
SmokeViz: Using Pseudo-Labels to Develop a Deep Learning Dataset of Wildfire Smoke Plumes in Satellite Imagery

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The increase in the frequency of wildfires on a global scale underscores the need
2 for advancements in fire monitoring techniques for disaster management, environmental
3 protection and to mitigate negative health outcomes. This research
4 introduces an innovative, data-driven framework that leverages the semi-supervised
5 method, pseudo-labeling, to generate smoke plume annotations in geostationary
6 satellite imagery. The primary objective is to refine an existing National Oceanic
7 and Atmospheric Administration smoke dataset that provides temporal and geo-
8 graphical information on individual smoke plumes but at variable and, primarily,
9 low temporal resolution. To do this, we use deep learning and pseudo-labels to
10 pinpoint the singular, most representative, satellite image that optimally illustrates
11 the smoke annotation within the given time window. By identifying the most
12 representative imagery of smoke plumes for a given smoke annotation, the study
13 seeks to create an accurate and relevant machine learning dataset. The resulting
14 dataset is anticipated to be an instrumental tool in developing further machine
15 learning models, such as an automated system capable of real-time monitoring and
16 annotation of smoke plumes directly from streaming satellite imagery.

17

1 Introduction

18 In recent years, the escalation of wildfire incidents worldwide has become a prominent environmental
19 and public health concern. The combustion process in wildfires releases smoke containing fine
20 particulate matter (PM2.5) and harmful gases, posing severe hazards to human health and air quality.
21 These risks underscore the necessity for efficient and effective monitoring methods to mitigate the
22 adverse health impacts associated with wildfire smoke.

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

29 Satellite imagery, equipped with advanced sensors, such as the Advanced Baseline Imager (ABI) on
30 the Geostationary Operational Environmental Satellites (GOES), have revolutionized environmental
31 monitoring. These tools enable the detailed observation of smoke plumes, their particulate density,
32 and the extent of smoke spread. These satellite-based systems offer the capabilities to provide critical
33 insights into the concentration and movement of smoke particulates, facilitating accurate and timely
34 assessments of air quality.

35 The integration of satellite imagery in wildfire smoke monitoring is not only instrumental in providing
 36 real-time data but also plays a significant role in public health planning and response. By mapping
 37 the spread and density of smoke, health authorities can issue timely warnings, implement evacuation
 38 protocols, and deploy resources effectively to mitigate health risks. Furthermore, long-term data
 39 gathered from satellite observations can aid in understanding the broader impacts of wildfire smoke
 40 on public health, influencing policy decisions and preventive measures.
 41 Currently, multi-channel thresholding is a popular method to distinguish smoke pixels from pixels
 42 containing dust, clouds or other phenomenon with similar signatures [24]. The method uses historical,
 43 labeled data to extract optimal radiance values for each channel that corresponds with the labeled
 44 class. These methods are tuned to particular biogeographies and often have issues with generalization
 45 to new locations with varying fuel types [13].
 46 In contrast to the numerical thresholding approach, human visual inspection of satellite imagery is
 47 another commonly used method for smoke identification. Trained analyst will inspect imagery and
 48 label the smoke by hand. This method is not as scalable as an automated approach and is limited by
 49 the availability of analysts and their time.
 50 To address these challenges we can look towards innovative approaches and technological advancements
 51 in computer vision. Machine learning methods have shown potential in improving the accuracy
 52 and efficiency of satellite-based wildfire smoke detection and monitoring. For instance, SmokeNet,
 53 uses a convolutional neural network (CNN) based framework to determine if a scene of MODIS
 54 imagery contains smoke [1]. Another study also used a CNN to identify smoke on a pixel-wise basis
 55 using imagery from Himiwiari-8 [8]. Additionally, Wen et al. developed a CNN architecture that
 56 takes GOES-East imagery as input and National Oceanic and Atmospheric Administration (NOAA)
 57 generated annotations for the target labels during training [22].
 58 The success of deep learning methods, such as CNNs, relies heavily on the availability of a large,
 59 representative dataset [20]. As laid out in table 1, existing methods use relatively small number of
 60 samples, from 57 [21] to 6825 [22], where one sample represents a satellite image with a singular time
 61 and geolocation. In contrast, benchmark datasets for image classification contain tens of thousands
 62 (CIFAR-10 and MNIST) to millions (CIFAR-100 and ImageNet) of data samples. Keeping in mind
 63 the correlation between both the quality and quantity of data with model performance, we introduce
 64 the largest known smoke dataset, SmokeViz, containing over 120,000 samples.

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

Reference	Method	# Samples	Satellite	# Channels	Level
[1]	CNN	6255	MODIS	5	image
[22]	CNN	6825	GOES-East	5	pixel
[8]	CNN	975	Himiwari-8	7	pixel
[21]	U-Net	47	Landsat-8	13	pixel
SmokeViz	U-Net	120,000	GOES-East/West	3	pixel

65 An approach to increase the number of labeled samples in a dataset, semi-supervised learning
 66 leverages a labeled dataset to generate new labels for an often larger, but unlabeled, dataset. Pseudo-
 67 labeling, a form of semi-supervised learning, uses labeled data to train an initial model, then runs
 68 that model on unlabeled data to predict pseudo-labels, and finally trains a new model using the
 69 pseudo-labels [9]. We introduce a variation of pseudo-labeling not to increase the size, but to increase
 70 the quality of our dataset by using pseudo-labels to choose the best satellite image out of a given
 71 time-window to represent each smoke plume annotation.

72 2 Methods

73 Dataset

74 The initial data source, discussed in further detail in the next section, is uniquely characterized by
 75 each annotation having corresponding imagery ranging between 1-60 frames, where each frame
 76 captures 5 minutes of exposure. Additionally, we have two satellites that overlap in coverage area,

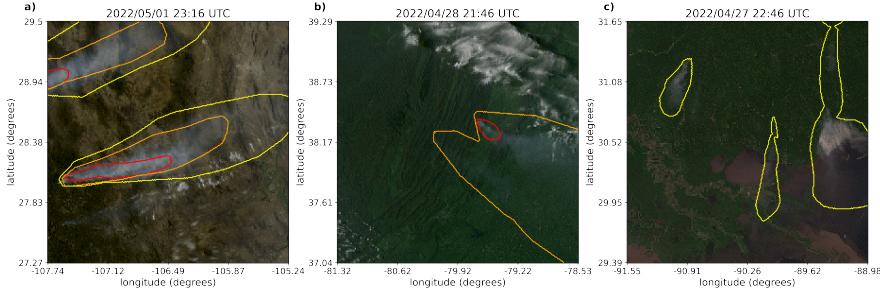


Figure 1: Satellite imagery captured by GOES-East within a few days of each other. The yellow, orange and red contours indicate the extent of Light, Medium and Heavy smoke. a) shows a canonical example of a smoke plume. b) and c) show variations in the density labels.

77 GOES-East and GOES-West, effectively doubling the number of frames for a single annotation. We
 78 apply pseudo-labeling to develop a dataset that has a one-to-one annotation-to-image ratio, where we
 79 choose the satellite image that has the maximum overlap between the geolocation of smoke in the
 80 imagery and the analyst annotation.

81 Dataset development came in three stages. First, we use the physics of light scattering to determine
 82 which singular satellite image would be in the optimal configuration for smoke detection. Second, we
 83 used that dataset to train an initial model that will identify smoke in satellite imagery. Third, we use
 84 that initial model to label each satellite image in a given annotation’s time-window and the optimal
 85 satellite image is chosen based on which image’s pseudo-labels has the greatest overlap with the
 86 analyst annotation for the given location and densities of smoke.

87 **Smoke Labels**

88 NOAA manages environmental satellite programs such as the Hazard Mapping System (HMS)
 89 [11, 19]. The HMS program is an operational system that uses an aggregation of satellite data to
 90 generate active fire and smoke data that is used in applications such as air quality assessments and
 91 serves as verification and validation for NOAA’s smoke forecasting model, HYSPLiT, [16]. To train
 92 our model, we implement a supervised learning framework that uses the HMS analyst smoke product
 93 as truth labels during the model training process.

94 HMS smoke analysis data gives the coordinates of the smoke perimeter and classifies the smoke by
 95 density within a given time window. The time windows can range from instantaneous (same start/end
 96 time) to lengths of 5 hours. While the bounds of the smoke annotations can change within the larger
 97 time spans, the analyst is making an approximation that should reflect the smoke coverage over the
 98 duration of the window. The density information is qualitatively determined by the analyst based on
 99 smoke opacity and categorized as either light, medium or heavy as seen in figure 1a.

100 **Thermometer Encoding Smoke Densities**

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

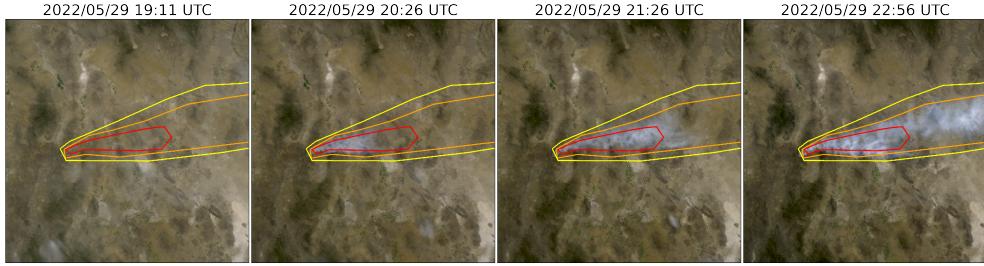


Figure 2: True Color GOES-East imagery from May 2022, Southeast New Mexico (31°N , 100°W) during the start of the Foster Fire. The red, orange and yellow lines represent the heavy, medium and low density HMS smoke annotations that span 19:10–23:00 UTC.

114 In addition to the densities being ordered and categorical, the differences between the density
 115 categories are not evenly distributed by a given metric, such as particulate matter per square meter.
 116 The intervals between densities being unknown along with the hierarchical nature of the density labels
 117 makes the labels ordinal instead of just categorical. This data property allows us to use thermometer
 118 encoding [3], which leverages the idea that heavy density smoke includes both medium and light
 119 density smoke, that heavy density smoke is closer to medium than it is to light and automatically
 120 weights the loss functions and incorporates the ranked ordering of the densities. As seen in Table 2,
 121 one-hot encoding, commonly used for categorical data, doesn't take ordinal properties of the data
 122 into consideration.

Table 2: 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]

123 Time Windows For Smoke Annotations

124 In order to take into account movement characteristics to help identify smoke, analysts use multi-
 125 frame animations of the satellite imagery. The resulting annotations often have large time windows
 126 over multiple hours to represent one smoke plume annotation. Since the goal of these annotations is
 127 to show the general coverage over that time span, as shown in figure 2, the smoke boundaries don't
 128 often match up with the satellite imagery over the entire time window. One way to approach this
 129 problem would be to use all the satellite images the analysts used as input. Since the timespans are
 130 non-uniform, this would vary the length in imagery inputs into the model, which would be difficult
 131 with a CNN architecture. Moreover, this would require a large amount of additional memory and
 132 computational resources. Instead of using the original analysts' many satellite image inputs to one
 133 annotated output, we develop a one-to-one input-to-output by finding the optimal singular satellite
 134 image input to represent the annotation. As discussed in the next section, we do this by making
 135 physics-driven choices on which satellite and timestamp would give the optimal angle between the
 136 sun and satellite that would produce the strongest smoke signature for the geolocation and timestamp
 137 of the smoke plume.

138 Satellite Imagery

139 The Geostationary Operational Environmental Satellites (GOES) are operated by the NOAA in order
 140 to support meteorology research and forecasting for the United States. We use the latest operational
 141 satellites, GOES-16 (East), 17 and 18 (West) that carry the ABI, that measure 16 bands between the

Table 3: To create a true color image, we use the following bands from the ABI Level 1b CONUS (ABI-L1b-RadC) product.

band	description	center wavelength	spatial resolution (km)
C01	blue visible	0.47	1
C02	red visible	0.64	0.5
C03	veggie near infrared	0.865	1

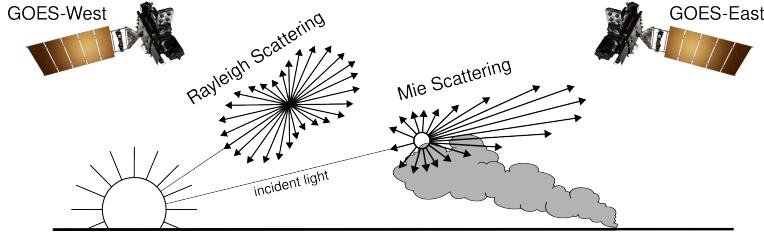


Figure 3: If the particle size is $< \frac{1}{10}$ the wavelength of the interacting light, then the primary scattering will be Rayleigh. Mie scattering is the predominant scattering mechanism when the particle size is larger than wavelength of light. This schematic demonstrates that when the sun is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

142 visible and infrared wavelengths. In improvement to the GOES predecessors, imagery is collected
 143 every 5 minutes for the contiguous United States and every 10 minutes for the full disk. We use bands
 144 1-3 (Table 3) as input to Satpy’s composite algorithm to develop a true color image representation,
 145 similar to what is used as input by HMS analysts [14] and [2].

146 **Mie-Derived Dataset**

147 We used a physics-informed approach in selecting the initial dataset, \mathcal{D}_M , we call the Mie-derived
 148 dataset, for training an initial parent model, f_p . Prior GOES ABI datasets for machine learning
 149 applications often include data from only one of the two GOES-series satellites, commonly opting
 150 for GOES-East [22], [12], [10]. Rather than using one satellite or the cumulative data from both
 151 GOES-West and GOES-East images, we select between one or the other based on the solar zenith
 152 angle. For smoke identification, this approach can achieve a much higher signal-to-noise than imaging
 153 the earth’s surface from an arbitrary angle. The elastic scattering of light is the primary mechanism
 154 to account for - while the atmosphere is composed of molecules with size $< 1\text{nm}$, smoke particles
 155 can vary from $100\text{ nm} - 10\text{ }\mu\text{m}$ in diameter, d . The GOES ABI covers spectral bands from $0.47\text{ }\mu\text{m} -$
 156 $13.3\text{ }\mu\text{m}$, so atmospheric and smoke particle sizes occupy two very different regimes with respect
 157 to the imaging wavelength λ . In the extreme limit of $\lambda \gg d$, the physics of scattering of light off a
 158 small sphere is captured by Rayleigh scattering. This process has two critical consequences: (1) the
 159 scattering cross section of light is strongly wavelength dependent (scaling with λ^{-4}), meaning that
 160 photons with wavelength closer to the ultraviolet are scattered more strongly than infrared photons. (2)
 161 the scattering cross section scales with an angular dependent cross section of $(1 + \cos^2 \theta)$. Scattered
 162 photons follow the emission distribution of a radiating dipole, scattering more strongly in the forward
 163 and backwards directions ($\theta = 0, \pi$) than orthogonal to the direction of propagation ($\theta = \pi/2, 3\pi/2$),
 164 see figure 3 for Rayleigh scattering schematic.

165 The significance of these scalings is that the observer, or detector, will receive blue photons in most
 166 directions orthogonal to the source. Equivalently, photons traveling colinearly with line of sight to
 167 the emission source will mostly have wavelengths in the infrared band. In the converse regime of
 168 $d > \lambda$, the elastic scattering of light against matter is modeled through Mie scattering. In comparison
 169 to Rayleigh scattering, Mie scattering is largely wavelength independent and has a more complicated
 170 radiation pattern where the cross section has a maximal amplitude in the forward direction. An
 171 observer downstream of this scatterer will collect more photons than one positioned directly behind it.
 172 In the context of smoke identification, a sunrise or sunset will lead to a higher Mie scattered signal in
 173 GOES-West and GOES-East respectively, as shown with a smoke plume producing a stronger signal
 174 in GOES-East imagery near sunset in figure 2.

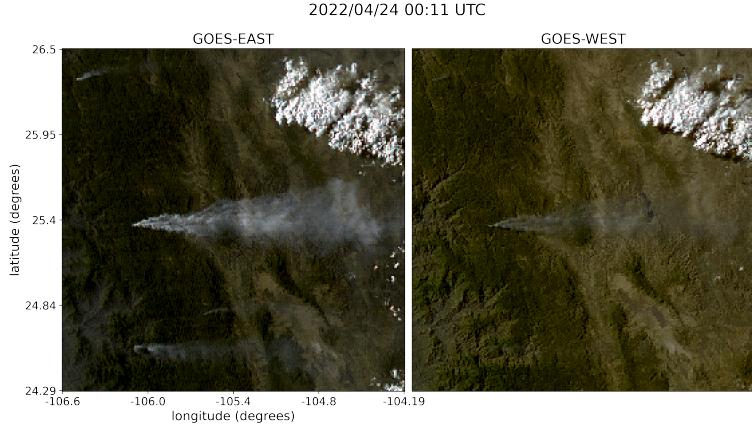


Figure 4: True Color GOES-East (left) and GOES-West (right) imagery from April 24th, 2022 in Durango, Mexico. The images were taken ~ 0.5 hours before sunset (01:43 UTC) for this geolocation and time of year.

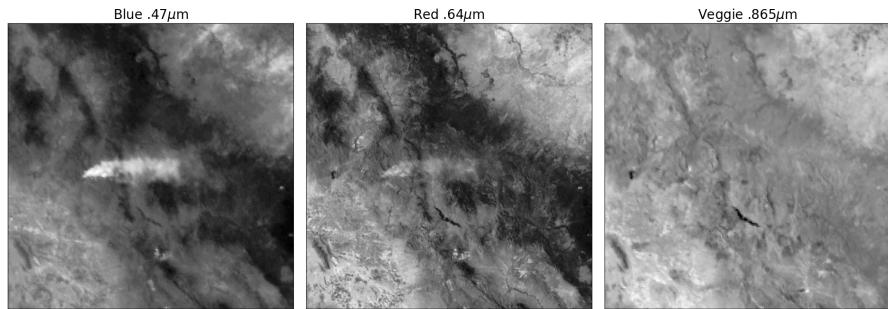


Figure 5: Three bands of GOES-East data are the raw input to generate the True Color image shown in figure 4. These plots show variations in signal-to-noise ratio for smoke detection in relation to the wavelength, λ , of light being measured.

175 Smoke identification therefore amounts to extracting a signal of $d > \lambda$ photons from the $\lambda \gg d$
 176 background. Positioning a detector along line of sight to the scatterer will result in a higher signal
 177 from smoke particles (figure 3). Filtering the imaged wavelength can enhance this signal; photons
 178 collected in the blue spectrum will have a naturally lower background along the line of sight to the
 179 illumination source do their high level of Rayleigh scattering as. Therefore, as demonstrated in figure
 180 5, this configuration results in the highest signal to noise imaging for smoke particles.

181 Based solely on these criteria, the optimal strategy would be to pull data from GOES-West right after
 182 sunrise and from GOES-East right before sunset. Another factor to consider is that the time when the
 183 sun is in optimal alignment with the satellite for smoke detection coincides with when solar zenith
 184 angle is maximized. Larger angles between the satellite and sun result in an increase in noise due
 185 to increased atmospheric interactions [18]. This is shown in figure 6, while we optimize for smoke
 186 signal detection, due to the high solar zenith angle, we introduce atmospheric interaction noise that
 187 obfuscate the smoke signal. To reduce the noise from large solar zenith angles, if given multiple
 188 options to choose from, we choose the image with the largest solar zenith angle that is below 80
 189 degrees.

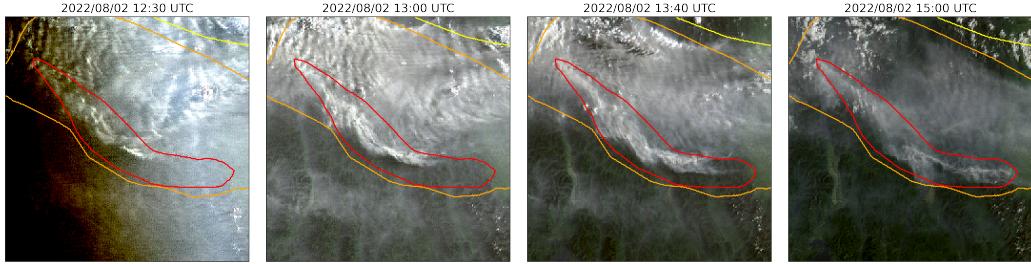


Figure 6: A smoke annotation projected onto GOES-West imagery from August 2022 that spans from 11:00 UTC to 15:00 UTC, sunrise on August 2nd, 2022 at coordinates ($49^{\circ}24'N$, $115^{\circ}29'W$) was 12:15 UTC.

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

197 Machine Learning Model

198 We implement a deep learning architecture that uses the encoder from the ResNet model [6] and a
 199 semantic segmentation classifier from the U-Net model [17]. Transfer learning has shown to reduce
 200 the time and resources needed to train a model by leveraging information from pre-trained models
 201 [23], [15]. We initialize the values of our model weights using the pre-trained values originally
 202 trained on the ImageNet dataset [4], containing 1.2 million images and 1000 categories. Our model
 203 was developed using the Segmentation Models PyTorch package [7] that was written as a high level
 204 API for implementing models for semantic segmentation problems. We input 256x256x3 snapshots
 205 of True Color GOES imagery that contains smoke and output a 256x256x3 classification map that
 206 predicts if a pixel contains smoke and if so, what the density of that smoke is. As mentioned earlier,
 207 we apply the thermometer encoding shown in table 2 to encode the smoke densities and apply binary
 208 cross entropy as the loss function per density of smoke.

209 The \mathcal{D}_M dataset contained over 120,000 samples. To train f_p , we split \mathcal{D}_M into training (95,000
 210 samples), validation (12,000 samples) and testing (12,000) datasets. Training data contains data from
 211 the years 2018, 2019, 2020, 2021 and 2023 while the data from 2022 is split into validation and
 212 testing data by taking data from alternating 10 days of the year. In order to make sure we include
 213 the monthly variations in wildfire trends over a full year, we split 2022 data up by every 10 days.
 214 This allowed us to: (1) allocate an additional full year of data for the training set, (2) show yearlong
 215 trends in both the validation and testing sets and (3) keep the validation and testing datasets relatively
 216 independent from one another.

217 We trained the parent model, f_p , for 10 epochs, then ran f_p on all images, x_N , within the analyst
 218 time-window for each annotation to select image that's pseudo-label best matched the HMS smoke
 219 annotation, y^a . An image, x^* , would have the potential be included in \mathcal{D}_{PL} only if it generated the
 220 highest Intersection over Union (IoU) value between the image's l^* and y^a over all x_N . The IoU
 221 metric is given by the ratio of area of overlap to the area of union as shown in equation 1.

$$IoU = \frac{|y^a \cap l^*|}{|y^a| \cup |l^*|} \quad (1)$$

222 To determine which image, x , out of the relevant imagery, x_N , for the given time window best
 223 represents the analyst annotation, y^a , we run f_p on each x to generate a pseudo-label, l^* . The output

Table 4: IoU results per density of smoke and over all densities.

Density	f_p		f_c	
	\mathcal{D}_M	\mathcal{D}_{PL}	\mathcal{D}_M	\mathcal{D}_{PL}
Light	0.394	0.551	0.418	0.538
Medium	0.283	0.392	0.340	0.411
Heavy	0.233	0.290	0.270	0.325
Overall	0.365	0.510	0.396	0.503

Table 5: IoU results per density of smoke and over all densities.

Density	IoU \mathcal{D}_M	IoU \mathcal{D}_{PL}
Light	0.394	0.551
Medium	0.283	0.392
Heavy	0.233	0.290
Overall	0.365	0.510

of f_p , l^* , 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. l^* serve as pseudo-labels for each density of smoke and are compared to the analyst annotations, y^a . To compare l^* and y^a , we calculate the IoU using the total set of pixels for l^* at that density of smoke and the entire set of pixels for y^a for a particular smoke density in each image. The image with the highest IoU score is chosen as the image, x^* , that best represents the analyst smoke annotation, y^a . Often used for pseudo-labeling, a confidence threshold value is defined to determine if a pseudo-label should be included in a dataset [5]. We chose a confidence threshold that would include the sample, x^* , in \mathcal{D}_{PL} if the maximum overall IoU (equation 2) between l^* and x^a over all densities was over 0.1.

Finally, we use \mathcal{D}_{PL} to train an additional child model, f_c . We use the same dataset split method and model setup but change \mathcal{D}_M to \mathcal{D}_{PL} to train the model over 10 epochs.

Results

To interpret the performance of f_p , we report the IoU metrics in table 5 that were computed by running f_p and f_c on \mathcal{D}_M and \mathcal{D}_{PL} . For each density, we calculate the IoU using the total set of pixels that f_p 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.

$$IoU_{overall} = \frac{\sum_{i=light}^{heavy} |y_i^a \cap l_i^*|}{\sum_{i=light}^{heavy} |y_i^a| \cup |l_i^*|} \quad (2)$$

An illustration of an improvement in the dataset is evident in Figure 7 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 Mei-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 result of this study is a representative dataset that can be used to train machine learning models for various wildfire smoke applications. The end 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 monitoring and as data provided to public health officials for air quality assessments.

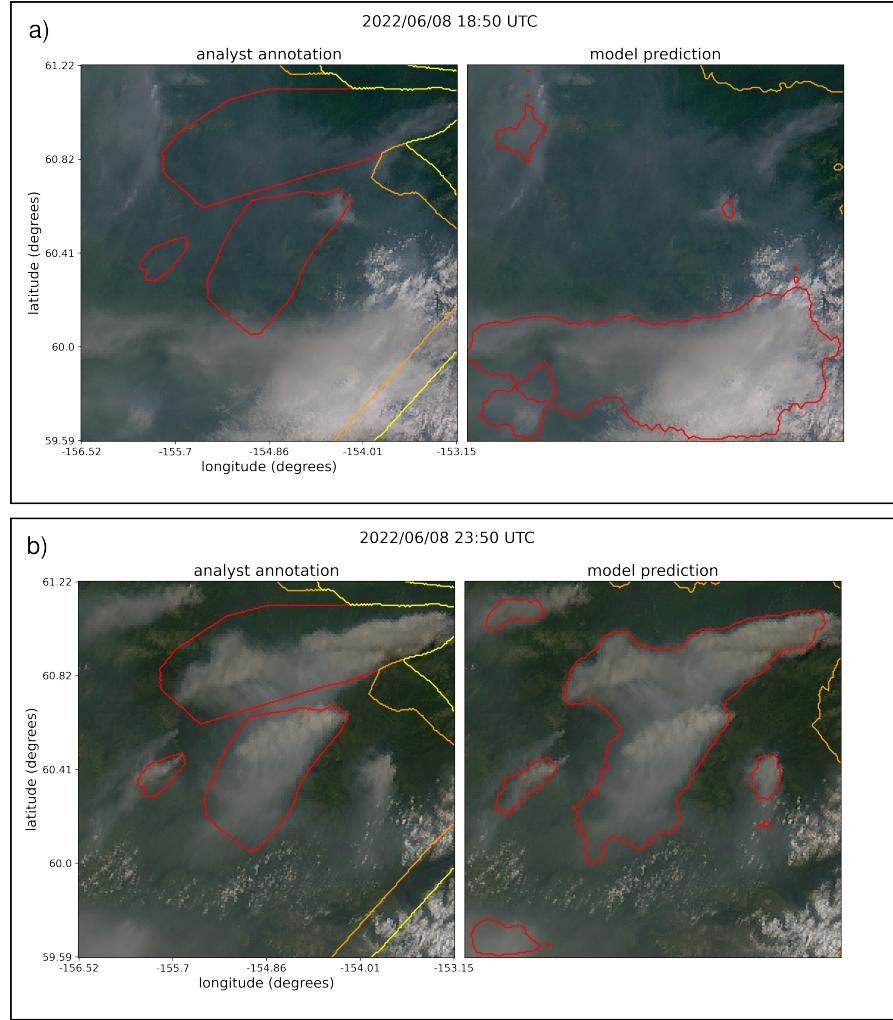


Figure 7: GOES-West imagery showing smoke on June 8th, 2022 in Alaska where, at the coordinates (61°03'N, 156°07'W), daylight was between 12:43-7:53 UTC. The HMS smoke annotations displayed span from 18:50 to 23:50 UTC. a) shows the imagery that was selected using the Mie-derived data selection process b) shows the image that had the highest IoU score between the pseudo-label and the analyst annotation.

252 3 Limitations

253 By selecting for imagery data based on which image had the highest IoU score, it could be anticipated
 254 that the results in shown in table 5 would show higher IoU scores for the pseudo-labeled dataset. We
 255 expect that a model trained on a dataset that better represents the truth labels would perform better.

256 4 Conclusion

257 In this study, we have refined an existing dataset originally curated by NOAA's HMS team, trans-
 258 forming it from a many-to-one imagery-to-annotation format to a, more succinct, one-to-one satellite
 259 image-to-annotation dataset. The initial HMS dataset primarily provided a general approximation
 260 of where smoke had been present for a given time window, though it did not guarantee the actual
 261 existence of smoke in the labeled pixels during the given times. Our goal was to create a dataset
 262 that could be used, along with additional applications, to train a model to detect wildfire smoke in
 263 real-time on an image-by-image level. The Mie-derived dataset seleciton process determines that if

264 smoke is present, what timestamp within the analyst timewindow would give the highest smoke
265 signal-to-noise ratio. While optimizing for being able to detect smoke, if it is present, the Mie-dataset
266 selection had no metric to determine if the smoke was effectively present in the selected image. Since
267 many of the images within the HMS time-window either contained no smoke at all or the smoke was
268 not contained within the geospatial bounds of the annotations, the Mie-derived dataset contained
269 a large number of mislabeled samples. Discrepancies between data and labels can be detrimental
270 towards the model's capacity to improve on feature representations in the target domain. During
271 model training, the penalization of accurate predictions can inadvertently introduce biases towards
272 misclassifying noise as meaningful signal.

273 To improve the dataset's capacity to accurately represent wildfire smoke plumes, we train a machine
274 learning model using the Mie-derived dataset and run it on the relevant satellite images within the
275 time-frame. The image with the maximum IoU score between the model's smoke predictions, or
276 pseudo-label, and the analyst smoke annotations are used to create the pseudo-label dataset.

277 5 Acknowledgments and Disclosure of Funding

278 This work was partially supported by the NOAA Global Systems Laboratory and Cooperative Institute
279 for Research in Environmental Sciences at the University of Colorado Boulder. We thank Wilfrid
280 Schroeder and the Hazard Mapping Systems team for giving guidance on how they created their
281 smoke plume dataset. This work utilized the Alpine high performance computing resource at the
282 University of Colorado Boulder. Alpine is jointly funded by the University of Colorado Boulder, the
283 University of Colorado Anschutz, Colorado State University, and the National Science Foundation
284 (award 2201538).

285 References

- 286 [1] R. Ba, C. Chen, J. Yuan, W. Song, and S. Lo. Smokenet: Satellite smoke scene detection using
287 convolutional neural network with spatial and channel-wise attention. *Remote Sensing*, 11(14):
288 1702, 2019.
- 289 [2] M. Bah, M. Gunshor, and T. Schmit. Generation of goes-16 true color imagery without a green
290 band. *Earth and Space Science*, 5(9):549–558, 2018.
- 291 [3] J. Buckman, A. Roy, C. Raffel, and I. Goodfellow. Thermometer encoding: One hot way to
292 resist adversarial examples. In *International conference on learning representations*, 2018.
- 293 [4] J. Deng, K. Li, M. Do, H. Su, and L. Fei-Fei. Construction and Analysis of a Large Scale Image
294 Ontology. Vision Sciences Society, 2009.
- 295 [5] R. E. Ferreira, Y. J. Lee, and J. R. Dórea. Using pseudo-labeling to improve performance of
296 deep neural networks for animal identification. *Scientific Reports*, 13(1):13875, 2023.
- 297 [6] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition, 2015.
- 298 [7] P. Iakubovskii. Segmentation models pytorch. https://github.com/qubvel/segmentation_models.pytorch, 2019.
- 300 [8] A. Larsen, I. Hanigan, B. J. Reich, Y. Qin, M. Cope, G. Morgan, and A. G. Rappold. A deep
301 learning approach to identify smoke plumes in satellite imagery in near-real time for health risk
302 communication. *Journal of exposure science & environmental epidemiology*, 31(1):170–176,
303 2021.
- 304 [9] D.-H. Lee. Pseudo-label : The simple and efficient semi-supervised learning method for deep
305 neural networks. *ICML 2013 Workshop : Challenges in Representation Learning (WREPL)*, 07
306 2013.
- 307 [10] Y. Lee, C. D. Kummerow, and I. Ebert-Uphoff. Applying machine learning methods to detect
308 convection using geostationary operational environmental satellite-16 (goes-16) advanced
309 baseline imager (abi) data. *Atmospheric Measurement Techniques*, 14(4):2699–2716, 2021.

- 310 [11] D. McNamara, G. Stephens, M. Rumsby, and T. Kasheta. The hazard mapping system (hms) -
 311 noaa's multi-sensor fire and smoke detection program using environmental satellites. *Conference*
 312 *on Satellite Meteorology and Oceanography*, 01 2004.
- 313 [12] T. C. Phan and T. T. Nguyen. Remote sensing meets deep learning: exploiting spatio-temporal-
 314 spectral satellite images for early wildfire detection. 2019.
- 315 [13] T. Randriambelo, S. Baldy, M. Bessafi, M. Petit, and M. Despinoy. An improved detection
 316 and characterization of active fires and smoke plumes in south-eastern africa and madagascar.
 317 *International Journal of Remote Sensing*, 19(14):2623–2638, 1998.
- 318 [14] M. Raspaud, D. Hoesel, A. Dybbroe, P. Lahtinen, A. Devasthale, M. Itkin, U. Hamann, L. Ø.
 319 Rasmussen, E. S. Nielsen, T. Leppelt, et al. Pytroll: An open-source, community-driven python
 320 framework to process earth observation satellite data. *Bulletin of the American Meteorological
 321 Society*, 99(7):1329–1336, 2018.
- 322 [15] A. S. Razavian, H. Azizpour, J. Sullivan, and S. Carlsson. Cnn features off-the-shelf: an
 323 astounding baseline for recognition, 2014.
- 324 [16] G. D. Rolph, R. R. Draxler, A. F. Stein, A. Taylor, M. G. Rumsby, S. Kondragunta, J. Zeng,
 325 H.-C. Huang, G. Manikin, J. T. McQueen, et al. Description and verification of the noaa smoke
 326 forecasting system: the 2007 fire season. *Weather and Forecasting*, 24(2):361–378, 2009.
- 327 [17] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image
 328 segmentation, 2015.
- 329 [18] A. Royer, P. Vincent, and F. Bonn. Evaluation and correction of viewing angle effects on
 330 satellite measurements of bidirectional reflectance. *Photogrammetric engineering and remote
 331 sensing*, 51(12):1899–1914, 1985.
- 332 [19] W. Schroeder, M. Rumsby, I. Csizsar, L. Giglio, E. Prins, C. Schmidt, and J. Morisette.
 333 Validation analyses of an operational fire monitoring product: The hazard mapping system.
 334 *International Journal of Remote Sensing*, 29(20):6059–6066, 2008.
- 335 [20] C. Sun, A. Shrivastava, S. Singh, and A. Gupta. Revisiting unreasonable effectiveness of data in
 336 deep learning era, 2017.
- 337 [21] Z. Wang, P. Yang, H. Liang, C. Zheng, J. Yin, Y. Tian, and W. Cui. Semantic segmentation and
 338 analysis on sensitive parameters of forest fire smoke using smoke-unet and landsat-8 imagery.
 339 *Remote Sensing*, 14(1):45, 2022.
- 340 [22] J. Wen and M. Burke. Wildfire smoke plume segmentation using geostationary satellite imagery.
 341 *ArXiv*, abs/2109.01637, 2021. URL [https://api.semanticscholar.org/CorpusID:
 342 237416777](https://api.semanticscholar.org/CorpusID:237416777).
- 343 [23] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural
 344 networks?, 2014.
- 345 [24] T. X.-P. Zhao, S. Ackerman, and W. Guo. Dust and smoke detection for multi-channel imagers.
 346 *Remote Sensing*, 2(10):2347–2368, 2010. ISSN 2072-4292. doi: 10.3390/rs2102347. URL
 347 <https://www.mdpi.com/2072-4292/2/10/2347>.