

LATEX Author Guidelines for CVPR Proceedings

Anonymous CVPR submission

Paper ID *****

Abstract

The global increase in the frequency and intensity of wildfires underscores the need for advancements in fire monitoring techniques. In order to investigate deep learning approaches for detecting and tracking wildfires and the related human health impacts, we present *SmokeViz*, a large scale machine learning dataset of smoke plumes in satellite imagery. To build the dataset, we refine a set of human-generated annotations created by analysts at the National Oceanic and Atmospheric Administration. Each annotation gives a general temporal and geographical approximation of smoke plumes but at variable and, primarily, low temporal resolution. We present an innovative solution for refining the temporal and spatial resolution in the given analyst annotations by leveraging the semi-supervised method, pseudo-labeling. Unlike typical pseudo-labeling applications that aim to increase the number of labeled samples, the objective is to use pseudo-labels to refine an existing but coarse-grained set of annotations. We train a deep learning model to generate pseudo-labels that pinpoint the singular, most representative, satellite image to match the smoke annotation within the given temporal range. By identifying the most representative imagery of smoke plumes for a given smoke annotation, the study seeks to create an accurate and relevant machine learning dataset. The resulting *SmokeViz* dataset is anticipated to be an instrumental tool in developing further machine learning models and is publically available at [aws download link].

health effects and premature mortality[14]. These risks underscore the necessity for efficient and effective monitoring methods to mitigate the adverse health impacts associated with wildfire smoke.

Traditionally, wildfire monitoring has relied on ground-based methods, such as forest service patrols, manned lookout towers, and aviation surveillance [3]. While these methods provide valuable localized insights, they are constrained by geographical and logistical limitations, often failing to deliver timely and comprehensive data, especially over large and remote areas. In contrast, satellite imagery offers a vantage point that overcomes these limitations, providing continuous, wide-area coverage and real-time data crucial for assessing and responding to the health risks posed by wildfire smoke.

Satellite imagery, equipped with state-of-the-art sensors, such as the Advanced Baseline Imager (ABI) on the Geostationary Operational Environmental Satellites (GOES) [15], have revolutionized environmental monitoring. Compared to orbiting satellites such as the Suomi or Sentinel satellites, geostationary satellites maintain constant observation over a fixed area. GOES offers the advantage being able to reliably and consistently capture the dynamic behavior of wildfire smoke plumes. In turn, GOES capabilities can provide critical insights into the concentration and movement of smoke particulates, facilitating real-time assessments of air quality.

Integrating satellite imagery into wildfire smoke monitoring provides real-time data that can improve the timeliness of public health planning and response. By mapping the spread and density of smoke, health authorities can issue prompt warnings, implement evacuation protocols, and deploy resources effectively to mitigate health risks. Furthermore, long-term data gathered from satellite observations can aid in understanding the broader impacts of wildfire smoke on public health, influencing policy decisions and preventive measures.

In addition, models for real-time smoke dispersion currently have no smoke analysis product available for data assimilation [1, 18]. This can cause delayed start up times for the smoke to begin being modeled and can result in further

1. Introduction

In part, due to public policy, the average levels of fine particulate matter ($PM_{2.5}$) in the US have generally been declining over the past few decades[2]. Despite those improvements, the contribution of wildfire smoke to $PM_{2.5}$ concentrations in the US has been calculated to have more than doubled between 2010 to 2020, accounting for up to half of the overall $PM_{2.5}$ exposure in Western regions [7]. Increases in $PM_{2.5}$ due to wildfire smoke are concerning since ambient $PM_{2.5}$ exposure is a leading environmental risk factor for adverse

038

039

040

041

042

043

044

045

046

047

048

049

050

051

052

053

054

055

056

057

058

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079 down-the-line errors. Providing a real-time data assimila-
 080 tion smoke product solely dependent on incoming satellite
 081 imagery has the potential to improve existing smoke disper-
 082 sion models.

083 2. Related Work

084 2.1. Numerical

085 Currently, multi-channel thresholding is a popular method
 086 to distinguish smoke pixels from pixels containing dust,
 087 clouds or other phenomenon with similar signatures [39].
 088 Thresholds are determined by using historical, labeled data
 089 to extract optimal radiance values for each channel that cor-
 090 responds with the labeled class. These methods are tuned to
 091 particular biogeographies and often have issues with gener-
 092 alization to new locations with varying fuel types [27].

093 2.2. Analyst

094 In contrast to the numerical thresholding approach, human
 095 visual inspection of satellite imagery is another commonly
 096 used method for smoke identification. Trained analyst in-
 097 spect satellite imagery and label the smoke by hand. An ex-
 098 ample of hand-labeled annotations is the National Oceanic
 099 and Atmospheric Administration (NOAA) Hazard Mapping
 100 System (HMS) fire and smoke product [24, 30]. For the
 101 HMS smoke product, trained satellite analysts use move-
 102 ment characteristics to help identify smoke by scanning
 103 through a time series of satellite imagery. When visual in-
 104 spection indicates smoke, the analyst will draw a polygon
 105 that corresponds to the geolocation and density of smoke.
 106 By design of the product, the HMS annotations have vary-
 107 ing time resolution and are released on a rolling but unde-
 108 fined schedule ranging from one to multiple times a day as
 109 observation conditions permit. This method is potentially
 110 not as scalable as an automated approach and is limited by
 111 the availability of analysts and their time.

112 NOAA manages environmental satellite programs such
 113 as the HMS program, the HMS program is an operational
 114 system that uses an aggregation of satellite data to generate
 115 active fire and smoke data. To train our model, we imple-
 116 ment a supervised learning framework that uses the HMS
 117 analyst smoke product as truth labels during the model
 118 training process.

119 HMS smoke analysis data gives the coordinates of the
 120 smoke perimeter as a polygon and classifies the smoke by
 121 density within a given time window. The time windows can
 122 range from instantaneous (same start/end time) to lengths of
 123 22 hours. While the true bounds of the smoke can change
 124 within the larger time spans, the analyst is making an ap-
 125 proximation that should reflect the smoke coverage over the
 126 duration of the time window. The density information is
 127 qualitatively determined by each analyst based on the ap-
 128 parent smoke opacity in the satellite imagery and cate-

129 rized as either light, medium or heavy as seen in figure 1a
 130 [25].

131 2.3. Deep Learning

132 To address the challenges associated with thresholding and
 133 manual labels, we can look towards innovative approaches
 134 and recent technological advancements in computer vision.
 135 Machine learning methods have shown potential in improv-
 136 ing the accuracy and efficiency of satellite-based wildfire
 137 smoke detection and monitoring. For instance, SmokeNet,
 138 uses a convolutional neural network (CNN) based frame-
 139 work to determine if a scene of MODIS satellite imagery
 140 contains smoke [4]. Another study, that looked at a singu-
 141 lar wildfire event, also used a CNN to identify smoke on a
 142 pixel-wise basis using imagery from Himawari-8 [21]. Ad-
 143 ditionally, Wen et al. developed a CNN architecture that
 144 takes GOES-East imagery as input and the HMS-generated
 145 annotations for the target labels during training [36].

146 The success of deep learning methods, such as CNNs,
 147 relies heavily on the availability of a large, representative
 148 dataset [33]. As laid out in table 1, prior studies use rela-
 149 tively small numbers of samples, from 47 [35] to 6825 [36],
 150 where one sample represents a satellite image with a singu-
 151 lar time and/or geolocation. In contrast, benchmark datasets
 152 for image classification contain tens of thousands (CIFAR-
 153 10 and MNIST) to millions (CIFAR-100 and ImageNet)
 154 of data samples [10, 11, 20]. Keeping in mind the corre-
 155 lation between both the quality and quantity of data with
 156 model performance, we introduce the largest known smoke
 157 dataset, SmokeViz, containing over 130,000 samples.

158 Semi-supervised learning is an approach that can be used
 159 to increase the number of labeled samples in a dataset. This
 160 is done by leveraging a labeled dataset to generate new la-
 161 bels for an often larger, but unlabeled, dataset. Pseudo-
 162 labeling, a form of semi-supervised learning, uses labeled
 163 data to train an initial model, then runs that model on unla-
 164 beled data to predict pseudo-labels, and finally trains a new
 165 model using the pseudo-labels [22]. We introduce a varia-
 166 tion of pseudo-labeling, not to increase the size, but to in-
 167 crease the quality of our dataset by generating pseudo-labels
 168 to select the best satellite image out of a given time-window
 169 to represent each smoke plume annotation.

170 3. Methods

171 3.1. Datasets

172 In order to take into account movement characteristics to
 173 help identify smoke, analysts use multi-frame animations
 174 of the satellite imagery. The resulting annotations primar-
 175 ily have time windows over multiple hours, with an average
 176 of 3 hours of imagery represents one smoke plume anno-
 177 tation. Since the goal of these annotations is to show the
 178 general coverage over that time span, as shown in figure ??,

Table 1. Comparison of different studies including method used, dataset size, satellite source, number of channels used and if classification is performed at a pixel or image level.

| Reference | Method | # Samples | Satellite | # Channels | Level |
|-----------|--------|-----------|----------------|------------|-------|
| [4] | CNN | 6255 | MODIS | 5 | image |
| [36] | CNN | 6825 | GOES-East | 5 | pixel |
| [21] | CNN | 975 | Himiwari-8 | 7 | pixel |
| [35] | U-Net | 47 | Landsat-8 | 13 | pixel |
| SmokeViz | U-Net | 207,106 | GOES-East/West | 3 | pixel |

179 the smoke boundaries don't often match up with the satellite
 180 imagery over the entire time window. One way to approach
 181 this problem would be to use all the satellite images the
 182 analysts used as input. Since the timespans are non-uniform,
 183 this would vary the length in imagery inputs into the model,
 184 which would be difficult with a CNN architecture. More-
 185 over, this would require a large amount of additional mem-
 186 ory and computational resources. Instead of using the origi-
 187 nal analysts' many satellite image inputs to one annotated
 188 output, we develop a one-to-one input-to-output by finding
 189 the optimal singular satellite image input to represent the
 190 annotation.

191 For the set of smoke annotations, \mathcal{Y} , $y \in \mathcal{Y}$ uses one or
 192 more $x \in \mathcal{X}$ where \mathcal{X} is the entire set of satellite imagery
 193 corresponding to the set of time windows defined by the la-
 194 bels. In order to develop a one-to-one data-to-label dataset,
 195 we apply pseudo-labeling to develop a subset of \mathcal{X} , denoted
 196 as \mathcal{X}_p , that has a one-to-one ratio such that $|\mathcal{X}_p| = |\mathcal{Y}|$,
 197 where we choose the satellite image that has the maximum
 198 overlap between the geolocation of smoke in the imagery
 199 and the analyst annotation.

200 But in order to create pseudo-labels we need an initial
 201 parent model, f_o . To train f_o , we need a way of choosing
 202 $x \in \mathcal{X}$ that has a higher chance than random selection of
 203 being representative of y . Discussed in further detail in the
 204 Mie-Derived Dataset subsection, we do this by making a se-
 205 ries of physics-driven choices on which satellite and times-
 206 tamp would give the optimal angle between the sun, smoke
 207 and satellite to produce the strongest smoke signature for
 208 the geolocation and timestamp of the smoke plume. This
 209 dataset, \mathcal{X}_M tells us that if there is smoke present during
 210 the entire time window, which timestamp would give the
 211 highest smoke signal-to-noise ratio.

212 But more importantly than knowing the timestamp for
 213 maximum signal-to-noise, we want to know which image
 214 actually has smoke present within the smoke label bound-
 215 aries. We used \mathcal{X}_M to train f_o , to identify smoke in satellite
 216 imagery, and then use that f_o to create pseudo-labels of each
 217 satellite image in a given annotation's time-window. From
 218 those results, the optimal satellite image is chosen based on
 219 which image's pseudo-labels has the greatest overlap with
 220 the analyst annotation.

Table 2. To create a true color image, we use the following bands from the ABI Level 1b CONUS (ABI-L1b-RadC) product.

| band | description | center $\lambda(\mu\text{m})$ | resolution (km) |
|------|--------------|-------------------------------|-----------------|
| C01 | blue visible | 0.47 | 1 |
| C02 | red visible | 0.64 | 0.5 |
| C03 | veggie NIR | 0.865 | 1 |

3.1.1. Satellite Imagery

The GOES satellites are operated by NOAA in order to support meteorology research and forecasting for the United States. We use the latest operational satellites, GOES-16 (East), 17 and 18 (West) that each carry the ABI, that measure 16 bands between the visible and infrared wavelengths. In improvement to the GOES predecessors, imagery is collected every 5 minutes for the contiguous United States and every 10 minutes for the full disk. Using PyTroll, a Python framework for processing satellite data [28], we input bands 1-3 (Table 2) to a GOES specific true color composite algorithm [5] to develop a, 1km resolution, true color image representation, similar to what is seen by HMS analysts. As discussed in further detail in the next section, the highest signal-to-noise ratio will come from the smallest wavelengths of light, larger wavelengths have lower smoke signal and higher noise (figure 5). For that reason, we only include the first 3 out of 16 available bands of data.

3.1.2. Mie-Derived Dataset

We used a physics-informed approach in selecting the initial GOES dataset, \mathcal{X}_M , which we call the Mie-derived dataset, for training an initial parent model, f_o , where if \mathcal{X} represents all the GOES imagery corresponding to the HMS smoke annotation time window, $\mathcal{X}_M \subset \mathcal{X}$. Prior GOES ABI datasets for machine learning applications often include data from only one of the two GOES-series satellites, commonly opting for GOES-East [36], [26], [23]. Rather than using one satellite or the cumulative data from both GOES-West and GOES-East images, we select between one or the other based on the solar zenith angle. For smoke identification, this approach can achieve a much higher signal-to-noise than imaging the earth's surface from an arbitrary

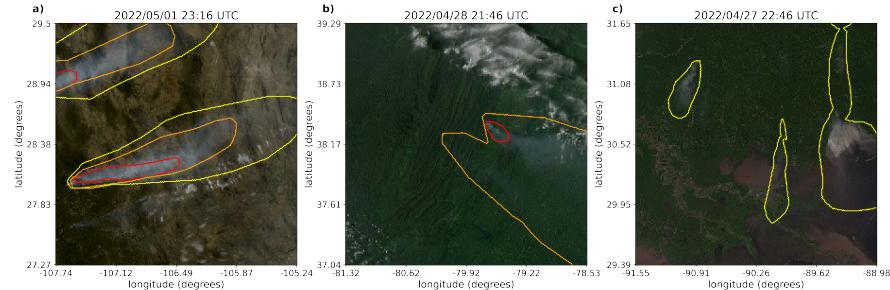


Figure 1. Satellite imagery captured by GOES-East within a few days of each other. The yellow, orange and red contours indicate the extent of light, medium and heavy density smoke. a) shows a canonical example of a smoke plume. b) and c) show observable variations in the density labels.

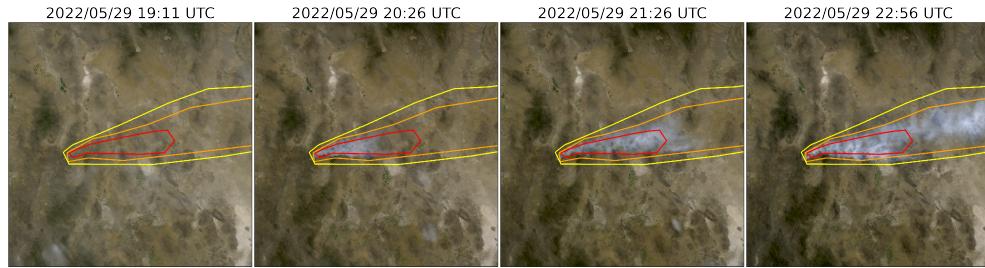


Figure 2. True color GOES-East imagery from May 2022, Southeast New Mexico (31°N , 100°W) during the start of the Foster Fire. The red, orange and yellow lines represent the heavy, medium and low density HMS smoke annotations that span 19:10–23:00 UTC.

angle. The elastic scattering of light is the primary mechanism to account for - while the atmosphere is composed of molecules with size $< 1\text{nm}$, smoke particles can vary from $100\text{ nm} - 10\text{ }\mu\text{m}$ in diameter, d . The GOES ABI covers spectral bands from $0.47\text{ }\mu\text{m} - 13.3\text{ }\mu\text{m}$, so atmospheric and smoke particle sizes occupy two very different regimes with respect to the imaging wavelength λ . In the extreme limit of $\lambda \gg d$, the physics of scattering of light off a small sphere is captured by Rayleigh scattering. This process has two critical consequences: (1) the scattering cross section of light is strongly wavelength dependent (scaling with λ^{-4}), meaning that photons with wavelength closer to the ultraviolet are scattered more strongly than infrared photons. (2) the scattering cross section scales with an angular dependent cross section of $(1 + \cos^2 \theta)$. Scattered photons follow the emission distribution of a radiating dipole, scattering more strongly in the forward and backwards directions ($\theta = 0, \pi$) than orthogonal to the direction of propagation ($\theta = \pi/2, 3\pi/2$), see figure 3 for a Rayleigh scattering schematic.

The significance of these scalings is that the observer, or detector, will receive blue photons in most directions

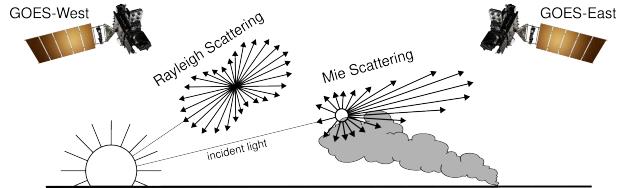


Figure 3. If the particle size is $< \frac{1}{10}$ the wavelength of the interacting light, then the primary scattering will be Rayleigh. Mie scattering is the predominant scattering mechanism when the particle size is larger than the wavelength of light. This schematic demonstrates that when the sun is setting in the West, the Mie scattering will predominately forward scatter towards GOES-East.

orthogonal to the source. Equivalently, photons traveling colinearly with line of sight to the emission source will mostly have wavelengths in the infrared band. In the converse regime of $d > \lambda$, the elastic scattering of light against matter is modeled through Mie scattering. In comparison to Rayleigh scattering, Mie scattering is largely wavelength-independent and has a more complicated radiation pattern where the cross section has a maximal amplitude in the forward direction. An observer downstream of this scatterer

275
276
277
278
279
280
281
282
283



Figure 4. True color GOES-West (left) and GOES-East (right) imagery from April 24th, 2022 in Durango, Mexico. The images were taken ~ 1.5 hours before sunset (01:43 UTC) for this geolocation and time of year.

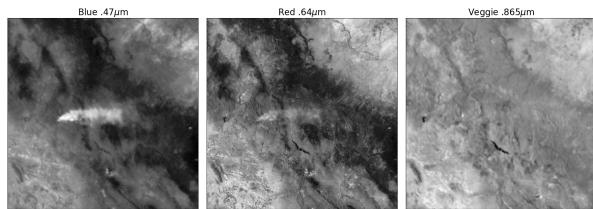


Figure 5. Three bands of GOES-East data are the raw input to generate a true color image. These plots show variations in the signal-to-noise ratio for smoke detection in relation to the λ of light being measured.

will collect more photons than one positioned directly behind it. In the context of smoke identification, a sunrise or sunset will lead to a higher Mie scattered signal in GOES-West and GOES-East respectively, as shown with a smoke plume producing a stronger signal in GOES-East imagery near sunset in figure 4.

Smoke identification therefore amounts to extracting a signal of $d > \lambda$ photons from the $\lambda \gg d$ background. Positioning a detector along line of sight to the scatterer will result in a higher signal from smoke particles (figure 3). Filtering the imaged wavelength can enhance this signal; photons collected in the blue spectrum will have a naturally lower background along the line of sight to the illumination source due to their high level of Rayleigh scattering as. Therefore, as demonstrated in figure 5, this configuration results in the highest signal to noise imaging for smoke particles.

Based solely on these criteria, the optimal strategy would be to pull data from GOES-West right after sunrise and from GOES-East right before sunset. Another factor to consider is that the time when the sun is in optimal alignment with the satellite for smoke detection coincides with when solar zenith angle is close to 90°. Larger angles between the

satellite and sun result in an increase in noise due to increased atmospheric interactions [29]. To reduce the noise from large solar zenith angles, if given multiple frames to choose from, we choose the image with the largest solar zenith angle that is $< 88^\circ$.

The resulting image selection process takes into account atmospheric properties and light scattering physics to generate an estimate of which singular satellite image within the analyst time-window could give the highest smoke signal-to-noise ratio. The resulting Mie-derived dataset, $\mathcal{X}_M = \{X_M, Y\}$, was then used to train a model, f_o , that would generate N pseudo-labels, y^* , for every sample, where N is determined by how many images, taken at a 10 minute interval, fit within the analyst time-window for that sample. Chosen from the N images, x_p is the image with the highest alignment between the f_o prediction of smoke, y^* , in the image and the HMS analysts' annotation y .

3.1.3. Thermometer Encoding Smoke Densities

One of the challenges introduced with using human generated qualitative smoke densities was that, as seen in figure 1b and 1c, there are variations in what is labeled as heavy or light density smoke. More generally, reproducing qualitative metrics with quantitative algorithms is a challenging problem, but we apply mathematical approaches that mitigate some of the underlying complications of our specific problem. Despite the smoke densities introduce qualitative complexities, we decided that the density approximations were important to use in our dataset because of the differences in signatures the densities produce. Within the satellite imagery, the appearance of a light density smoke plume will look significantly different than a heavy density smoke plume as seen in figure 1. Additionally, a light density smoke plume is expected to be more challenging to detect since it is easier for it to be misclassified as not smoke. During the training process, the separate density categories allows us to deferentially weight the penalization given to the model for incorrect classifications based on category. For example, the model can be given a small penalization for misclassifying light smoke as not smoke while given a higher penalization for misclassifying heavy smoke as not smoke.

In addition to the densities being ordered and categorical, the differences between the density categories are not evenly distributed by a given metric, such as PM_{2.5} density. The intervals between densities being unknown along with the hierarchical nature of the density labels makes the labels ordinal instead of just categorical. This data property allows us to use thermometer encoding [6], which leverages the idea that heavy density smoke includes both medium and light density smoke, that heavy density smoke is closer to medium than it is to light, and automatically weights the loss functions and incorporates the ranked ordering of the densities. As seen in Table 3, one-hot encoding, commonly

359 used for categorical data, doesn't take ordinal properties of
 360 the data into consideration.

391

392

393

394

395

396

397

398

399

400

401

Table 3. A comparison of one-hot encoding used for categorical data to thermometer encoding for ordinal data.

| category | one-hot | thermometer |
|----------|---------|-------------|
| No Smoke | [0 0 0] | [0 0 0] |
| Light | [0 0 1] | [0 0 1] |
| Medium | [0 1 0] | [0 1 1] |
| Heavy | [1 0 0] | [1 1 1] |

361 3.1.4. Pseudo-label Dataset

402

We implement a deep learning architecture that uses the encoder from EfficientNetV2 [34] and a semantic segmentation classifier from the DeepLabV3 model [9]. Transfer learning has shown to reduce the time and resources needed to train a model by leveraging information from pre-trained models [37], [31]. We initialize the values of our model weights using the pre-trained values originally trained on the ImageNet dataset [10], containing 1.2 million images and 1000 categories. Our model was developed using the Segmentation Models PyTorch package [16] that was written as a high level API for implementing models for semantic segmentation problems. We input 256x256x3 snapshots of 1km resolution true color GOES imagery that contains smoke and output a 256x256x3 classification map that predicts if a pixel contains smoke and if so, what the density of that smoke is. As mentioned earlier, we apply the thermometer encoding shown in table 3 to encode the smoke densities and apply binary cross entropy as the loss function per density of smoke.

403

404

405

406

407

The dataset, \mathcal{X}_M , contains 207,106 samples as shown in the dataset split in table 4.

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

Table 4. Dataset split for \mathcal{X}_M and \mathcal{X}_p , samples for 2024 go up to November 1st. We use an entire year of data for both validation and testing sets to capture year-long wildfire trends.

| dataset | \mathcal{X}_M | \mathcal{X}_p | years |
|------------|-----------------|-----------------|-----------------|
| training | 165,609 | 144,225 | 2018-2021, 2024 |
| validation | 20,056 | 19,223 | 2023 |
| testing | 21,541 | 16,855 | 2022 |

To determine which image out of the relevant imagery for the given time window best represents the analyst annotation, we implement a greedy algorithm by running f_o on each x to generate a pseudo-label, y^* . The output of f_o , y^* give predictions on if smoke is in the image, and if there is smoke, where the smoke is in that image and the density of that smoke. y^* serve as pseudo-labels for each density of smoke and are compared to the analyst annotations, y . To

compare y^* and y , we calculate the IoU using the total set of pixels for y^* at that density of smoke and the entire set of pixels for y for a particular smoke density in each image as shown in equation 1. The image with the highest IoU score is chosen as the image, x_p , that best represents the analyst smoke annotation, y . Often used for pseudo-labeling, a confidence threshold value is defined to determine if a pseudo-label should be included in a dataset [13]. We chose a confidence threshold that would include the sample, x_p , in \mathcal{X}_p if the maximum overall IoU (equation 1) between y^* and y over all densities was over 0.01.

$$IoU_{\text{overall}} = \frac{\sum_{i=\text{light}}^{\text{heavy}} |y_i \cap y_i^*|}{\sum_{i=\text{light}}^{\text{heavy}} |y_i| \cup |y_i^*|} \quad (1)$$

We use \mathcal{X}_p to train an additional child model, f_c in order to assess if training with \mathcal{X}_p can produce a more robust semantic segmentation model compared to training on \mathcal{X}_M . We use the same dataset split method and model setup but change \mathcal{X}_M to \mathcal{X}_p to train f_c .

3.2. Benchmark Models

While this dataset is anticipated to be primarily useful for solving various wildfire smoke applications, this dataset could be a uniquely insightful test case for remote sensing semantic segmentation. Many deep learning satellite image datasets are focused on objects with sharp contrasts such as crops [17], human infrastructure [40], or even clouds over oceans [19] [32], but smoke has indistinct boundaries that often fade both spatially and temporally.

We benchmark the SmokeViz dataset, \mathcal{X}_p by varying the semantic segmentation classification heads. We train Linknet [8], PSPNet [38] and Ma-net [12] using the same encoder used for f_c and f_o , EfficientNetV2. Each model is trained over 100 epochs using a batch size of 32 and the Adam optimizer on 8 Nvidia P100 GPUs allocating 10GB of memory over 12 hours of allotted training time.

4. Results

To interpret the performance of f_o , we report the IoU metrics in table 5 that were computed by running f_o and f_c on \mathcal{X}_M and \mathcal{X}_p . For each density, we calculate the IoU using the total set of pixels that f_o predicts as that density of smoke and the entire set of pixels labeled by the analyst as a particular smoke density over all imagery contained in the testing dataset. Additionally, we compute the overall IoU for all densities by first computing the number of pixels that intersect their corresponding density and divide that by the total number of pixels that make up the union of model predicted and analyst labeled smoke in the testing dataset.

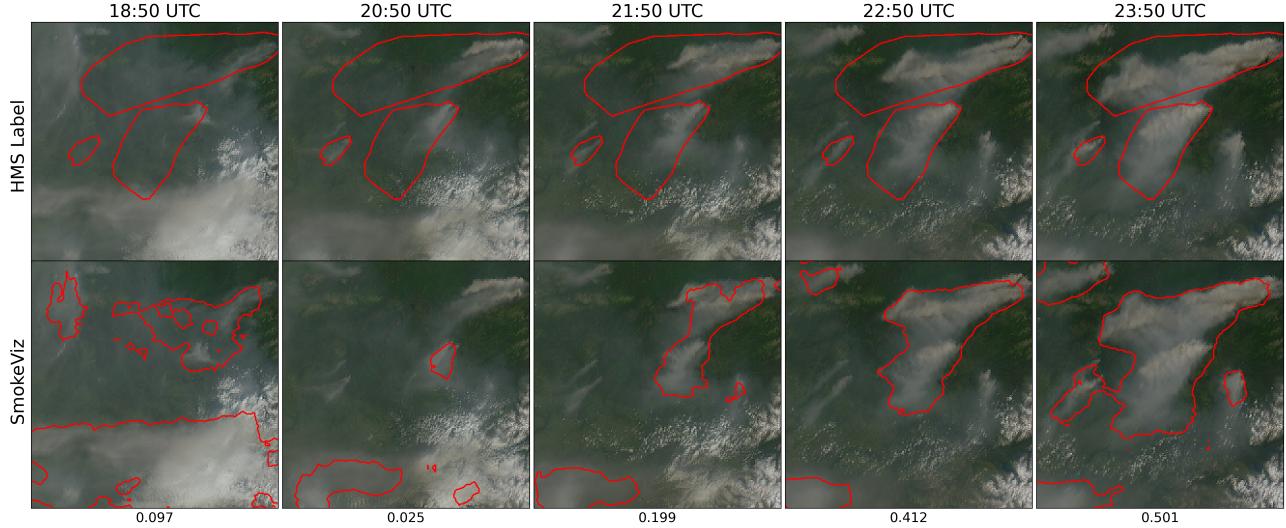


Figure 6. GOES-West imagery showing smoke on June 8th, 2022 in Alaska where, at this geolocation, daylight was between 12:43-7:53 UTC. The HMS smoke annotations (top row) span from 18:50 to 23:50 UTC and are compared to the f_o generated pseudo-labels (bottom row). The first column would be the GOES imagery selected for \mathcal{X}_M since it is closest to sunrise. The last column was selected for \mathcal{X}_p since it had the highest IoU value between the pseudo-label and analyst annotation.

Table 5. IoU results per density of smoke and over all densities using f_o and f_c with \mathcal{X}_M and \mathcal{X}_p .

| | f_o | | f_c | |
|---------|-----------------|-----------------|-----------------|-----------------|
| | \mathcal{X}_M | \mathcal{X}_p | \mathcal{X}_M | \mathcal{X}_p |
| Heavy | 0.278 | 0.368 | 0.218 | 0.411 |
| Medium | 0.310 | 0.417 | 0.319 | 0.484 |
| Light | 0.480 | 0.585 | 0.491 | 0.660 |
| Overall | 0.430 | 0.533 | 0.438 | 0.607 |

An illustration of a pseudo-label picked image better representing the analyst annotation when compared to the Mie-derived image selection is evident in Figure 6, where the heavy density smoke IoU increases from 0.01 to 0.59. The analyst annotation for these densities cover 5 hours of imagery, the Mie-derived selection optimizes for the image closest to sunrise while the pseudo-label image selection chooses the image with the highest overlap between the pseudo-label and the analyst annotation. The figure also illustrates how using a deep learning model can provide higher time resolution and give a dynamic representation of smoke over time.

To get an idea on how f_c compares to the HMS analyst annotations we show a series of samples from \mathcal{X}_p in figure 6. The examples give a qualitative representation of how the predictions from f_c can provide more detailed boundaries of smoke densities than the HMS annotations do.

The results for the benchmarking models (table 6) show similar performance across the models. DeepLabV3+ (f_c)

Table 6. Comparison of semantic segmentation model IoU performance on \mathcal{X}_p .

| | DLV3+ | MANet | PSPNet | Linknet |
|---------|-------|-------|--------|---------|
| Heavy | 0.411 | 0.336 | 0.355 | 0.324 |
| Medium | 0.484 | 0.487 | 0.502 | 0.456 |
| Light | 0.662 | 0.675 | 0.690 | 0.662 |
| Overall | 0.607 | 0.615 | 0.626 | 0.601 |

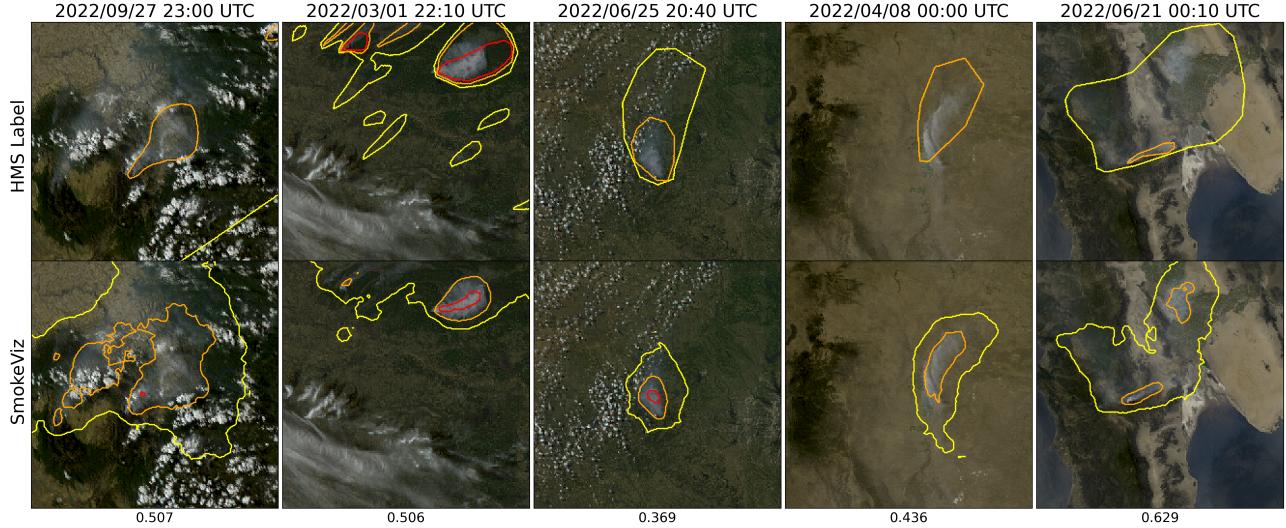
gives the highest heavy density smoke IoU value, while PSPNet gives the highest overall IoU score.

5. Limitations

One of the concerns that comes with using pseudo-labeling methods is that you can perpetuate biases from the parent model into subsequent child models. Due to the increase in detectable forward scattered light off smoke particular matter, we expect the model to have a bias towards producing a higher success rate for smoke detection at larger solar zenith angles. The original HMS annotations do not distinguish by type of fire and include a large representation of controlled agricultural burns. This can be a limitation to consider if the dataset is being trained to target detection of large wildfires. All these limitations are discussed and analyzed further in the Appendix. Additional work should be done to analyze the performance of SmokeViz derived models on dust vs smoke.

455
456

457
458
459
460
461
462
463
464
465
466
467
468
469
470
471

Figure 7. Examples of HMS annotations (top row) vs f_c output (bottom row) on \mathcal{X}_p samples.

472

6. Conclusion

In this study, we have refined an existing dataset originally curated by NOAA’s HMS team, transforming it from a many-to-one imagery-to-annotation format to a more succinct, one-to-one satellite image-to-annotation dataset. The initial HMS dataset provided a general approximation of where smoke had been present for a given time window, though it did not guarantee the actual existence of smoke in the labeled pixels during the given times. Our goal was to create a dataset that could be used, along with additional applications, to train a model to detect wildfire smoke in real-time on an image-by-image level. The Mie-derived dataset selection process determined that if smoke was present, what timestamp within the analyst time window would the give the highest smoke signal-to-noise ratio. While optimizing for being able to detect smoke, if it is present, the Mie-dataset selection had no metric to determine if the smoke was effectually present in the selected image. Since many of the images within the HMS time-window either contained no smoke at all or the smoke was not contained within the geospatial bounds of the annotations, the Mie-derived dataset contained a large number of mislabeled samples. Discrepancies between data and labels can be detrimental towards the model’s capacity to improve on feature representations in the target domain. During model training, the penalization of accurate predictions can inadvertently introduce biases towards misclassifying noise as meaningful signal.

To improve the dataset’s capacity to accurately represent wildfire smoke plumes, we train a parent machine learning model, f_o , using the Mie-derived dataset, \mathcal{X}_M , and run it on the relevant satellite images within the time-frame. The

image with the maximum IoU score between the model’s smoke predictions, or pseudo-label, and the analyst smoke annotations are used to create the pseudo-label generated dataset, \mathcal{X}_p . We then train a child model, f_c , using \mathcal{X}_p and test f_o and f_c on both the 2022 testing sets from \mathcal{X}_M and \mathcal{X}_p . The results reported in table 5 suggest that \mathcal{X}_p was able to train a better performing model, f_c , that gave higher IoU metrics on both dataset’s testing sets in comparison to the original parent model, f_o .

The result of this study is a representative dataset, SmokeViz, that can be used to train machine learning models for various wildfire smoke applications. A future goal is to produce a robust and reliable machine learning based approach for detecting wildfires using satellite imagery. That information can be used for wildfire detection and monitoring in along with a highly needed smoke product for data assimilation into smoke dispersion models. Additionally, this dataset can be used as a benchmark for how well remote sensing segmentation models can perform on dispersed edges such as smoke. On a broader scale, we show how pseudo-labeling can be used to optimize a dataset when the resolution for the data and corresponding labels do not match. This could be useful in similar applications involving time-series/video data with a singular label where the data can be compressed while still remaining representative of the label.

References

- [1] Ravan Ahmadov, Haiqin Li, Johana Romero-Alvarez, Jordan Schnell, Sudheer Bhimireddy, Eric James, Ka Yee Wong, Ming Hu, Jacob Carley, Partha Bhattacharjee, et al. Forecasting smoke and dust in noaa’s next-generation high-resolution coupled numerical weather prediction model. Technical re-

- 536 port, Copernicus Meetings, 2024. 1
- 537 [2] Joseph E Aldy, Maximilian Auffhammer, Maureen Cropper,
538 Arthur Fraas, and Richard Morgenstern. Looking back at 50
539 years of the clean air act. *Journal of Economic Literature*, 60
540 (1):179–232, 2022. 1
- 541 [3] Robert S Allison, Joshua M Johnston, Gregory Craig, and
542 Sion Jennings. Airborne optical and thermal remote sensing
543 for wildfire detection and monitoring. *Sensors*, 16(8):1310,
544 2016. 1
- 545 [4] Rui Ba, Chen Chen, Jing Yuan, Weiguo Song, and Siu-
546 ming Lo. Smokenet: Satellite smoke scene detection using
547 convolutional neural network with spatial and channel-wise
548 attention. *Remote Sensing*, 11(14):1702, 2019. 2, 3
- 549 [5] MK Bah, MM Gunshor, and TJ Schmit. Generation of goes-
550 16 true color imagery without a green band. *Earth and Space
551 Science*, 5(9):549–558, 2018. 3
- 552 [6] Jacob Buckman, Aurko Roy, Colin Raffel, and Ian Good-
553 fellow. Thermometer encoding: One hot way to resist ad-
554 versarial examples. In *International conference on learning
555 representations*, 2018. 5
- 556 [7] Marshall Burke, Anne Driscoll, Sam Heft-Neal, Jiani Xue,
557 Jennifer Burney, and Michael Wara. The changing risk and
558 burden of wildfire in the united states. *Proceedings of the
559 National Academy of Sciences*, 118(2):e2011048118, 2021.
560 1
- 561 [8] Abhishek Chaurasia and Eugenio Culurciello. Linknet: Ex-
562 ploiting encoder representations for efficient semantic seg-
563 mentation. In *2017 IEEE visual communications and image
564 processing (VCIP)*, pages 1–4. IEEE, 2017. 6
- 565 [9] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos,
566 Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image
567 segmentation with deep convolutional nets, atrous convolu-
568 tion, and fully connected crfs. *IEEE transactions on pattern
569 analysis and machine intelligence*, 40(4):834–848, 2017. 6
- 570 [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li,
571 and Li Fei-Fei. Imagenet: A large-scale hierarchical image
572 database. In *2009 IEEE conference on computer vision and
573 pattern recognition*, pages 248–255. Ieee, 2009. 2, 6
- 574 [11] Li Deng. The mnist database of handwritten digit images for
575 machine learning research [best of the web]. *IEEE signal
576 processing magazine*, 29(6):141–142, 2012. 2
- 577 [12] Tongle Fan, Guanglei Wang, Yan Li, and Hongrui Wang.
578 Ma-net: A multi-scale attention network for liver and tumor
579 segmentation. *IEEE Access*, 8:179656–179665, 2020. 6
- 580 [13] Rafael EP Ferreira, Yong Jae Lee, and João RR Dórea. Us-
581 ing pseudo-labeling to improve performance of deep neural
582 networks for animal identification. *Scientific Reports*, 13(1):
583 13875, 2023. 6
- 584 [14] Emmanuela Gakidou, Ashkan Afshin, Amanuel Alemu Aba-
585 jobir, Kalkidan Hassen Abate, Cristiana Abbafati, Kaja M
586 Abbas, Foad Abd-Allah, Abdishakur M Abdulle, Se-
587 maw Ferede Abera, Victor Aboyans, et al. Global, regional,
588 and national comparative risk assessment of 84 behavioural,
589 environmental and occupational, and metabolic risks or clus-
590 ters of risks, 1990–2016: a systematic analysis for the global
591 burden of disease study 2016. *The Lancet*, 390(10100):
592 1345–1422, 2017. 1
- 593 [15] Steven J Goodman, Timothy J Schmit, Jaime Daniels, and
594 Robert J Redmon. *The GOES-R series: a new generation of
595 geostationary environmental satellites*. Elsevier, 2019. 1
- 596 [16] Pavel Iakubovskii. Segmentation models pytorch. https://github.com/qubvel/segmentation_models.pytorch, 2019. 6
- 597 [17] J Jakubik, S Roy, C Phillips, P Fraccaro, D Godwin, B
598 Zadrozny, D Szwarcman, C Gomes, G Nyirjesy, B Edwards,
599 et al. Foundation models for generalist geospatial artificial
600 intelligence. arxiv 2023. *arXiv preprint arXiv:2310.18660*.
601 6
- 602 [18] Eric James, Ravan Ahmadov, and Georg A Grell. Real-
603 time wildfire smoke prediction in the united states: The hrrr-
604 smoke model. In *EGU General Assembly Conference Ab-
605 stracts*, page 19526, 2018. 1
- 606 [19] Asanobu Kitamoto, Jared Hwang, Bastien Vuillod, Lucas
607 Gautier, Yingtao Tian, and Tarin Clanuwat. Digital typhoon:
608 Long-term satellite image dataset for the spatio-temporal
609 modeling of tropical cyclones. *Advances in Neural Infor-
610 mation Processing Systems*, 36, 2024. 6
- 611 [20] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple
612 layers of features from tiny images. 2009. 2
- 613 [21] Alexandra Larsen, Ivan Hanigan, Brian J Reich, Yi Qin,
614 Martin Cope, Geoffrey Morgan, and Ana G Rappold. A
615 deep learning approach to identify smoke plumes in satel-
616 lite imagery in near-real time for health risk communication.
617 *Journal of exposure science & environmental epidemiology*,
618 31(1):170–176, 2021. 2, 3
- 619 [22] Dong-Hyun Lee. Pseudo-label : The simple and efficient
620 semi-supervised learning method for deep neural networks.
621 *ICML 2013 Workshop : Challenges in Representation Learn-
622 ing (WREPL)*, 2013. 2
- 623 [23] Yoonjin Lee, Christian D Kummerow, and Imme Ebert-
624 Uphoff. Applying machine learning methods to detect
625 convection using geostationary operational environmental
626 satellite-16 (goes-16) advanced baseline imager (abi) data.
627 *Atmospheric Measurement Techniques*, 14(4):2699–2716,
628 2021. 3
- 629 [24] Donna McNamara, George Stephens, Mark Ruminski, and
630 Tim Kasheta. The hazard mapping system (hms) - noaa's
631 multi-sensor fire and smoke detection program using envi-
632 ronmental satellites. *Conference on Satellite Meteorology
633 and Oceanography*, 2004. 2
- 634 [25] NOAA. Hazard mapping system fire and smoke product. 2
- 635 [26] Thanh Cong Phan and Thanh Tam Nguyen. Remote sens-
636 ing meets deep learning: exploiting spatio-temporal-spectral
637 satellite images for early wildfire detection. 2019. 3
- 638 [27] T Randriambelo, S Baldy, M Bessafi, Michel Petit, and Marc
639 Despinoy. An improved detection and characterization of
640 active fires and smoke plumes in south-eastern africa and
641 madagascar. *International Journal of Remote Sensing*, 19
642 (14):2623–2638, 1998. 2
- 643 [28] Martin Raspaud, David Hoese, Adam Dybbroe, Panu Lahti-
644 nen, Abhay Devasthale, Mikhail Itkin, Ulrich Hamann,
645 Lars Ørum Rasmussen, Esben Stigård Nielsen, Thomas Lep-
646 pelt, et al. Pytroll: An open-source, community-driven
647 python framework to process earth observation satellite data.
648
- 649

- 650 *Bulletin of the American Meteorological Society*, 99(7):
 651 1329–1336, 2018. 3
- 652 [29] Alain Royer, Pierre Vincent, and Ferdinand Bonn. Evaluation
 653 and correction of viewing angle effects on satellite mea-
 654 surements of bidirectional reflectance. *Photogrammetric en-*
 655 *gineering and remote sensing*, 51(12):1899–1914, 1985. 5
- 656 [30] W Schroeder, M Ruminski, I Csiszar, L Giglio, E Prins, C
 657 Schmidt, and J Morisette. Validation analyses of an oper-
 658 ational fire monitoring product: The hazard mapping sys-
 659 tem. *International Journal of Remote Sensing*, 29(20):6059–
 660 6066, 2008. 2
- 661 [31] Ali Sharif Razavian, Hossein Azizpour, Josephine Sullivan,
 662 and Stefan Carlsson. Cnn features off-the-shelf: an astound-
 663 ing baseline for recognition. In *Proceedings of the IEEE con-*
 664 *ference on computer vision and pattern recognition work-*
 665 *shops*, pages 806–813, 2014. 6
- 666 [32] Bjorn Stevens, Sandrine Bony, Hélène Brogniez, Laureline S
 667 Hentgen, Cathy Hohenegger, Christoph Kiemle, Tristan S
 668 L’Ecuyer, Ann Kristin Naumann, Hauke Schulz, Pier A
 669 Siebesma, et al. Sugar, gravel, fish and flowers: Mesoscale
 670 cloud patterns in the trade winds. *Quarterly Journal of the*
 671 *Royal Meteorological Society*, 146(726):141–152, 2020. 6
- 672 [33] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhi-
 673 nav Gupta. Revisiting unreasonable effectiveness of data in
 674 deep learning era. In *Proceedings of the IEEE interna-*
 675 *tional conference on computer vision*, pages 843–852, 2017. 2
- 676 [34] Mingxing Tan and Quoc Le. Efficientnetv2: Smaller models
 677 and faster training. In *International conference on machine*
 678 *learning*, pages 10096–10106. PMLR, 2021. 6
- 679 [35] Zewei Wang, Pengfei Yang, Haotian Liang, Chang Zheng,
 680 Jiyuan Yin, Ye Tian, and Wenbin Cui. Semantic segmentation
 681 and analysis on sensitive parameters of forest fire smoke us-
 682 ing smoke-unet and landsat-8 imagery. *Remote Sensing*, 14
 683 (1):45, 2022. 2, 3
- 684 [36] Jeff Wen and M. Burke. Wildfire smoke plume seg-
 685 mentation using geostationary satellite imagery. *ArXiv*,
 686 abs/2109.01637, 2021. 2, 3
- 687 [37] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson.
 688 How transferable are features in deep neural networks? *Ad-*
 689 *vances in neural information processing systems*, 27, 2014.
 690 6
- 691 [38] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang
 692 Wang, and Jiaya Jia. Pyramid scene parsing network. In
 693 *Proceedings of the IEEE conference on computer vision and*
 694 *pattern recognition*, pages 2881–2890, 2017. 6
- 695 [39] Tom X.-P. Zhao, Steve Ackerman, and Wei Guo. Dust and
 696 smoke detection for multi-channel imagers. *Remote Sensing*,
 697 2(10):2347–2368, 2010. 2
- 698 [40] Stefano Zorzi, Shabab Bazrafkan, Stefan Habenschuss, and
 699 Friedrich Fraundorfer. Polyworld: Polygonal building ex-
 700 traction with graph neural networks in satellite images. In
 701 *Proceedings of the IEEE/CVF Conference on Computer Vi-*
 702 *sion and Pattern Recognition*, pages 1848–1857, 2022. 6