and also how/if this affects the retrieval of the layer height.

We have added this information in Section 2.1:

*The DOAS vertical column SO<sub>2</sub> retrieval requires knowledge of a prior SO<sub>2</sub> profile to convert the slant columns into vertical columns. Because this profile shape is generally not known at the time of the observation and it is also not known whether the observed SO<sub>2</sub> is of volcanic origin or from pollution (or both) the TROPOMI algorithm calculates four vertical columns for different hypothetical SO<sub>2</sub> profiles. One vertical column is provided for anthropogenic SO<sub>2</sub> with the prior SO<sub>2</sub> profile taken from the TM5 CTM and three for volcanic scenarios assuming the SO<sub>2</sub> is either located in the boundary layer, in the mid-troposphere (around 7 km) or in the stratosphere (around 15 km). These volcanic prior profiles are box profiles of 1 km thickness, located at the corresponding altitudes. The NRT CAMS system uses the mid-troposphere product. TROPOMI SO<sub>2</sub> data are provided with averaging kernels based on the prior hypothetical SO<sub>2</sub> profiles (i.e. the 1 km box profiles centred around the assumed SO<sub>2</sub> altitude for the volcanic columns). However, as these do not provide any real information about the altitude of the volcanic plume they are not used in the CAMS system. More information about the NRT TROPOMI SO<sub>2</sub> retrieval can be found in the TROPOMI ATBD. For the TROPOMI data (and also the other SO<sub>2</sub> products used in this paper) observation errors as given by the data providers are used.*

We have addressed the other comments in the list of specific comments below.

Can you provide some indications as to how much more expensive (in terms of run time) the model runs are with the higher spectral resolutions used? For example, it would be very interesting to know the difference in run time for each of the experiments in Table 3.

The experiments shown in Table 3 all use the same spectral resolution (model at T511, minimisations at T159/T255) so it does not make sense to add any run time information. Compared to the operational configuration which uses T95/T159 spectral resolutions in the minimisation the numerical cost for one analysis cycle is increased by about 20-30%, with the largest increase from the second minimisation which is about 50% more expensive when going from T159 to T255. We have added this text:

*'The numerical cost of one analysis cycle increases by about 20-30% when the spectral resolution of the minimisation is increased in this way, with the largest increase coming from the second minimisation which is about 50% numerically more expensive.'*

Specific comments

L30 - the last sentence of the abstract: It would be good to include here something about the increase in skill time scales by including the LH information. I would also include a couple more key results here; that including LH information leads to higher modelled TCSO<sub>2</sub> values in better agreement with the satellite observations, but that plume area and burden are overestimated also when including LH data and that the reason for this overestimation is explored.

We have added this to the abstract:

*Including the layer height information leads to higher modelled TCSO<sub>2</sub> values in better agreement with the satellite observations. However, the plume area and SO<sub>2</sub> burden are generally overestimated in the CAMS analysis also when LH data are used. The main reason for this overestimation is the coarse horizontal resolution used in the minimisations. By assimilating the SO<sub>2</sub> layer height data the CAMS system can predict the overall location of the Raikoke SO<sub>2</sub> plume up to 5 days in advance for about 20 days after the initial eruption which is better than what is obtained with the operational CAMS configuration (without prior knowledge of the plume height) where the forecast skill drops much more for longer forecast lead-times.*

L40: "SO<sub>2</sub> in the aircraft cabin is the biggest issue leading to respiratory problems for passengers and crew". Respiratory problems related to SO<sub>2</sub> depend on the SO<sub>2</sub>

concentrations/dose and air quality standards for SO<sub>2</sub> exist. Potential problems also depend on people's underlying health problems like asthma. It doesn't therefore always lead to respiratory problems as this sentence seem to indicate.

We have changed the sentence to: '*SO<sub>2</sub> in the aircraft cabin is the biggest issue and can lead to respiratory problems....*'

L80: You use the different terms 'injection height' / 'plume height' / 'layer height' but not consistently and the difference between them (if any) is not explained. Personally, I'd use injection height as above the volcano and plume/layer height for the cloud altitude away from the vent, but it might be best to keep to as few terms as possible throughout the paper.

In L80 it makes sense to keep injection height, because this is what is determined by Flemming and Inness (2013). We have modified the rest of the paper to not use injection height and only use either layer height or plume height.

L135-147 (section 2.2): I miss some details on how the retrieval of the LH is done and what it relies on besides the exact wavelength ranges used. In which cases does it work well and which not (see related comment on L416). Why does it not work well below 20 DU? Also, what does this LH mean physically? You later use the height of the modelled maximum concentration as the model equivalent, would be good to comment on this here to justify that that is appropriate.

For SO<sub>2</sub> columns below 20 DU the error of the retrieval of the layer height gets larger, making the data less accurate and less useful. The retrieval algorithm is documented in detail in other papers that we refer to. We have added a sentence in Section 2.2 and also a reference to a new validation paper by Koukouli et al. which has just been submitted to ACP:

*For low SO<sub>2</sub> columns, high-altitude layer heights cannot be retrieved and the retrieval is biased towards low layer heights (Hedelt et al., 2018). Therefore, the use of the data in the CAMS system is restricted to values > 20 DU. More details about the retrieval algorithm can be found in Hedelt et al. (2018) and Koukouli et al. (2021). Koukouli et al. (2021) compared the S5P LH data with IASI observations for the 2019 Raikoke, the 2020 Nishinoshima and the 2021 La Soufrière-St Vincent eruptive periods and found good agreement with a mean difference of  $\sim 0.5 \pm 3$  km, while for the 2020 Taal eruption, a larger difference of between 3 and 4  $\pm 3$  km was found.*

L207/section 3.2: The reader needs to know quite a bit about 4DVar assimilation systems to follow this section. It would be good to expand a little in particular on those aspects which you later explore in more detail: background error covariance matrix and the minimisations. Also, observations errors are not mentioned at all, how are errors in the observations taken into account?

The treatment of the background error formulation is already described in detail in section 3.2.1. We have added more information about 4D-var in Section 3.2:

*In the CAMS 4D-Var a cost function that measures the differences between the model's background fields and the observations is minimized to obtain the best possible forecast through the length of the assimilation window by adjusting the initial conditions.*

We have added a sentence in Section 2.1 to state that we use the observation errors provided by the data providers.

L260: "calculate the SO<sub>2</sub> column not between the surface and the top of the atmosphere, but between the pressure values that correspond to the bottom and the top of the retrieved volcanic SO<sub>2</sub> layer. The depth of this layer is currently set in the FP\_ILM retrieval as 2 km, which corresponds to the uncertainty of the retrieved layer height." I

am a little confused about this. Does it mean you use a fixed plume thickness of 2 km to calculate the modelled total columns, i.e., that you only calculate the SO<sub>2</sub> column between the bottom of the plume (retrieved LH – 2 km) and the retrieved LH? What if there is a much thicker plume say several km thick, then the calculation of the SO<sub>2</sub> column loading will miss a large fraction of the SO<sub>2</sub> in the vertical by only summing only over the LH-2km depth.

You are right, the observations assume a depth of the layer of 2 km and we use this to calculate the model equivalent in the observation operator. If the SO<sub>2</sub> layer was deeper this would not be accounted for in our method, but as that information is also not available from the observations we use it is not possible to include it based on assimilation of the SO<sub>2</sub> LH product alone. Additional data (lidar?) which would give vertically resolved information would be needed. Some vertical variation in the SO<sub>2</sub> loading will be achieved if parts of the plume have different altitudes as this information will be available from the observations. As you can see in Figure 3 there is quite a spread in retrieved LH for the eruption and that information will be brought into the SO<sub>2</sub> analysis. Apart from that, we will depend on vertical transport to modify the vertical SO<sub>2</sub> distribution.

We have added this sentence to the paper to document the limitation: *‘One limitation of this method is that the SO<sub>2</sub> LH product gives the plume altitude with an accuracy of 2 km, but does not give a value for the lower vertical boundary of the SO<sub>2</sub> plume, and for a thick plume part of the SO<sub>2</sub> loading could be missed in the calculation of the model equivalent. However, as the model’s background SO<sub>2</sub> concentrations in the free troposphere are low this should not be a big issue in the column calculation. Also, some vertical variation of the SO<sub>2</sub> loading will be achieved if parts of the plume have different altitudes, and Figure 3 shows that this is indeed the case for the Raikoke eruption.’*

L262: “This approach mimics the procedure of using averaging kernels with box profiles given for the SO<sub>2</sub> layer.” I don’t understand how this mimic the use of averaging kernels because if applying an averaging kernel sensitivity, the model data would be multiplied with a different sensitivity/ averaging kernel (AK) value at different vertical levels. Please elaborate. Ideally the satellites vertical AK profiles should be applied to the model data prior to any comparisons to the satellite data – this AK profile can vary from one satellite pixel to the next.

It mimics the use of the averaging kernels that are supplied with the volcanic SO<sub>2</sub> data, which are supplied with the TROPOMI data and are box profiles (see our reply to reviewer 1 above) that represent typical volcanic SO<sub>2</sub> profiles and are used in the AMF calculation to calculate the vertical columns.

We have added information about the AK in section 2.1 and changed the sentence to: *This approach mimics the procedure of using TROPOMI SO<sub>2</sub> averaging kernels which are box profiles, but for the retrieved layer and not an assumed hypothetical volcanic SO<sub>2</sub> profile (see TROPOMI SO<sub>2</sub> ATBD, <http://www.tropomi.eu/documents/>).*

L278: “The ‘dip’ in the TROPOMI SO<sub>2</sub> burden after the initial peak is an artefact that results from missing observations in the TROPOMI NRT data.” This ‘dip’ is not seen in the equivalent time series shown in the de Leeuw paper (their Fig 11) which also show TROPOMI data (different retrieval method). What is the cause of these ‘missing observations’?

In the NRT data we are using there is a data gap in the area of highest SO<sub>2</sub> values (also visible in Figure 9c2) on 25 June. We do not know why. Possible cloud/ ash contamination or data being flagged because of ‘unrealistically’ high SO<sub>2</sub> columns? De Leeuw et al. (2021) do not mention that they used NRT TROPOMI SO<sub>2</sub> data, so we assume they used an offline product for which that problem might have been corrected. We have added this information in the paper:

*The 'dip' in the TROPOMI SO<sub>2</sub> burden after the initial peak is an artefact that results from missing observations in the TROPOMI NRT data on 25 June 2019 in the area of highest SO<sub>2</sub> values (also visible in Figure 9c2 below).*

L300 / Table 3: From the order of the experiments given in the table I expected first the difference between the BLexp and LHexp to be discussed, however the LH50/100/250 cases are first discussed. Perhaps guide the reader at the start of the section to say which experiments are compared first and why. Similarly, would be good there to guide the reader to say that the BLexp and LHexp will be further explored later to assess the skill timescales to see if using a more realistic height rather than the default 5 km will improve the forecasts – a key point and question for the paper.

We think this is addressed in the paper when introducing Table 3 because we already have a paragraph describing the experiments: *'...listed in Table 3. The baseline experiment (BLexp) which assimilated NRT TROPOMI TCSO<sub>2</sub> data with the operational CAMS configuration and the layer height experiment (LHexp) which uses the FP\_ILM S5P LH data with a horizontal background error correlation length of 100 km and background error standard deviation values of 0.7e-7 kg/kg are the main experiments used in this paper (Section 4.3 below) to assess if the assimilation of the SO<sub>2</sub> LH data using a more realistic height rather than the default 5 km improves the CAMS SO<sub>2</sub> analyses and forecasts. The other LH experiments assess the impact of using different horizontal SO<sub>2</sub> background error correlation length scales and various SO<sub>2</sub> background error standard deviation values.'* We have added the green part to make it even clearer.

We can change the order of the entries in the table if the editor deems this necessary.

L415: "TROPOMI NRT lower detection limit": is this a true detection limit from the sensor/retrieval or do you mean rather than you applied a lower DU threshold (5 DU) for the NRT TROPOMI data compared to the SP ILM SO<sub>2</sub>LH retrieval data (20 DU)? Not clear to me if this is a direct 'detection limit' or more a 'chosen threshold' based on various limitations (not necessarily a detection limit). For the 5 DU threshold you mention this is applied to avoid assimilating SO<sub>2</sub> from outgassing volcanoes which are covered by SO<sub>2</sub> emissions in the CAMS model. Also see related question below.

This is a real detection limit. The plots of the NRT TROPOMI data show all available volcanic NRT SO<sub>2</sub> data even though only values > 5DU are assimilated.

L416: "FP\_ILM SO<sub>2</sub>LH retrieval (v3.1) does not provide reliable information for TCSO<sub>2</sub> < 20 DU and therefore only picks up those parts of the plume that are associated with the highest SO<sub>2</sub> load" The work 'information' is ambiguous. Does it mean that both the retrieved column load values and the layer height values are not reliable under 20 DU, or is it only the retrieved layer height data which is not reliable under 20 DU? Maybe to add in section 2.2.

It means the layer height retrieval has too large an error to be useful. We have added more information in Section 2.2. See reply to L135-147 above.

L425/ Figure 9: Would be useful if the figure caption could explain why the NRT TROPOMI data differ to the SO<sub>2</sub>LH TROPOMI data (i.e., DU levels used/displayed).

We already mention this in the text, but have now added in the caption of Fig 9: *'In panels (c)-(e) all available observations are shown, illustrating that the SO<sub>2</sub> LH product only picks up those parts of the plume that are associated with the highest SO<sub>2</sub> load.'*

L430: It is not directly explained why the SO<sub>2</sub> burden is so much larger (2-3 Tg) for the LHexp compared with BLexp. Is it because of higher TCSO<sub>2</sub> values as well as overestimating the plume area? 2-3 Tg is quite a lot higher than the total burden values

from the satellite data.

The overestimation of the plume area in LHexp is actually less than in BLexp for >5 DU. The larger overestimation of the burden in LHexp is likely the result of differences in the background error standard deviation values and the fact that lower SO<sub>2</sub> columns that could correct an overestimation in parts of the plume are not assimilated. We have added this sentence to the paper:

*The larger overestimation of the SO<sub>2</sub> burden in LHexp is the result of differences in the background error standard deviation values used in the experiments and of the fact that lower SO<sub>2</sub> columns, which could correct an overestimation in parts of the plume, are not assimilated.*

L575-L590: It would be good to compare these skill time scales to what was found by de Leeuw et al for the NAME model (skill for 12-17 days for the low-density (<1 DU) parts of the SO<sub>2</sub> cloud and 2-4 days for the denser parts (>20 DU) of the SO<sub>2</sub> cloud).

We have added a sentence to this section:

*Leeuw et al. (2021), using the Met Office's Numerical Atmospheric-dispersion Modelling Environment (NAME) dispersion model, found skill timescales of 12–17 days for low density (> 1 DU) parts of the Raikoke SO<sub>2</sub> cloud and shorter skill timescales of 2–4 days for the denser parts of the cloud (>20 DU). It is interesting to see skill timescales of similar magnitude to the ones obtained in our study even though the method is different. Leeuw et al. (2021) initialized the NAME dispersion model with eruption source parameters and then followed the evolution of the SO<sub>2</sub> cloud, while we use data assimilation to update the location of the plume daily and provide daily SO<sub>2</sub> forecasts with a maximum length of 5 days.*

#### **Technical comments**

Figure text and labels need to be increased as on a print-out version some figures (especially figures 4, 12, 13, 16,17) are very hard or near to impossible to read.

We have improved the figures.

Figure 3: Suggest changing the colour scale as there are very few values >100 DU so hard to distinguish the dots. Is this showing values only >20 DU as you mention the retrieval is accurate only for larger DU values.

We have changed the colour scale to only show values up to 250 DU.

The de Leeuw reference should be updated to the final revised version for 2021.

Done.

The reference Prata et al. 2019 is used in the main text but is not in the reference list.

Added.

Figure 18 could be removed as there are many figures and the difference to Fig 17 is not very big so describing by words in text should be sufficient.

We have removed Figure 18 but kept the text referring to the GOME-2 result.