# OpenContrails: Benchmarking Contrail Detection on GOES-16 ABI

Joe Yue-Hei Ng, Kevin McCloskey, Jian Cui, Vincent R. Meijer, Erica Brand, Aaron Sarna, Nita Goyal, Christopher Van Arsdale, Scott Geraedts

*Abstract*—Contrails (condensation trails) are line-shaped ice clouds caused by aircraft and are likely the largest contributor of aviation-induced climate change. Contrail avoidance is potentially an inexpensive way to significantly reduce the climate impact of aviation. An automated contrail detection system is an essential tool to develop and evaluate contrail avoidance systems. In this paper, we present a human-labeled dataset named *OpenContrails* to train and evaluate contrail detection models based on GOES-16 Advanced Baseline Imager (ABI) data. We propose and evaluate a contrail detection model that incorporates temporal context for improved detection accuracy. The human labeled dataset and the contrail detection outputs are publicly available on Google Cloud Storage at *gs://goes_contrails_dataset*.

*Index Terms*—Contrails, GOES-16, ABI, benchmark datasets, geostationary satellites, machine learning, neural network, image segmentation.

## I. INTRODUCTION

Contrails, or condensation trails, are cirrus clouds created by aircraft when flying through cold and humid regions. By trapping infrared radiation which would have otherwise escaped into space, these contrails warm the earth. Recent works have suggested that the warming produced by global aviation through these contrails is comparable to the warming produced through $CO_2$ emissions [1], [2]. Studies have suggested that the majority of contrail warming can be attributed to a small fraction of flights, and their warming impact can be significantly reduced by using a different route [3], [4]. However considerable uncertainties exist both around the overall impact of contrails [1], and how well their formation and persistence can be predicted in advance [5]. More complete knowledge of where contrails are actually forming is key to reducing this uncertainty, and this can be obtained by creating an automated contrail detector with high spatio-temporal coverage.

Detecting contrails from satellite images is challenging because of their visual similarity to natural cirrus. Contrails are formed as line-shaped ice cloud and then slowly deform over time, becoming difficult to distinguish from natural cirrus. In this work we focus on detecting young contrails that are still linear in shape. While computer vision has made substantial progress thanks to the rise of large scale datasets, like ImageNet [6], and the use of neural networks, recognition on domain specific problems like infrared satellite images are still
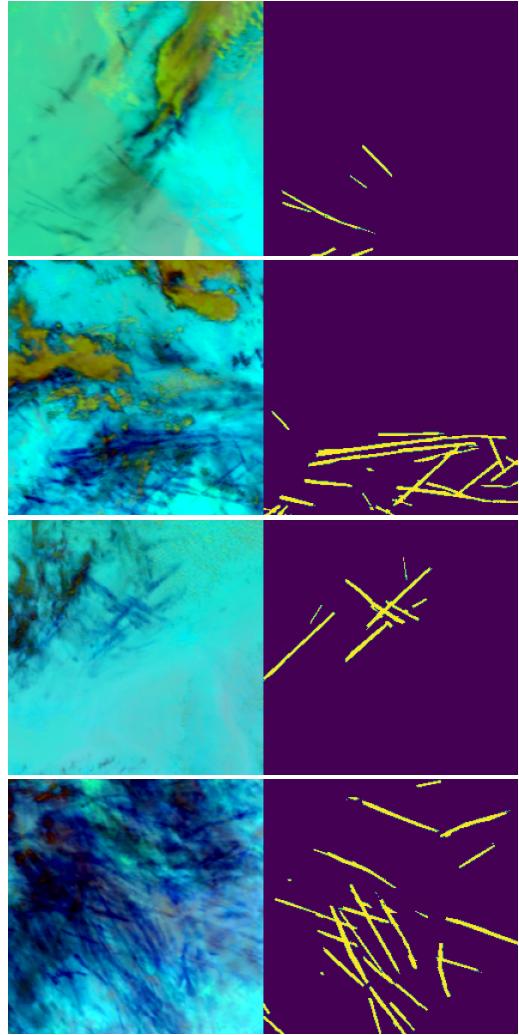
Fig. 1. Example images and labels in the dataset. Left: False color images from GOES-16 satellite; Right: Human labeled binary contrail mask.

challenging due to the lack of large in-domain datasets for pre-training. Apart from that, unlike traditional object recognition tasks where textures and colors are informative, contrails often look similar to or even indistinguishable from natural cirrus. In addition to color and texture, human experts must take into account the overall shape of the contrails (linear) and their temporal evolution (unlike natural cirrus, contrails appear very quickly and spread out rapidly over time).

Challenging computer vision tasks have often benefited

from the existence of high-quality datasets, such as the MNIST dataset [7] for handwriting recognition. Such datasets lower the barrier to entry for different research groups to try to solve the problem, and provide a standard way to compare new models to prior work to assess whether a given new technique has led to improvement. In this work we provide a public dataset, *OpenContrails*, to facilitate model comparison and reproducible research.

While previous works have attempted to train contrail detection models with deep convolutional neural networks [8], [9], the performance of those models are still significantly behind that of a human expert, and therefore improvement in model performance is possible. One possible reason why humans can do better is that during labelling, humans can evaluate the temporal behavior of the contrails, whereas existing computer models perform detection based on only one frame. We develop a model that takes multiple input frames into account for temporal context. We show that such a model improves the the detection performance when compared to models that only process a single frame.

In this work, we present a dataset of human-labeled contrail detections on the geostationary satellite GOES-16, which provides high spatio-temporal coverage, and we make it publicly available. The proposed *OpenContrails* dataset contains high quality per-pixel human label contrail masks for each image in the dataset. We train a neural network for detecting contrails and show that the resulting model produces high quality contrail detection outputs. We run the contrail detection model on multiple years of available GOES-16 images and confirm previous findings on contrail research, such as contrail coverage patterns and diurnal effects. We envision the *OpenContrails* dataset as the benchmark for the contrail detection task, and hope it accelerates contrail research by making large-scale contrail detection results widely available.

### A. Related Work

The effect of contrails on climate change has been widely studied in previous literature [10], [11], [2], [1]. Simulations indicate that a small set of flights are responsible for the majority of contrail warming, which suggests that it is possible to adjust a small number of flight routes to substantially reduce global warming [3], [4].

Empirical analyses of contrails from satellite imagery in the last two decades [12], [13], [14] were typically accomplished by the algorithm proposed by Mannstein et al [12]. Mannstein et al. proposed to detect contrails by applying a set of filters and thresholds on brightness temperature imagery. Zhang et al. trained convolutional neural networks to detect contrails on Himawari-8 satellite images [9] for evaluating contrail coverage. Meijer et al. also trained a neural network for contrail detection on GOES-16 images [8] to examine what effects the reduced flight traffic during the COVID-19 pandemic had on contrails. None of these works made the dataset or contrail detections publicly available, and while Meijer et al. only labeled the continental US (CONUS) region, we sample scenes to be labeled from the majority of the GOES-16 full-disk viewing extent.

McCloskey et al. released a contrail dataset on Landsat-8 data [15]. With 30 meter pixels in its infrared channels, Landsat-8 images have sufficient spatial resolution to identify even very young contrails. However due to its low-earth sun-synchronous orbit and mission goals, imagery is only available for scenes that are mostly land, mostly in the daytime, at fixed local times of day and with a 16-day repeat cycle. Thus Landsat-8 imagery has clear limits for conducting large scale analysis for contrail research. We construct our dataset on GOES-16 ABI images [16] taken from geostationary orbit, which nominally have 2km pixels in the infrared channels but have much higher spatio-temporal coverage. Consequently, our dataset and contrail detection model can be used as a foundation for contrail warming impact assessment and validating contrail avoidance experiments in the western hemisphere.

Deep neural networks have been extensively used for remote sensing applications [17], [18]. The task of contrail detection can be cast as image segmentation, on which deep neural network have shown great success. Here we employ promising image segmentation techniques including using the ResNet backbone [19] and DeeplabV3+ [20], [21] architectures.

## II. DATASET

We build our dataset using GOES-16 Advanced Baseline Imager (ABI) imagery [16], specifically brightness temperatures calculated from the Level 1B radiances using the Planck constants provided by [22]. GOES-16 views the North and South American region with a full-disk image taken every 10 minutes since April 2019 (15 minute interval before April 2019). The nominal pixel size of GOES-16 ABI is $2 \times 2$km at nadir. This relatively coarse resolution means that contrails cannot be seen when they initially form, instead we must wait some time for them to spread out enough to be visible. This is not necessarily a disadvantage: the warming effects of contrails are dominated by contrails which persist for hours [4], so the inability to detect shorter-lived contrails may not hinder our ability to assess and prevent the vast majority of contrail warming.

We generate the examples by randomly sampling image patches from the GOES-16 viewable extent from April 2019 to April 2020. To avoid large viewing angles, we restrict the sampling of image patches to be between -50 to 50 degree latitude and -135 to -30 degree longitude. Uniform sampling of GOES-16 images would result in very few positive examples, as contrails occur infrequently. To increase the number of positive examples in the dataset, we downsampled regions unlikely to contain contrails. We did this in three ways. First, we obtained aircraft flight tracks obtained from terrestrial ADS-B data licensed for publication from FlightAware, LLC (https://flightaware.com). We advected these flight tracks according to European Centre for Medium-Range Weather Forecasts (ECMWF) ERA5 [23] wind data, using the Runge-Kutta method [24], for four hours. We kept only $5\%$ of the images which contained no advected flight tracks, and only $20\%$ of the images which contained less than 10 tracks. Second, we calculated the relative humidity over ice from ECMWF ERA5 temperature, pressure, and specific humidity

values. Since contrails can only persist if this humidity is above $100\%$ [25], we kept only $5\%$ of images for which no location in the image had relative humidity $> 90\%$ for any of the altitudes where contrails might form ($7000 - 12000m$). Finally, we applied the Mannstein et al. [12] contrail detection algorithm tuned as in Mccloskey et al. [15] to achieve very high recall, as opposed to high precision ([15] showed that a similar method applied to images of contrails in Landsat-8 detected every contrail, but for every contrail it detected it also detected around 20 false positives.) We kept only $5\%$ of the images that failed this screen. Note that the FlightAware data does not contain all flights and the ECMWF wind and humidity data are not always accurate. Therefore, the filters are not expected to be perfectly accurate, and it is still necessary to have human labelers annotate the images.

To further boost the number of positives in the dataset, we also included some GOES-16 ABI imagery at locations in the US where Google Street View images of the sky contained contrails. To define when a Street View image of the sky contained a contrail, we used 64-dimensional image feature vectors derived from image-text data, created with an approach similar to that used by Juan et al. [26]. We applied a threshold to the cosine similarity of the Street View image feature vector and that of a seed image of a contrail taken from the ground; if it was similar enough, GOES-16 imagery at that time and location were sampled for human labeling of contrails. These additional labeled images are only in the training set: because Street View cars operate on days with sunnier weather, it may be easier than usual to identify contrails in the GOES-16 imagery of those locations. We found in our experiment that these additional training data slightly improve the contrail detection model performance on the validation set. Note that there are no Street View images in the dataset or used in models reported here: they were only used to identify contrail-rich GOES-16 scenes for human labeling.

We want to detect contrails throughout both day and night, so we show imagery to human labelers in an "ash" false color scheme that combines three longwave GOES-16 brightness temperatures [27]. The red, blue and green channels are represented by the $12\mu$m, difference between $12\mu$m and $11\mu$m, and difference between $11\mu$m and $8\mu$m respectively. This color scheme is chosen to help identify contrails by highlighting ice-clouds as darker colors. An example of the ash color scheme and the true color RGB image is shown in Figure 2. Each image patch in the dataset corresponds to an approximately $500 \times 500$ km region and is reprojected to the Universal Transverse Mercator (UTM) coordinate system (with the zone decided by the northwest corner of the region) using bilinear resampling. Labelers are shown image patches with size of $281 \times 281$. As the labels near the boundary are often noisy due to the lack of spatial context, we center crop to $256 \times 256$ pixels on the images and contrail masks for training and evaluation.

Since the temporal context is important for recognizing contrails, we show labelers 5 images (50 minutes) before and 2 images (20 minutes) after the image being labeled. An example image sequence of contrail evolution is shown in Figure 4, which shows how contrails appear spontaneously and spread out over time, while natural cirrus often occur for a longer time
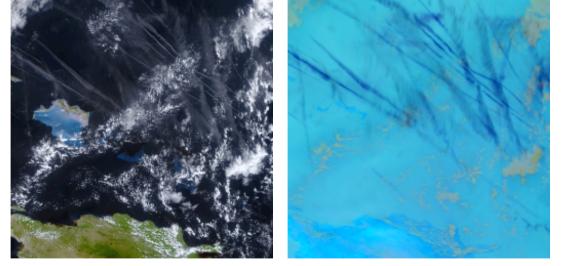


Fig. 2. Left: True color image; Right: false color image in ash color scheme.



Fig. 3. Example GOES-16 false color image with advected flight density shown to labelers.

and may get sharper over time instead of spreading out. The advected flight density is also shown to the labelers to assist labeling. It is generated from FlightAware flight data and is advected using the same method as for screening (ECMWF ERA5 wind vector data, with the Runge-Kutta method). We spread out the advected flight density of flights which are older (at the time the image was taken) to reflect how uncertainty in the wind vector data makes us less certain about where a contrail will ultimately advect to. We ask the labelers to still label the contrails if flights are not shown in the flight density only when they are very confident. An example GOES-16 "ash" image and the corresponding advected flight density is shown in Figure 3.

We ask human labelers to draw polygons on the images to enclose the contrail pixels in a web-based labeling interface. The labelers are instructed on how to label contrails through a set of slides containing guidelines, including examples of contrails and non-contrail cirrus.

We provide detailed instructions for labelers to identify and label contrails, for example: contrails are typically seen as dark linear objects in the satellite images using the ash color scheme; contrails appear as sharp lines and then spread out over time; and static structures across multiple frames are likely to be ground based objects such as roads, rivers or coastlines. We also provide the following additional labeling guidelines:

1) Contrail must contain at least 10 pixels to be labeled;
2) At some time in their life, contrails must be at least 3x longer than their width;
3) Contrails must either appear suddenly or enter from the sides of the image;
4) Contrails should be visible in at least two images.

Labelers were asked to initially label a small set of images, and then authors of this work gave feedback on their initial
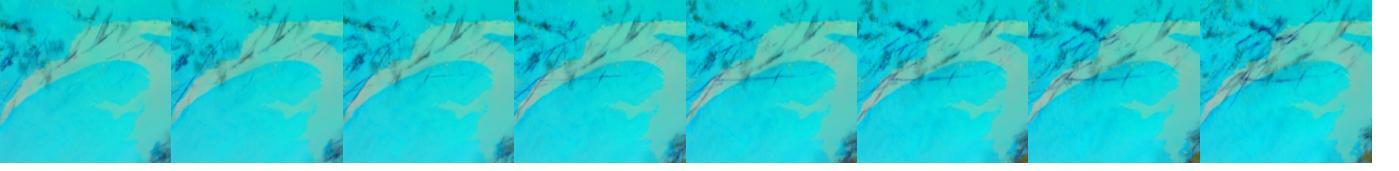
Fig. 4. Example sequence showing evolution of contrails. Contrails appear quickly and start as straight lines. They deform and fade over time slowly from satellite observations.

label quality. The process iterates until a labeler generates labels with satisfactory accuracy. As we observe variance across individual labelers, each GOES-16 scene is labeled by 4 human labelers, and pixels are considered positive if at least 3 labelers considered the pixel a contrail.

Even when instructions and trainings are provided, identifying and labeling contrails on satellite images is still challenging. We compared the labeler performance compared to the groundtruth (labeled by authors of this work) on about 200 examples. Individual labelers give about 64% precision and 76% recall measured per-pixel. The aggregated labels from 4 labelers with majority vote gives about 80% precision and 78% recall.

The full dataset contains 20,544 examples in the train set and 1,866 examples in the validation set. The examples are randomly partitioned except for the satellites scenes that were identified as likely to have contrails by Google Street View, which are only included in the training set. 9,283 of the training examples contain at least one annotated contrail. About 1.2% of the pixels in the training set are labeled as contrails. The dataset contains a wide variety of times and locations, as shown in Figure 5 and 6. The examples are not uniformly distributed in space and time as the images are sampled to include more contrail examples as described above.

Contrails occur more often in cloudy scenes. We estimate the cloud cover fraction using the GOES-16 ABI L2 Cloud Top Phase product and regard all non-clear-sky pixels as clouds. We show the example counts in our dataset by the cloud coverage fraction in Figure 7. Our dataset contains more cloudy scenes than clear sky images, and the positives examples are roughly proportionally distributed at different cloud coverage fractions. It is possible that contrails are harder to be identified in cloudy scenes. Figure 8 shows the labeler agreement at different cloud cover fractions in the scenes, where agreement is defined as the number of contrail pixels labeled by the majority of the labelers divided by the number of contrail pixels labeled by at least one labelers. Labelers agree slightly less in scenes with higher cloud cover fractions.

The dataset is provided in TFRecord format. Each example includes an image with size of $256 \times 256$. Each example is provided with multiple GOES-16 ABI brightness temperatures and brightness temperature differences. The dataset will be uploaded to Google Cloud bucket and made publicly available.

## A. Evaluation Metrics

The best evaluation metric for contrail detection depends on the downstream application. A per-pixel metric would be useful for assessing the energy forcing of contrails as
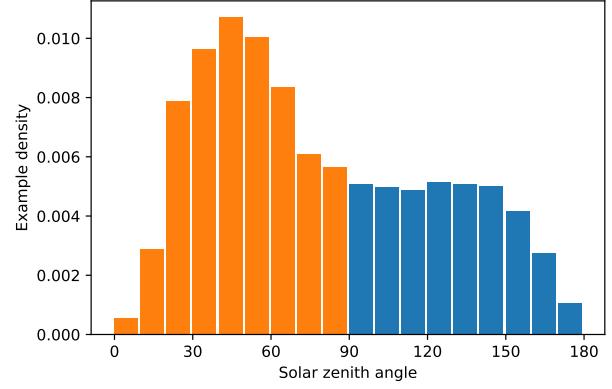


Fig. 5. Distribution of solar zenith angle in the dataset. Orange and blue bars correspond to daytime and nighttime examples respectively.
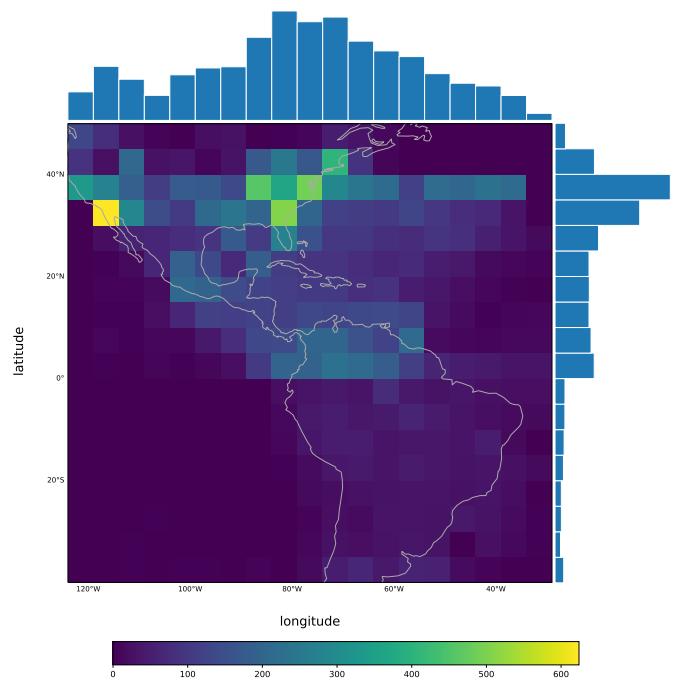


Fig. 6. Spatial distribution according to the image center in the dataset. The top and the right histograms correspond to the marginal distributions along the latitude and longitude axes.
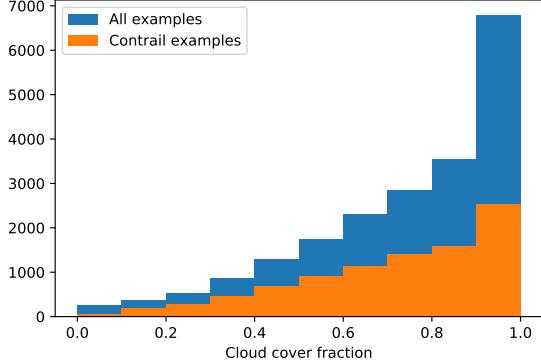
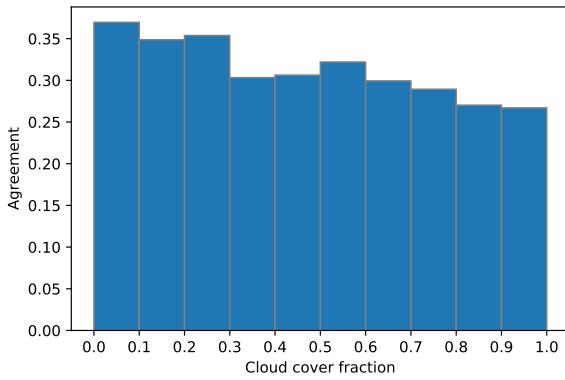Fig. 7. Statistics of cloud coverage in the dataset.



Fig. 8. Labeler agreement at different cloud cover fractions.

it captures the area covered by contrails, and a per-contrail metric would be more useful for associating the contrails to flights and computing the contrail lengths. Here we suggest two evaluation metrics for contrail detection.

Contrail detection can be considered as an image segmentation task which is traditionally evaluated with per-pixel metrics. Typical semantic segmentation metrics like intersection-over-union (IoU) and DICE coefficient are sensitive to the choice of a single binarization threshold. Since the detection outputs could be useful at different operating points for different applications, we evaluate the per-pixel performance using area under the precision recall curve (AUC-PR). We use 10,000 thresholds sampled uniformly between 0 and 1 to compute the precision and recall for each pixel to obtain the AUC-PR.

We also evaluate contrail detections as a line segment detection task. The contrail detection system should output a set of line segments for each image. Similar to [15], we evaluate the resulting linear contrails with a precision recall curve. The contrail detected is considered a correct detection if its angle is within 10 degrees from the groundtruth and within 10km (approximately 5 pixels) mean distance to the contrail [15].

## III. MODELS

We use a convolutional neural network model for contrail detection. The model outputs a score (between 0 to 1) for each pixel indicating the confidence that pixel is part of a contrail. The per-pixel contrail detection outputs are then postprocessed to produce line segments representing the contrail instances.

To consistently identify contrails including at night we train on infrared channels represented as brightness temperatures as input. Our contrail detection model takes 6 channels as input: $8\mu m$, $10\mu m$, $11\mu m$, $12\mu m$, difference between $12\mu m$ and $11\mu m$, and difference between $11\mu m$ and $8\mu m$. Each input channel is standardized by subtracting the global mean and dividing by the global variance of the channel before feeding it into the network.

### A. Single-frame Model

We first train a binary image segmentation model for classifying each pixel as contrail or background. We employ the semantic segmentation architecture DeeplabV3+ [21] and experimented with different image backbones. Following [21], we use a dilated ResNet [20] by modifying the last residual block with dilated convolutions, resulting in an output stride of 16 in the output feature maps. An Atrous Spatial Pyramid Pooling (ASPP) module [21] follows, to incorporate information at multiple scales. Low level feature maps are then combined with upsampled higher level feature maps to preserve the resolution of the feature maps. The final segmentation head consists of three convolutional layers and a 2-channel output representing 'background' or 'contrail'. The model is trained with a per-pixel cross-entropy loss.

It has been shown that the performance of the original ResNet can be improved by various recent techniques [28]. We follow the best practice documented in [28] and use an improved architecture with ResNet-D stem [29] and Squeeze-and-Excitation blocks [30]. We set the squeeze-and-excitation ratio to 0.25. Stochastic depth [31] was used to regularize the network with 0.2 drop rate.

The output resolution of the segmentation model can be smaller than the input images by a factor of the output stride of the backbone. We apply bilinear upsampling to the network prediction to resize the outputs to the groundtruth image resolution during both training and evaluation.

### B. Multi-frame Model

To incorporate temporal context into our model we extend the DeeplabV3+ model to incorporate multiple input frames. We use a spatio-temporal encoder (based on 3D convolutions) to incorporate temporal information.

Inspired by the wide success of convolutional networks in video classification, we use the inflated 3D (I3D) convnet as the backbone of our detection model. Our inflated convnet is based on the ResNet architecture [19]. Since 3D convolution is computationally expensive, we implement the inflated ResNet block by factoring spatial and temporal convolutions similar to [32], [33]. A regular ResNet block consists of three convolution layers with $1 \times 1$, $k \times k$, and $1 \times 1$ kernel sizes. We inflate each ResNet block by expanding the first convolutional layer with a temporal component, while keeping the second and third convolutional layers with a temporal kernel size of 1. The inflated ResNet block thus has kernel sizes of $t \times 1 \times 1$,
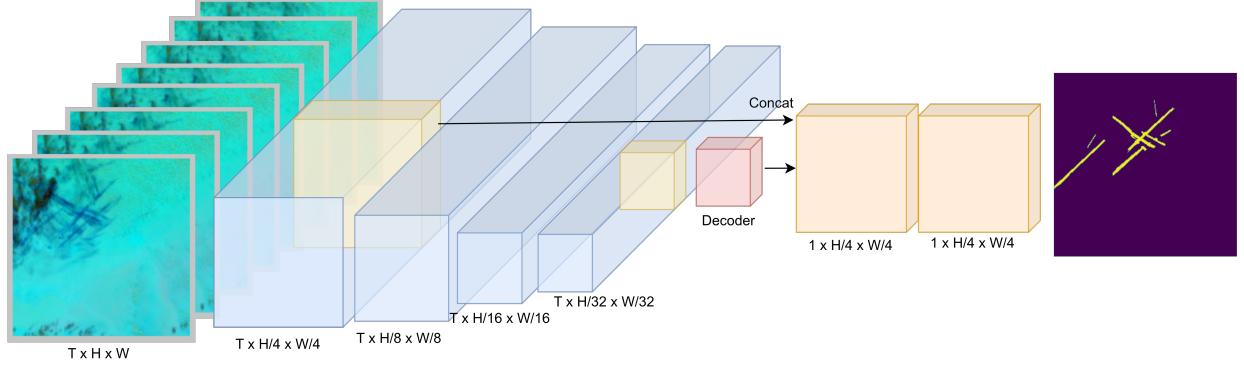
Fig. 9. Network architecture for multi-frame detection model. The model takes a sequence of frames as inputs, and take a slice of the feature maps that correspond to a single timestep, and then run the decoder.

$1 \times k \times k$, and $1 \times 1 \times 1$. In our implementation, we apply the temporal convolution with a kernel size of 3, at every ResNet block in stages 3 and 4.

Differently from video classification (where only a single label is returned for each example), fine-grained spatio-temporal information is required for precisely predicting the contrail detection masks. We augment the backbone network to have a temporal stride of 1 in the first max pooling and convolution layers in the stem. Thus, the temporal resolution is fixed throughout the backbone.

Since the dataset only has one labeled frame available for each example, we slice one frame of features to pass to the decoder. Similar to the single-frame model, we use the ASPP module for the decoder and stacked convolutional layers for classifier. The overall model architecture is illustrated in Figure 9.

### C. Training

We train the network using stochastic gradient descent with 10,000 iterations, using a batch size of 256 and a momentum of 0.9. All models are implemented in TensorFlow and trained on 16 Tensor Processing Unit V3 (TPUv3) chips. The learning rate is set to 0.2 and gradually decays to zero with half-cycle cosine learning rate decay. To improve training stability, we employ a linear warm up learning rate schedule, increasing the learning rate linearly from 0 for the first 500 iterations. L2 weight decay is used with a factor of 0.0001. Data augmentation is applied during training to reduce overfitting: a scale factor is chosen uniformly randomly at each iteration, and then the image is scaled according to the scale factor, and a crop of the scaled image is used for training. We found this significantly helps overall performance.

As our detection dataset is relatively small compared to natural image datasets like ImageNet [6], pre-training the neural network on larger datasets can often help — even pre-training on images which are very different from satellite images. We initialize the network with ImageNet pre-trained weights for the single-frame model. We have also experimented with weights pre-trained on image segmentation datasets like COCO but did not see improvement. For the multi-frame model, we pre-train the backbone model on the Kinetics-400 dataset [34] which consists of about 240,000 videos collected from YouTube: each video clip is 10 seconds long and is annotated with an action from 400 classes.

### D. Converting binary masks to line segments

The network outputs corresponding to pixel-wise contrail probabilities are further post-processed to produce the line segment outputs. We first choose a threshold to convert the pixel probabilities to binary masks, and then use off-the-shelf OpenCV's LineSegmentDetector, which is an implementation of [35], to detect line segments from the binary masks. We then iteratively merge line segments if one segment is part of (or a continuation of) another line segment. We found that this procedure works similarly well compared to a more complicated algorithm which convolves multiple line kernels at multiple scales as described in [12], [15], which requires tuning about 10 hyperparameters for the best performance.

## IV. RESULTS

### A. Main Results

We compare the single frame models and multi-frame model with different backbones. The results are summarized in Table I. All models achieve reasonable detection performance, showing that our dataset is sufficiently large for training contrail detection models. The multi-frame model slightly outperforms the single-frame based models, showing the model is able to use temporal context to improve detection.

TABLE I
DIFFERENT ARCHITECTURES FOR CONTRAIL DETECTION.

| Backbone | Pixel AUC-PR |
|---|---|
| ResNet-50 | 68.9 |
| ResNet-101 | 69.5 |
| Dilated ResNet-101 | 70.5 |
| ResNet3D-101 | **72.7** |

**Linear Contrail Detection.** To assess model accuracy on linear contrail objects, we linearize the output prediction mask as discussed in Section III-D. Before each linearization, we vary the thresholding value for the model output mask to generate multiple binary contrail masks at different confidence
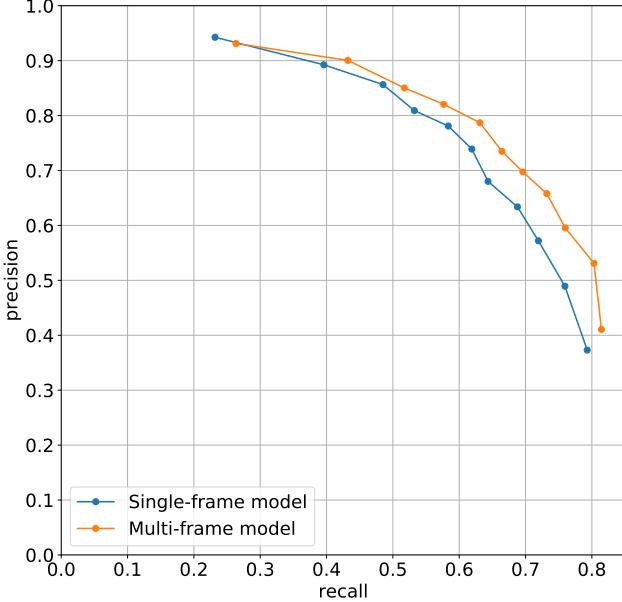
Fig. 10. Per-contrail precision/recall curves, showing better precision and recall for the multi-frame model compared to the best single-frame model we explored.
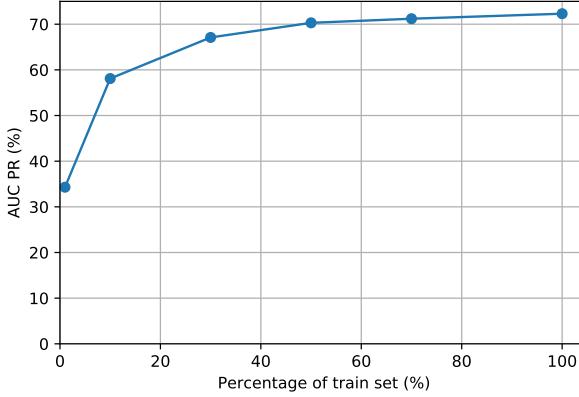


Fig. 11. Validation performance (pixel-wise AUC-PR) when the multi-frame model is trained with different numbers of labels.

levels. Figure 10 shows the precision recall curves from this process. Once again our multi-frame model shows better performance compared to the single-frame model.

McCloskey et al. [15] provided a similar precision recall curve on the contrail dataset based on the Landsat-8 satellite with the baseline Mannstein et al algorithm [12]. It is worth noting that our precision presented here is significantly higher. At 60% recall, The Mannstein et al. algorithm gave 15% precision on the Landsat-8 dataset, while our model gives more than 70% precision on our GOES-16 based dataset. While the evaluation datasets and models are different, this suggests that our model detects contrails more accurately for contrail research.

### B. Detailed analyses

**Input image size.** Using a larger input image size often improves performance on visual recognition tasks [36], [37].

This is because the network can use more computation on smaller patches, and the output resolution is increased. We perform bilinear resampling to enlarge the input image before feeding it to the network, during both training and evaluation. The output predictions are resampled to match the groundtruth resolution when computing the AUC-PR. Table II shows the results with different resolution. We see that larger input image sizes give slightly better results, at the cost of slower inference and more memory usage during training.

TABLE II
CONTRAIL DETECTION (MULTI-FRAME MODEL) WITH DIFFERENT INPUT
RESOLUTIONS.

| Input resolution | Single-frame Model Pixel AUC-PR | Multi-frame Model Pixel AUC-PR |
|---|---|---|
| 256x256 | 63.2 | 65.1 |
| 384x384 | 68,7 | 70.7 |
| 512x512 | 70.5 | **72.7** |
| 784x784 | 71.5 | - |
| 1024x1024 | 71.8 | - |

**Number of input frames.** Using multiple input frames improves the performance by providing more temporal context. The results are summarized in Table III. In general, we observe that the temporal context before the labeled frames is more important than after the labeled frames. This matches our intuition as contrails dissipate over time and become harder to recognize.

TABLE III
DETECTION PERFORMANCE WITH DIFFERENT NUMBERS OF INPUT
FRAMES.

| Size of temporal context | Pixel AUC-PR |
|---|---|
| $T_{\text{before}} = 0, T_{\text{after}} = 0$ | 71.4 |
| $T_{\text{before}} = 0, T_{\text{after}} = 1$ | 71.3 |
| $T_{\text{before}} = 1, T_{\text{after}} = 0$ | 71.7 |
| $T_{\text{before}} = 0, T_{\text{after}} = 3$ | 70.4 |
| $T_{\text{before}} = 2, T_{\text{after}} = 1$ | 71.7 |
| $T_{\text{before}} = 3, T_{\text{after}} = 0$ | **72.7** |
| $T_{\text{before}} = 4, T_{\text{after}} = 3$ | 72.0 |

**Performance across spatial locations.** We visualize the precision and recall (with 0.4 detection threshold) at different spatial locations in Figure 12 for each 10x10 degree gridbox that contains at least one positive examples in our validation set. The precision and recall estimates are noisier in South America due to smaller number of positive examples. In general, we observe no substantial spatial bias in different spatial locations, except with slightly worse performance in the northwest part of South America.

**Performance at different solar zenith angles.** The solar zenith angle is how many degrees away from "directly overhead" the sun is at any given moment; greater than 90 degrees indicates night scenes. Figure 13 shows the precision and recall of the detector at different solar zenith angles. We observe relatively uniform performance across different solar zenith angles, except in the zenith angle bucket with small number of examples which may lead to noisy results.

**Performance at different cloud cover fraction.** Contrails in cloudy scenes could be more difficult to recognize. We show the precision and recall (with 0.4 detection threshold) in
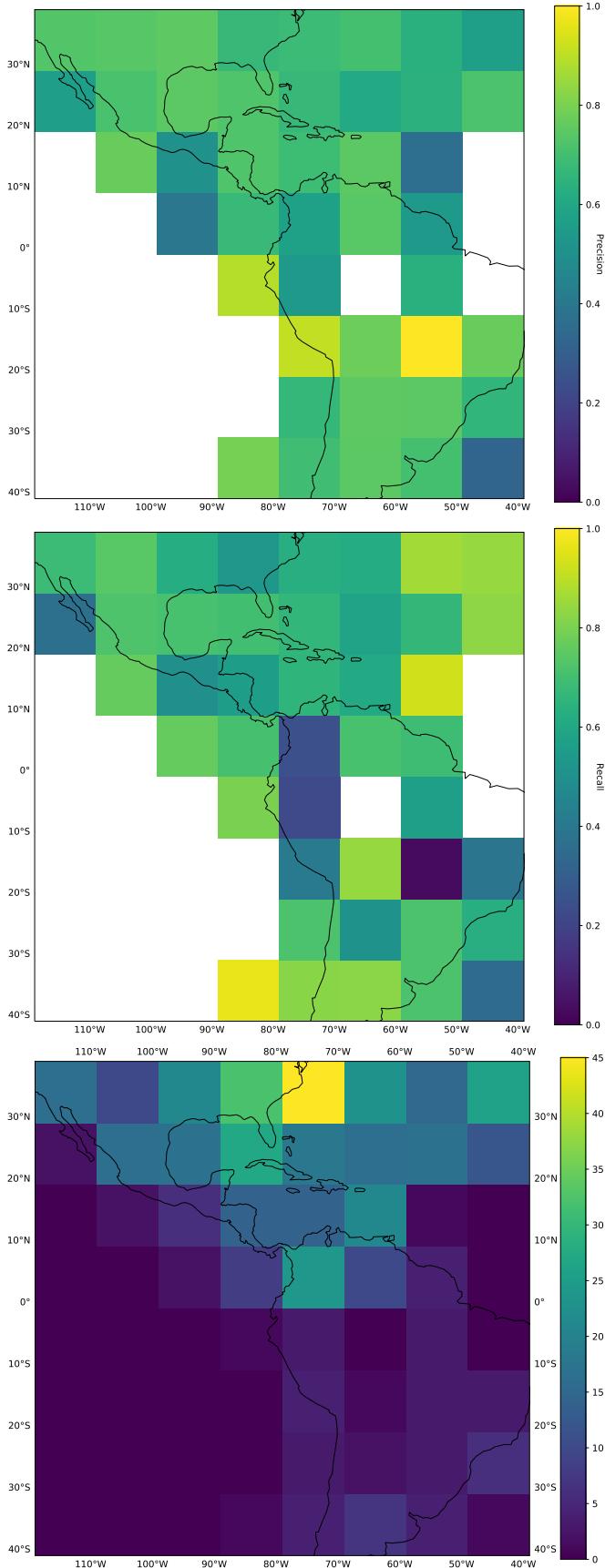
Fig. 12. Detector precision (top), recall (middle), and number of positive examples (bottom) at different spatial locations.

Figure 14. The precision is relatively similar across different cloud cover levels, while the recall slightly drops in higher cloud cover levels.

**Dataset size.** Larger training datasets often improve the quality of machine learning models. Figure 11 shows the effect of training with different amounts of training data. More training data improves the detection performance, but diminishing returns are observed at the current scale. This suggests that it may be more important at this point to improve the model architectures, input features, or learning algorithms than it is to collect more labeled data.

### C. Comparison with Meijer et al [8]

Meijer et al. [8] have also trained a contrail detection model on GOES-16 images. Their model is based on a U-Net architecture with ResNet-18 backbone, and they trained the model on their dataset which consists of 103 labeled images covering the continental US (CONUS) region, with each image having $3000 \times 2000$ pixels. We obtained the evaluation set and detection results from [8] and here we provide an inter-comparison of both models on both datasets. We report on the images from the year 2018 in the evaluation set from [8], which contains human labels on the CONUS region for 12 different GOES-16 scans spread throughout the year. We have verified that there is no overlap between the evaluation set and OpenContrail's training set. Since their dataset uses GOES-16 images from 2018 where full-disk images are taken every 15 minutes instead of 10 minutes, we evaluate our single frame model on this dataset instead of the multi-frame model.

Due to a discrepancy in our respective implementations of image projection, during initial comparisons we observed a slight misalignment between the Meijer et al. human labels when compared to our projections of the GOES-16 images. Therefore, we use a "relaxed" version of precision and recall as the metric, similar to [38]. We consider the relaxed precision represents the fraction of the predicted contrails that are within $\rho$ pixels of the true contrails, and the relaxed recall represents the fraction of true contrails that are within $\rho$ pixels of the predicted contrails. We use $\rho = 2$ in our evaluation.

We compare the detection results without model retraining. We select a threshold value of 0.35 to generate binary detection masks. The results are summarized in Table IV. Our model gives better accuracy than [8] on our dataset, which is unsurprising given our model is trained on data with a similar distribution. Despite the distribution shift of labels from our training set to the labels from [8], our model slightly outperforms their model trained on their dataset. This highlights the good quality and the degree of generalization of our contrail detection model trained on our new dataset.

### D. Visualization

Figure 18 shows some qualitative results of our detections. We observe good visual agreement between the model detections and human label masks, and the line segment detector appears able to convert the detection mask to linear segments reasonably well. The model sometimes misses some detections when there are many contrails in the scene. It is also clear
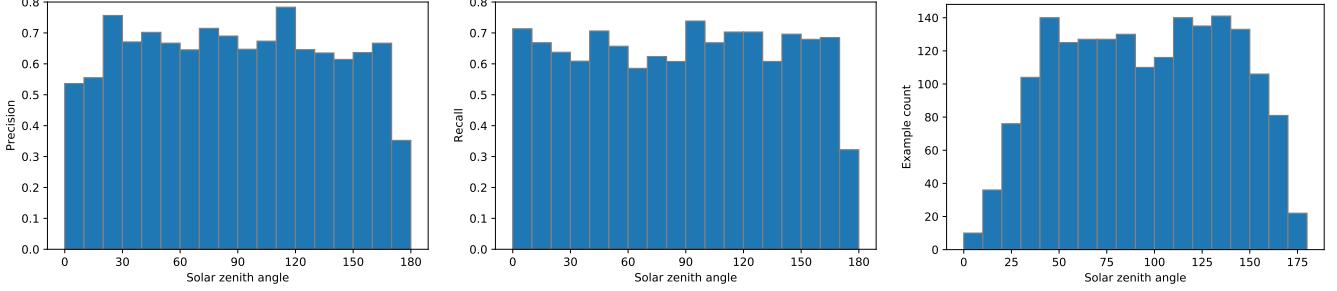
Fig. 13. Detector precision (left), recall (middle), and number of examples (right) at different solar zenith angle on the validation set.
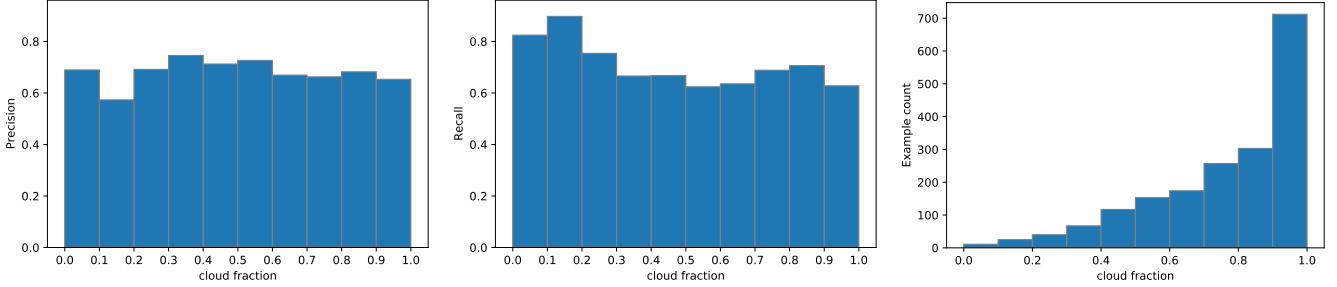


Fig. 14. Detector precision (left), recall (middle), and number of examples (right) at different cloud cover fractions on the validation set.

TABLE IV
MODEL COMPARISON WITH MEIJER ET AL'S DETECTION [8] ON DIFFERENT DATASETS. THE RELAXED PRECISION AND RECALL ($\rho = 2$) OF A FIXED THRESHOLD ARE REPORTED.

|  | Relaxed Precision | Relaxed Recall |
|---|---|---|
| Performance on Meijer et al's dataset [8] | | |
| Meijer et al [8] | 71.6 | 57.2 |
| Ours | **74.5** | **59.7** |
| Performance on our dataset | | |
| Meijer et al [8] | 66.9 | 51.6 |
| Ours | **91.3** | **75.1** |

the model generally does not assign much (if any) probability to anthropogenic contrail-cirrus clouds that are adjacent to still-linear contrails. This is unsurprising given the instructions we have provided to labelers (to only mark linear contrails), but it has an important implication for using the detection model as the basis of contrail warming assessments: we expect the warming from contrail pixels detected by the model to comprise a lower-bound on the actual contrail warming impact.

## V. CONTRAIL COVERAGE

To compute contrail coverage visuals and statistics we run our contrail detection model as an overlapping sliding window, and then stitch the single-frame model detection outputs for the CONUS region in year 2018 – 2019. To give an unbiased contrail coverage estimate, following Meijer et al. [8] we select a per-pixel detection threshold (for our model, 0.4) which gives a precision rate and recall rate on our validation set that are approximately equal. We show the average contrail coverage in 2018 – 2019 in the CONUS region and the corresponding flight density (computed from FlightAware data) in Figure 17. The contrails cluster along the high traffic flight routes, which
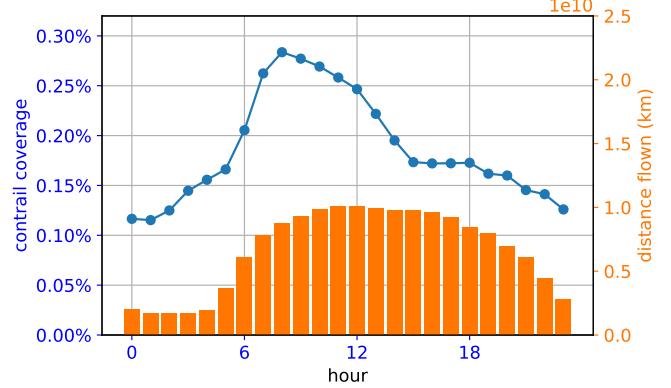


Fig. 15. Linear contrail coverage and flight density by time of day. The blue curve shows the linear contrail coverage derived from our contrail detector and the orange bar shows the flight densities.

confirms the pattern from [8]. The average linear contrail coverage in the considered domain is 0.19% which is similar to 0.17% reported in [8].

Figure 15 shows our estimated linear contrail coverage and flight density by the time of day in 2019. The local hour of the contrails are approximated by $\mathrm{hour_{UTC}} + \mathrm{longitude}/360 \times 24$. We observe a similar diurnal pattern as in previous work [8]: the contrail coverage is lowest around midnight with about 0.1% coverage. It peaks at around 8am with about 0.21% contrail coverage (as the distance flown goes up) but then drops in the afternoon despite the flight distance remaining high.

Figure 16 shows the contrail coverage vs flight distance by the day of year from 2018-2021. We observe similar trends to those observed by Meijer et al [8]. Contrail coverage drops
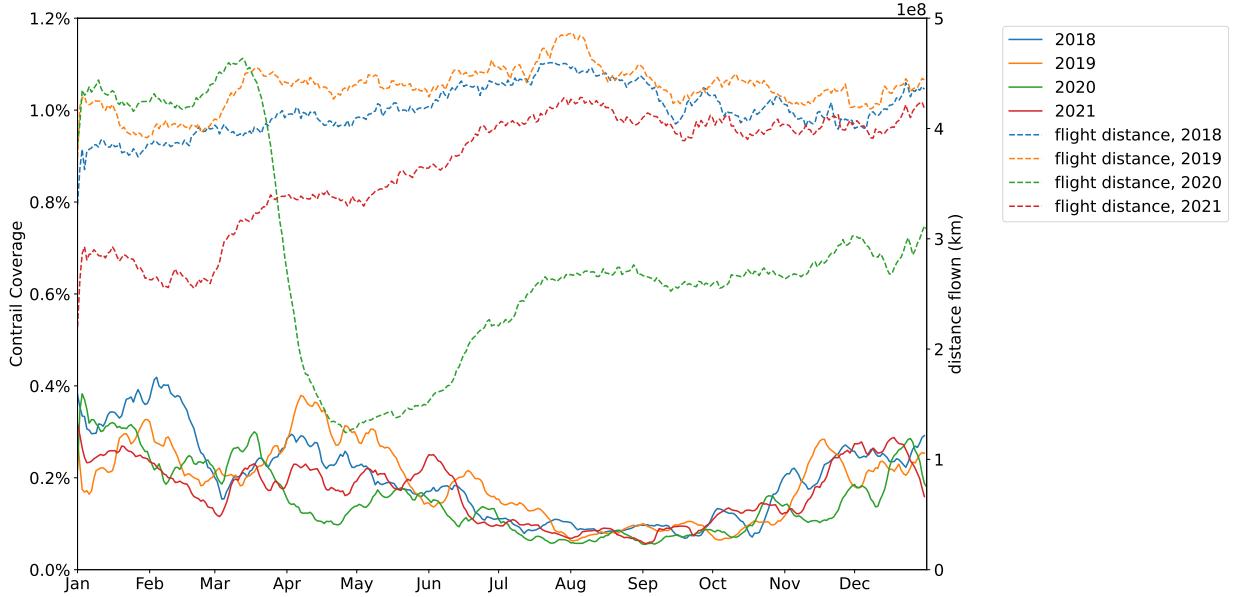
Fig. 16. Linear contrail coverage from our single-frame model with flight distance by day of year from 2018 to 2021. 15-day rolling means are displayed.

as the flight distance drops significantly in April 2020, due to COVID-19. We also observe a seasonal effect of contrail coverage, where fewer contrails are observed consistently in the summer months across multiple years, regardless of the flight distance flown.

## VI. DISCUSSION

While our dataset provides high quality contrail labels on GOES-16 satellite images, some contrails may not be observable from geostationary satellites and thus cannot be labeled by human raters. For example, young contrails may already be deformed (or dissipated) before they are visible from geostationary orbit. Also, contrails forming below existing optically thick cirrus are unlikely to be recognizable by human labelers; however it has been suggested that those contrails do not generate a significant amount of warming [39]. Further research is needed to empirically determine the fraction of contrails and their corresponding warming which can be detected in geoestationary satellite images.

Our current system detects contrails at each target frame independently, and does not provide correspondence of contrails between difference frames. Such correspondence could be useful for identifying the contrail age, and improving flight attribution (i.e., which flight caused the contrail). It is also currently estimated that a large fraction of contrail warming happens after contrails are old enough they are likely no longer linear [39]. Our proposed models and labels only include linear-shaped contrails, and do not identify deformed contrail-cirrus clouds which are no longer distinguishable from natural cirrus without tracing back to when they are formed. To fully assess the warming impact of contrails, it may be crucial to track contrail pixels for an extended period of time once they are detected. Such a contrail tracking system can be built based on the linear contrail models demonstrated in this work.

Due to the high cost of pixel-wise labeling of contrails, our proposed dataset is relatively small compared to large scale natural image datasets like ImageNet (20k vs 1M examples). Large scale labeled datasets for geostationary satellite images are unfortunately not currently available for model pre-training. Therefore, our proposed models are initialized on checkpoints trained on natural image or video datasets such as ImageNet or Kinetics, and their benefits may be limited by domain mismatch between internet images and satellite images. Self-supervised and semi-supervised machine learning techniques that leverage unlabeled data have been shown to improve machine learning model performance for remote sensing segmentation tasks [40], and should be explored here in the future.

Our work here can be extended to cover other geostationary satellites. We plan to collect labels on other geostationary satellites (e.g. Himawari-8 [41] and Meteosat-11 [42]) to cover other regions including Europe and Asia-Pacific. We believe that models trained on the proposed OpenContrails dataset based on GOES-16 can be used for transfer learning to obtain high quality contrail detectors on other satellites with fewer human labels.

## VII. CONCLUSION

We present here a large human-labeled contrail detection benchmark dataset based on GOES-16 ABI imagery. Using it, we explore several contrail detection model architectures, and achieve a new state of the art in linear contrail detection on GOES-16 ABI imagery. We show that leveraging temporal context improves contrail detection, and confirm previously reported seasonal and diurnal contrail coverage trends. The contrail detection dataset and the detection model outputs are made publicly available. We believe these contributions will advance contrail empirical studies including contrail detection, climate impact assessment, and mitigation.
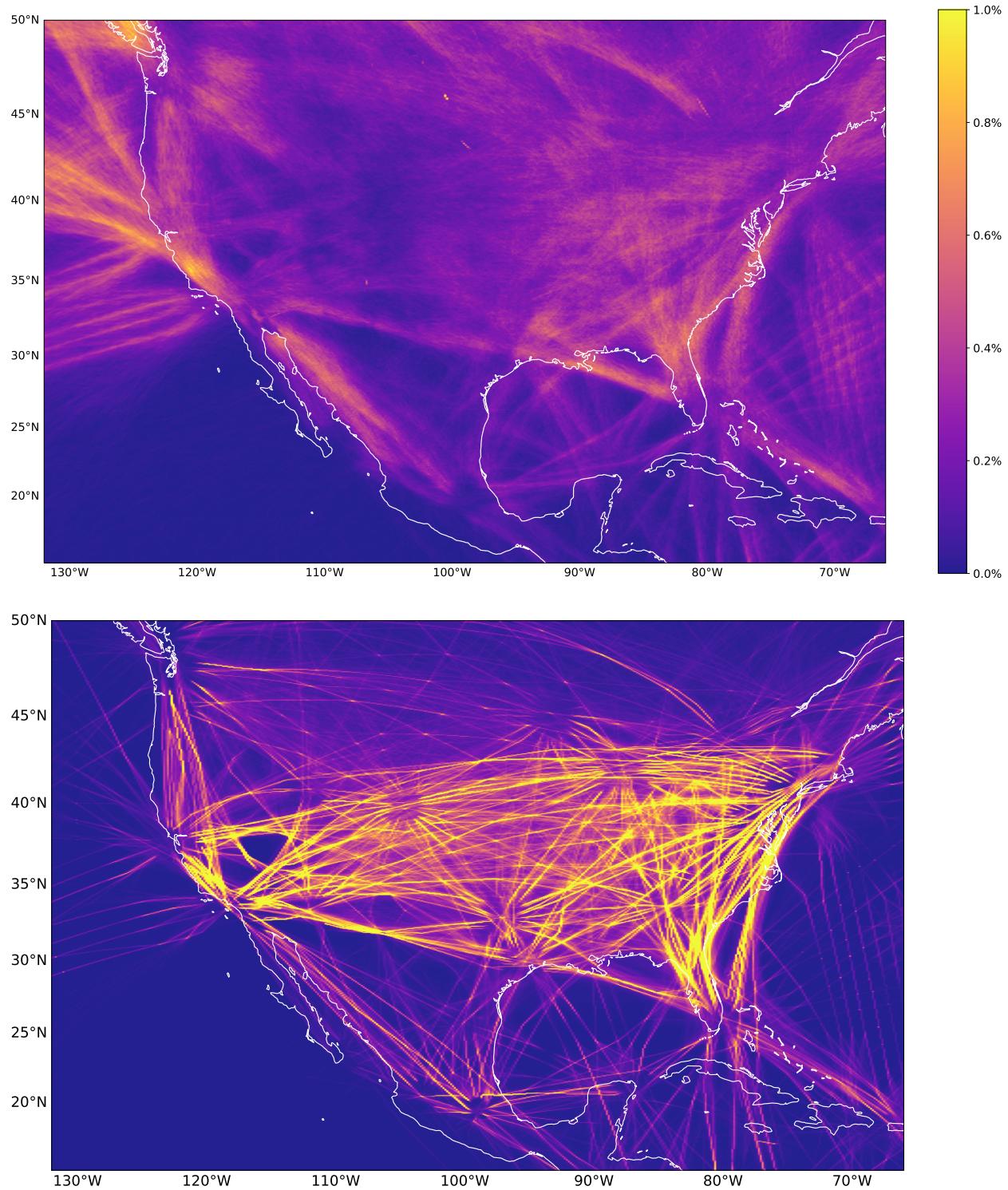
Fig. 17. Top: averaged detected contrail mask, 2018 − 2019. Bottom: flight waypoint density, 2018 − 2019.
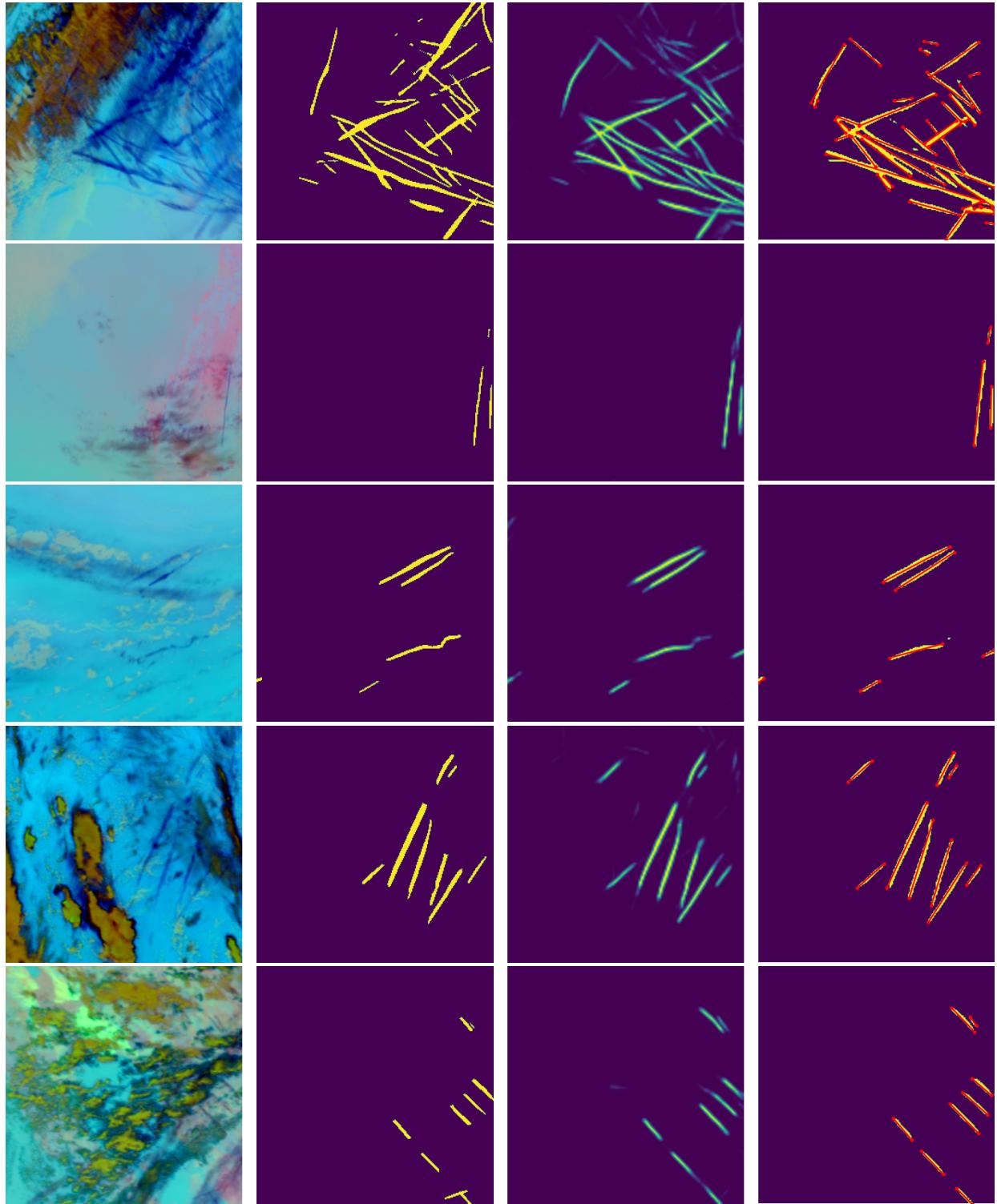
Fig. 18. Example detection outputs from our model. Columns from left to right: 1) False color images from GOES; 2) Ground truth contrail labels; 3) Pixel outputs from our model; 4) Binarized detection outputs (with threshold 0.35) and the red line segments corresponds to the linearization outputs from the line segment detector.

## VIII. Acknowledgement

## References

[1] D. S. Lee, D. Fahey, A. Skowron, M. Allen, U. Burkhardt, Q. Chen, S. Doherty, S. Freeman, P. Forster, J. Fuglestvedt *et al.*, "The contribution of global aviation to anthropogenic climate forcing for 2000 to 2018," *Atmospheric Environment*, vol. 244, p. 117834, 2021.

[2] M. Bickel, M. Ponater, L. Bock, U. Burkhardt, and S. Reineke, "Estimating the effective radiative forcing of contrail cirrus," *Journal of Climate*, vol. 33, no. 5, pp. 1991–2005, 2020.

[3] D. Avila, L. Sherry, and T. Thompson, "Reducing global warming by airline contrail avoidance: A case study of annual benefits for the contiguous united states," *Transportation Research Interdisciplinary Perspectives*, vol. 2, p. 100033, 2019.

[4] R. Teoh, U. Schumann, A. Majumdar, and M. E. Stettler, "Mitigating the climate forcing of aircraft contrails by small-scale diversions and technology adoption," *Environmental Science & Technology*, vol. 54, no. 5, pp. 2941–2950, 2020.

[5] K. Gierens, S. Matthes, and S. Rohs, "How well can persistent contrails be predicted?" *Aerospace*, vol. 7, no. 12, 2020. [Online]. Available: https://www.mdpi.com/2226-4310/7/12/169

[6] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proceedings of IEEE conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.

[7] L. Deng, "The mnist database of handwritten digit images for machine learning research," *IEEE Signal Processing Magazine*, no. 6, pp. 141–142, 2012.

[8] V. R. Meijer, L. Kulik, S. D. Eastham, F. Allroggen, R. L. Speth, S. Karaman, and S. R. Barrett, "Contrail coverage over the united states before and during the covid-19 pandemic," *Environmental Research Letters*, vol. 17, no. 3, p. 034039, 2022.

[9] G. Zhang, J. Zhang, and J. Shang, "Contrail recognition with convolutional neural network and contrail parameterizations evaluation," *SOLA*, vol. 14, pp. 132–137, 2018.

[10] C.-C. Chen and A. Gettelman, "Simulated radiative forcing from contrails and contrail cirrus," *Atmospheric Chemistry and Physics*, vol. 13, no. 24, pp. 12 525–12 536, 2013.

[11] U. Schumann, "A contrail cirrus prediction model," *Geoscientific Model Development*, vol. 5, no. 3, pp. 543–580, 2012.

[12] H. Mannstein, R. Meyer, and P. Wendling, "Operational detection of contrails from noaa-avhrr-data," *International Journal of Remote Sensing*, vol. 20, no. 8, pp. 1641–1660, 1999.

[13] R. Meyer, H. Mannstein, R. Meerkötter, U. Schumann, and P. Wendling, "Regional radiative forcing by line-shaped contrails derived from satellite data," *Journal of Geophysical Research: Atmospheres*, vol. 107, no. D10, pp. ACL–17, 2002.

[14] M. Vázquez-Navarro, H. Mannstein, and S. Kox, "Contrail life cycle and properties from 1 year of msg/seviri rapid-scan images," *Atmospheric Chemistry and Physics*, vol. 15, no. 15, pp. 8739–8749, 2015.

[15] K. McCloskey, S. Geraedts, B. Jackman, V. R. Meijer, E. Brand, D. Fork, J. C. Platt, C. Elkin, and C. Van Arsdale, "A human-labeled landsat-8 contrails dataset," in *International Conference on Machine Learning workshop on Tackling Climate Change with Machine Learning*, 2021.

[16] T. J. Schmit, P. Griffith, M. M. Gunshor, J. M. Daniels, S. J. Goodman, and W. J. Lebair, "A closer look at the abi on the goes-r series," *Bulletin of the American Meteorological Society*, vol. 98, no. 4, pp. 681–698, 2017.

[17] L. Ma, Y. Liu, X. Zhang, Y. Ye, G. Yin, and B. A. Johnson, "Deep learning in remote sensing applications: A meta-analysis and review," *ISPRS journal of photogrammetry and remote sensing*, vol. 152, pp. 166–177, 2019.

[18] X. X. Zhu, D. Tuia, L. Mou, G.-S. Xia, L. Zhang, F. Xu, and F. Fraundorfer, "Deep learning in remote sensing: A comprehensive review and list of resources," *IEEE Geoscience and Remote Sensing Magazine*, vol. 5, no. 4, pp. 8–36, 2017.

[19] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.

[20] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, "Rethinking atrous convolution for semantic image segmentation," *arXiv preprint arXiv:1706.05587*, 2017.

[21] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 801–818.

[22] GOES-R, "Goes-r series product definition and users' guide," 2019. [Online]. Available: https://www.goes-r.gov/products/docs/PUG-L2+-vol5.pdf

[23] H. Hersbach, B. Bell, P. Berrisford, S. Hirahara, A. Horányi, J. Muñoz-Sabater, J. Nicolas, C. Peubey, R. Radu, D. Schepers *et al.*, "The era5 global reanalysis," *Quarterly Journal of the Royal Meteorological Society*, vol. 146, no. 730, pp. 1999–2049, 2020.

[24] P. Bogacki and L. F. Shampine, "A 3 (2) pair of runge-kutta formulas," *Applied Mathematics Letters*, vol. 2, no. 4, pp. 321–325, 1989.

[25] K. Gierens, P. Spichtinger, and U. Schumann, *Atmospheric Physics: Background - Methods - Trends.* Springer, Heidelberg, 2012.

[26] D.-C. Juan, C.-T. Lu, Z. Li, F. Peng, A. Timofeev, Y.-T. Chen, Y. Gao, T. Duerig, A. Tomkins, and S. Ravi, "Graph-rise: Graph-regularized image semantic embedding," *arXiv preprint arXiv:1902.10814*, 2019.

[27] L. Kulik, "Satellite-based detection of contrails using deep learning," Ph.D. dissertation, Massachusetts Institute of Technology, 2019.

[28] I. Bello, W. Fedus, X. Du, E. D. Cubuk, A. Srinivas, T.-Y. Lin, J. Shlens, and B. Zoph, "Revisiting resnets: Improved training and scaling strategies," *Advances in Neural Information Processing Systems*, vol. 34, pp. 22 614–22 627, 2021.

[29] T. He, Z. Zhang, H. Zhang, Z. Zhang, J. Xie, and M. Li, "Bag of tricks for image classification with convolutional neural networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 558–567.

[30] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 7132–7141.

[31] G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Q. Weinberger, "Deep networks with stochastic depth," in *Proceedings of the European conference on computer vision (ECCV)*. Springer, 2016, pp. 646–661.

[32] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, 2018, pp. 7794–7803.

[33] Z. Qiu, T. Yao, and T. Mei, "Learning spatio-temporal representation with pseudo-3d residual networks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5533–5541.

[34] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev *et al.*, "The kinetics human action video dataset," *arXiv preprint arXiv:1705.06950*, 2017.

[35] R. G. Von Gioi, J. Jakubowicz, J.-M. Morel, and G. Randall, "Lsd: A line segment detector," *Image Processing On Line*, vol. 2, pp. 35–55, 2012.

[36] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International Conference on Machine Learning*. PMLR, 2019, pp. 6105–6114.

[37] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2980–2988.

[38] V. Mnih and G. E. Hinton, "Learning to detect roads in high-resolution aerial images," in *Proceedings of the European conference on computer vision (ECCV)*. Springer, 2010, pp. 210–223.

[39] R. Teoh, U. Schumann, E. Gryspeerdt, M. Shapiro, J. Molloy, G. Koudis, C. Voigt, and M. Stettler, "Aviation contrail climate effects in the north atlantic from 2016–2021," *Atmospheric Chemistry and Physics Discussions*, pp. 1–27, 2022.

[40] C. Patel, S. Sharma, and V. Gulshan, "Evaluating self and semi-supervised methods for remote sensing segmentation tasks," *arXiv preprint arXiv:2111.10079*, 2021.

[41] K. Bessho, K. Date, M. Hayashi, A. Ikeda, T. Imai, H. Inoue, Y. Kumagai, T. Miyakawa, H. Murata, T. Ohno *et al.*, "An introduction to himawari-8/9—japan's new-generation geostationary meteorological satellites," *Journal of the Meteorological Society of Japan. Ser. II*, vol. 94, no. 2, pp. 151–183, 2016.

[42] J. Schmetz, P. Pili, S. Tjemkes, D. Just, J. Kerkmann, S. Rota, and A. Ratier, "An introduction to meteosat second generation (msg)," *Bulletin of the American Meteorological Society*, vol. 83, no. 7, pp. 977–992, 2002.