
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 global increase in the frequency and intensity of wildfires underscores the
2 need for advancements in fire monitoring techniques. In order to investigate
3 deep learning approaches for detecting and tracking wildfires and the related
4 human health impacts, we present SmokeViz, a large scale machine learning
5 dataset of smoke plumes in satellite imagery. To build the dataset, we refine a
6 set of human-generated annotations created by analysts at the National Oceanic
7 and Atmospheric Administration. Each annotation gives a general temporal and
8 geographical approximation of smoke plumes but at variable and, primarily, low
9 temporal resolution. We present an innovative solution for refining the temporal and
10 spatial resolution in the given analyst annotations by leveraging the semi-supervised
11 method, pseudo-labeling. Unlike typical pseudo-labeling applications that aim to
12 increase the number of labeled samples, the objective is to use pseudo-labels to
13 refine an existing but coarse-grained set of annotations. We train a deep learning
14 model to generate pseudo-labels that pinpoint the singular, most representative,
15 satellite image to match the smoke annotation within the given temporal range. By
16 identifying the most representative imagery of smoke plumes for a given smoke
17 annotation, the study seeks to create an accurate and relevant machine learning
18 dataset. The resulting SmokeViz dataset is anticipated to be an instrumental tool
19 in developing further machine learning models, such as an automated system for
20 the real-time monitoring and annotation of smoke plumes directly from streaming
21 satellite data.

22 1 Introduction

23 In part, due to public policy, the average levels of fine particulate matter ($PM_{2.5}$) in the US have
24 generally been declining over the past few decades[1]. Despite those improvements, the contribution
25 of wildfire smoke to $PM_{2.5}$ concentrations in the US has been calculated to have more than doubled
26 between 2010 to 2020, accounting for up to half of the overall $PM_{2.5}$ exposure in Western regions
27 [7]. Increases in $PM_{2.5}$ due to wildfire smoke are concerning since ambient $PM_{2.5}$ exposure is a
28 leading environmental risk factor for adverse health effects and premature mortality[12]. These risks
29 underscore the necessity for efficient and effective monitoring methods to mitigate the adverse health
30 impacts associated with wildfire smoke.

31 Traditionally, wildfire monitoring has relied on ground-based methods, such as forest service patrols,
32 manned lookout towers, and aviation surveillance [2]. While these methods provide valuable localized
33 insights, they are constrained by geographical and logistical limitations, often failing to deliver timely
34 and comprehensive data, especially over large and remote areas. In contrast, satellite imagery offers

35 a vantage point that overcomes these limitations, providing continuous, wide-area coverage and
36 real-time data crucial for assessing and responding to the health risks posed by wildfire smoke.

37 Satellite imagery, equipped with state-of-the-art sensors, such as the Advanced Baseline Imager
38 (ABI) on the Geostationary Operational Environmental Satellites (GOES) [13], have revolutionized
39 environmental monitoring. These tools enable the detailed observation of smoke plumes, their
40 particulate density, and the extent of smoke spread. These satellite-based systems offer the capabilities
41 to provide critical insights into the concentration and movement of smoke particulates, facilitating
42 real-time assessments of air quality.

43 The integration of satellite imagery in wildfire smoke monitoring is not only instrumental in providing
44 real-time data but also plays a significant role in public health planning and response. By mapping
45 the spread and density of smoke, health authorities can issue timely warnings, implement evacuation
46 protocols, and deploy resources effectively to mitigate health risks. Furthermore, long-term data
47 gathered from satellite observations can aid in understanding the broader impacts of wildfire smoke
48 on public health, influencing policy decisions and preventive measures.

49 Currently, multi-channel thresholding is a popular method to distinguish smoke pixels from pixels
50 containing dust, clouds or other phenomenon with similar signatures [35]. Thresholds are determined
51 by using historical, labeled data to extract optimal radiance values for each channel that corresponds
52 with the labeled class. These methods are tuned to particular biogeographies and often have issues
53 with generalization to new locations with varying fuel types [25].

54 In contrast to the numerical thresholding approach, human visual inspection of satellite imagery is
55 another commonly used method for smoke identification. Trained analyst inspect satellite imagery
56 and label the smoke by hand. An example of hand-labeled annotations is the National Oceanic
57 and Atmospheric Administration (NOAA) Hazard Mapping System (HMS) fire and smoke product
58 [22, 28]. For the HMS smoke product, trained satellite analysts use movement characteristics to
59 help identify smoke by scanning through a time series of satellite imagery. When visual inspection
60 indicates smoke, the analyst will draw a polygon that corresponds to the geolocation and density
61 of smoke. By design of the product, the HMS annotations have varying time resolution and are
62 released on a rolling but undefined schedule ranging from one to multiple times a day as observation
63 conditions permit. This method is potentially not as scalable as an automated approach and is limited
64 by the availability of analysts and their time.

65 To address the challenges associated with thresholding and manual labels, we can look towards
66 innovative approaches and recent technological advancements in computer vision. Machine learning
67 methods have shown potential in improving the accuracy and efficiency of satellite-based wildfire
68 smoke detection and monitoring. For instance, SmokeNet, uses a convolutional neural network
69 (CNN) based framework to determine if a scene of MODIS satellite imagery contains smoke [4].
70 Another study, that looked at a singular wildfire event, also used a CNN to identify smoke on a
71 pixel-wise basis using imagery from Himawari-8 [18]. Additionally, Wen et al. developed a CNN
72 architecture that takes GOES-East imagery as input and the HMS-generated annotations for the target
73 labels during training [33].

74 The success of deep learning methods, such as CNNs, relies heavily on the availability of a large,
75 representative dataset [30]. As laid out in table 1, prior studies use relatively small numbers of
76 samples, from 47 [32] to 6825 [33], where one sample represents a satellite image with a singular
77 time and/or geolocation. In contrast, benchmark datasets for image classification contain tens of
78 thousands (CIFAR-10 and MNIST) to millions (CIFAR-100 and ImageNet) of data samples [17],
79 [10], [9]. Keeping in mind the correlation between both the quality and quantity of data with model
80 performance, we introduce the largest known smoke dataset, SmokeViz, containing over 130,000
81 samples.

82 Semi-supervised learning is an approach that can be used to increase the number of labeled samples
83 in a dataset. This is done by leveraging a labeled dataset to generate new labels for an often larger,
84 but unlabeled, dataset. Pseudo-labeling, a form of semi-supervised learning, uses labeled data to
85 train an initial model, then runs that model on unlabeled data to predict pseudo-labels, and finally
86 trains a new model using the pseudo-labels [19]. We introduce a variation of pseudo-labeling, not to
87 increase the size, but to increase the quality of our dataset by generating pseudo-labels to select the
88 best satellite image out of a given time-window to represent each smoke plume annotation.

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
[33]	CNN	6825	GOES-East	5	pixel
[18]	CNN	975	Himawari-8	7	pixel
[32]	U-Net	47	Landsat-8	13	pixel
SmokeViz	U-Net	133,871	GOES-East/West	3	pixel

89 2 Methods

90 Dataset

91 The initial source for smoke labels, discussed in further detail in the HMS Smoke Labels section, is
 92 uniquely characterized by each annotation, y , having corresponding imagery ranging between 1-60
 93 frames, where each frame, x , captures 5 minutes of exposure. Additionally, we have two satellites
 94 that overlap in coverage area, GOES-East and GOES-West, effectively doubling the number of frames
 95 for a single annotation. For the set of smoke annotations, \mathcal{Y} , $y \in \mathcal{Y}$ uses one or more $x \in \mathcal{X}$ where
 96 \mathcal{X} is the entire set of satellite imagery corresponding to the set of time windows defined by the
 97 labels. We apply pseudo-labeling to develop a subset of \mathcal{X} , denoted as \mathcal{X}_p , that has a one-to-one
 98 annotation-to-image ratio such that $|\mathcal{X}_p| = |\mathcal{Y}|$, where we choose the satellite image that has the
 99 maximum overlap between the geolocation of smoke in the imagery and the analyst annotation.

100 Dataset development came in three stages. First, we create an initial dataset, \mathcal{X}_M , that leverages light
 101 scattering physics to determine which singular satellite image would be in the optimal configuration
 102 for smoke detection. Second, we used \mathcal{X}_M to train an initial parent model, f_o , that identifies smoke in
 103 satellite imagery. Third, we use f_o to label each satellite image in a given annotation’s time-window
 104 and the optimal satellite image is chosen based on which image’s pseudo-labels has the greatest
 105 overlap with the analyst annotation for the given location and densities of smoke.

106 HMS Smoke Labels

107 NOAA manages environmental satellite programs such as the HMS program, the HMS program is an
 108 operational system that uses an aggregation of satellite data to generate active fire and smoke data.
 109 To train our model, we implement a supervised learning framework that uses the HMS analyst smoke
 110 product as truth labels during the model training process.

111 HMS smoke analysis data gives the coordinates of the smoke perimeter as a polygon and classifies
 112 the smoke by density within a given time window. The time windows can range from instantaneous
 113 (same start and end time) to lengths of 5 hours. While the true bounds of the smoke can change
 114 within the larger time spans, the analyst is making an approximation that should reflect the smoke
 115 coverage over the duration of the time window. The density information is qualitatively determined
 116 by each analyst based on the apparent smoke opacity in the satellite imagery and categorized as either
 117 light, medium or heavy as seen in figure 1a [23].

118 Thermometer Encoding Smoke Densities

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

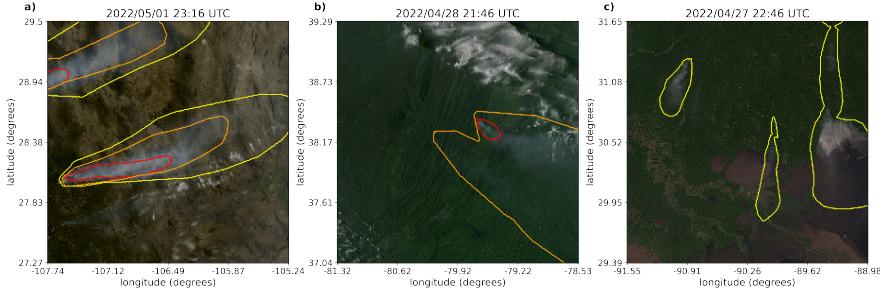


Figure 1: Satellite imagery captured by GOES-East within a few days of each other. The yellow, orange and red contours indicate the extent of Light, Medium and Heavy smoke. a) shows a canonical example of a smoke plume. b) and c) show observable variations in the density labels.

130 based on category. For example, the model can be given a small penalization for misclassifying light
 131 smoke as not smoke while given a higher penalization for misclassifying heavy smoke as not smoke.
 132 In addition to the densities being ordered and categorical, the differences between the density
 133 categories are not evenly distributed by a given metric, such as particulate matter per square meter.
 134 The intervals between densities being unknown along with the hierarchical nature of the density labels
 135 makes the labels ordinal instead of just categorical. This data property allows us to use thermometer
 136 encoding [6], which leverages the idea that heavy density smoke includes both medium and light
 137 density smoke, that heavy density smoke is closer to medium than it is to light, and automatically
 138 weights the loss functions and incorporates the ranked ordering of the densities. As seen in Table 2,
 139 one-hot encoding, commonly used for categorical data, doesn't take ordinal properties of the data
 140 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]

141 Time Windows For Smoke Annotations

142 In order to take into account movement characteristics to help identify smoke, analysts use multi-
 143 frame animations of the satellite imagery. The resulting annotations often have large time windows
 144 over multiple hours to represent one smoke plume annotation. Since the goal of these annotations is
 145 to show the general coverage over that time span, as shown in figure 2, the smoke boundaries don't
 146 often match up with the satellite imagery over the entire time window. One way to approach this
 147 problem would be to use all the satellite images the analysts used as input. Since the timespans are
 148 non-uniform, this would vary the length in imagery inputs into the model, which would be difficult
 149 with a CNN architecture. Moreover, this would require a large amount of additional memory and
 150 computational resources. Instead of using the original analysts' many satellite image inputs to one
 151 annotated output, we develop a one-to-one input-to-output by finding the optimal singular satellite
 152 image input to represent the annotation. Discussed in further detail in the next section, we do this
 153 by making physics-driven choices on which satellite and timestamp would give the optimal angle
 154 between the sun and satellite that would produce the strongest smoke signature for the geolocation
 155 and timestamp of the smoke plume.

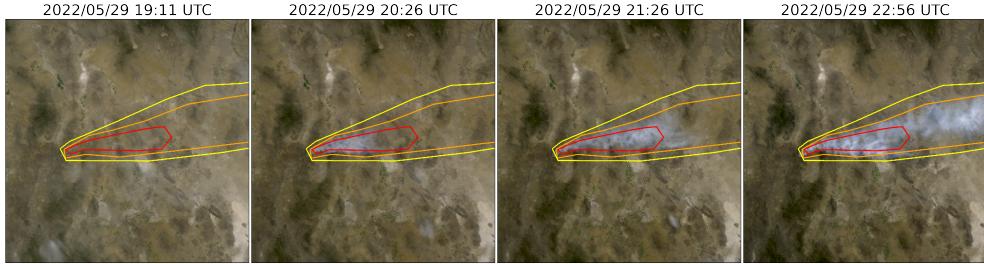


Figure 2: True color GOES-East imagery from May 2022, Southeast New Mexico (31°N , 100°W) during the start of the Foster Fire. The red, orange and yellow lines represent the heavy, medium and low density HMS smoke annotations that span 19:10–23:00 UTC.

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 (μm)	spatial resolution (km)
C01	blue visible	0.47	1
C02	red visible	0.64	0.5
C03	veggie near infrared	0.865	1

156 Satellite Imagery

157 The GOES satellites are operated by NOAA in order to support meteorology research and forecasting
 158 for the United States. We use the latest operational satellites, GOES-16 (East), 17 and 18 (West)
 159 that each carry the ABI, that measure 16 bands between the visible and infrared wavelengths. In
 160 improvement to the GOES predecessors, imagery is collected every 5 minutes for the contiguous
 161 United States and every 10 minutes for the full disk. Using PyTroll, a Python framework for
 162 processing satellite data [26], we input bands 1-3 (Table 3) to a GOES specific true color composite
 163 algorithm [5] to develop a, 1km resolution, true color image representation, similar to what is seen by
 164 HMS analysts. As discussed in further detail in the next section, the highest signal-to-noise ratio will
 165 come from the smallest wavelengths of light, higher wavelengths have lower smoke signal and higher
 166 noise (figure 5). For that reason, we only include the first 3 out of 16 available bands of data.

167 Mie-Derived Dataset

168 We used a physics-informed approach in selecting the initial GOES dataset, \mathcal{X}_M , which we call the
 169 Mie-derived dataset, for training an initial parent model, f_o , where if \mathcal{X} represents all the GOES
 170 imagery corresponding to the HMS smoke annotation time window, $\mathcal{X}_M \subset \mathcal{X}$. Prior GOES ABI
 171 datasets for machine learning applications often include data from only one of the two GOES-series
 172 satellites, commonly opting for GOES-East [33], [24], [20]. Rather than using one satellite or the
 173 cumulative data from both GOES-West and GOES-East images, we select between one or the other
 174 based on the solar zenith angle. For smoke identification, this approach can achieve a much higher
 175 signal-to-noise than imaging the earth’s surface from an arbitrary angle. The elastic scattering of
 176 light is the primary mechanism to account for - while the atmosphere is composed of molecules
 177 with size $< 1\text{nm}$, smoke particles can vary from $100\text{ nm} - 10\text{ }\mu\text{m}$ in diameter, d . The GOES ABI
 178 covers spectral bands from $0.47\text{ }\mu\text{m} - 13.3\text{ }\mu\text{m}$, so atmospheric and smoke particle sizes occupy two
 179 very different regimes with respect to the imaging wavelength λ . In the extreme limit of $\lambda \gg d$, the
 180 physics of scattering of light off a small sphere is captured by Rayleigh scattering. This process has
 181 two critical consequences: (1) the scattering cross section of light is strongly wavelength dependent
 182 (scaling with λ^{-4}), meaning that photons with wavelength closer to the ultraviolet are scattered more
 183 strongly than infrared photons. (2) the scattering cross section scales with an angular dependent
 184 cross section of $(1 + \cos^2 \theta)$. Scattered photons follow the emission distribution of a radiating dipole,

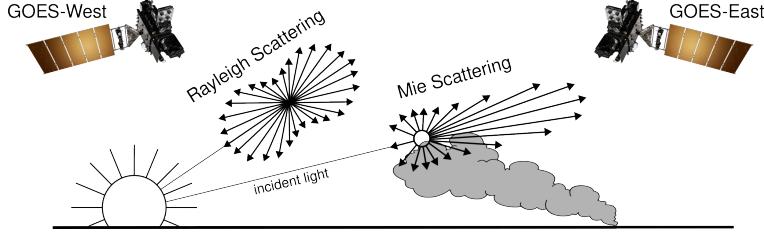


Figure 3: If the particle size is $< \frac{1}{10}$ the wavelength of the interacting light, then the primary scattering will be Rayleigh. Mie scattering is the predominant scattering mechanism when the particle size is larger than the wavelength of light. This schematic demonstrates that when the sun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

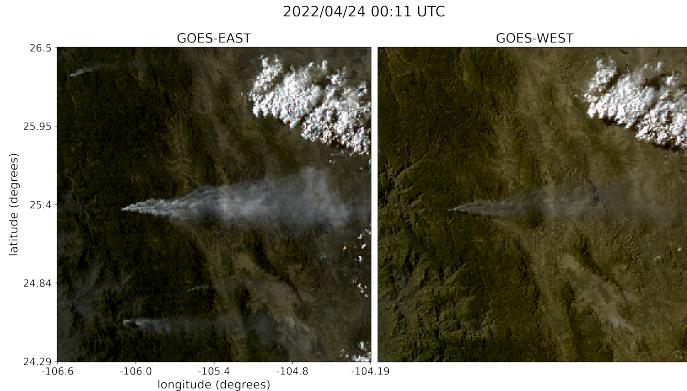


Figure 4: True color GOES-East (left) and GOES-West (right) imagery from April 24th, 2022 in Durango, Mexico. The images were taken ~ 0.5 hours before sunset (01:43 UTC) for this geolocation and time of year.

185 scattering more strongly in the forward and backwards directions ($\theta = 0, \pi$) than orthogonal to the
186 direction of propagation ($\theta = \pi/2, 3\pi/2$), see figure 3 for a Rayleigh scattering schematic.

187 The significance of these scalings is that the observer, or detector, will receive blue photons in most
188 directions orthogonal to the source. Equivalently, photons traveling colinearly with line of sight to
189 the emission source will mostly have wavelengths in the infrared band. In the converse regime of
190 $d > \lambda$, the elastic scattering of light against matter is modeled through Mie scattering. In comparison
191 to Rayleigh scattering, Mie scattering is largely wavelength-independent and has a more complicated
192 radiation pattern where the cross section has a maximal amplitude in the forward direction. An
193 observer downstream of this scatterer will collect more photons than one positioned directly behind it.
194 In the context of smoke identification, a sunrise or sunset will lead to a higher Mie scattered signal in
195 GOES-West and GOES-East respectively, as shown with a smoke plume producing a stronger signal
196 in GOES-East imagery near sunset in figure 4.

197 Smoke identification therefore amounts to extracting a signal of $d > \lambda$ photons from the $\lambda \gg d$
198 background. Positioning a detector along line of sight to the scatterer will result in a higher signal
199 from smoke particles (figure 3). Filtering the imaged wavelength can enhance this signal; photons
200 collected in the blue spectrum will have a naturally lower background along the line of sight to the
201 illumination source due to their high level of Rayleigh scattering as. Therefore, as demonstrated in figure
202 5, this configuration results in the highest signal to noise imaging for smoke particles.

203 Based solely on these criteria, the optimal strategy would be to pull data from GOES-West right after
204 sunrise and from GOES-East right before sunset. Another factor to consider is that the time when the
205 sun is in optimal alignment with the satellite for smoke detection coincides with when solar zenith
206 angle is maximized. Larger angles between the satellite and sun result in an increase in noise due
207 to increased atmospheric interactions [27]. This is shown in figure 6, while we optimize for smoke
208 signal detection, due to the high solar zenith angle, we introduce atmospheric interaction noise that

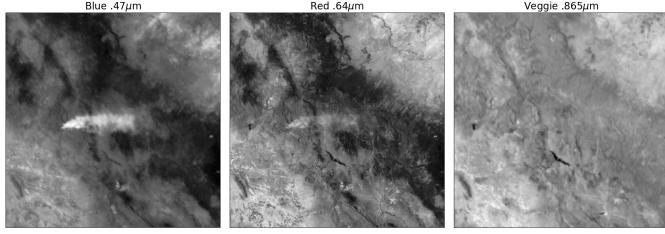


Figure 5: Three bands of GOES-East data are the raw input to generate a true color image. These plots show variations in the signal-to-noise ratio for smoke detection in relation to the wavelength, λ , of light being measured.

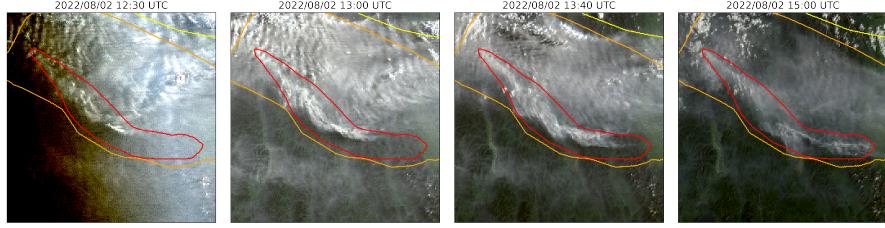


Figure 6: A smoke annotation projected onto GOES-West imagery from August 2022 that spans from 11:00 UTC to 15:00 UTC, sunrise on August 2nd, 2022 at coordinates ($49^{\circ}24'N$, $115^{\circ}29'W$) was 12:15 UTC.

209 obfuscate the smoke signal. To reduce the noise from large solar zenith angles, if given multiple
 210 frames to choose from, we choose the image with the largest solar zenith angle that is $< 80^{\circ}$.
 211 The resulting image selection process takes into account atmospheric properties and light scattering
 212 physics to generate an estimate of which singular satellite image within the analyst time-window could
 213 give the highest smoke signal-to-noise ratio. The resulting Mie-derived dataset, $\mathcal{X}_M = \{X_M, Y\}$,
 214 was then used to train a model, f_o , that would generate N pseudo-labels, y^* , for every sample,
 215 where N is determined by how many images, taken at a 10 minute interval, fit within the analyst
 216 time-window for that sample. Chosen from the N images, x_p is the image with the highest alignment
 217 between the f_o prediction of smoke, y^* , in the image and the HMS analysts' annotation y .

218 Machine Learning Model

219 We implement a deep learning architecture that uses the encoder from EfficientNetV2 [31] and a
 220 semantic segmentation classifier from the DeepLabV3 model [8]. Transfer learning has shown to
 221 reduce the time and resources needed to train a model by leveraging information from pre-trained
 222 models [34], [29]. We initialize the values of our model weights using the pre-trained values originally
 223 trained on the ImageNet dataset [9], containing 1.2 million images and 1000 categories. Our model
 224 was developed using the Segmentation Models PyTorch package [15] that was written as a high level
 225 API for implementing models for semantic segmentation problems. We input 256x256x3 snapshots of
 226 1km resolution true color GOES imagery that contains smoke and output a 256x256x3 classification
 227 map that predicts if a pixel contains smoke and if so, what the density of that smoke is. As mentioned
 228 earlier, we apply the thermometer encoding shown in table 2 to encode the smoke densities and apply
 229 binary cross entropy as the loss function per density of smoke.

230 The dataset, \mathcal{X}_M , contained over 130,000 samples. To train f_o , we split \mathcal{X}_M into training (118,691
 231 samples), validation (8,335 samples) and testing (7,474 samples) datasets. Training data contains

Table 4: 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
Light	0.394	0.551	0.437	0.583
Medium	0.283	0.392	0.345	0.431
Heavy	0.233	0.290	0.275	0.332
Overall	0.365	0.510	0.412	0.539

232 data from the years 2018, 2019, 2020, 2021 and 2023 while the data from 2022 is split into validation
 233 and testing sets by taking data from alternating 10 days of the year. In order to make sure we include
 234 the monthly variations in wildfire trends over a full year, we split 2022 data up by every 10 days.
 235 This allowed us to: (1) allocate an additional full year of data for the training set, (2) show yearlong
 236 trends in both the validation and testing sets and (3) keep the validation and testing datasets relatively
 237 independent from one another since only two out of every ten days of data will have adjacent days in
 238 validation and testing.

239 To determine which image out of the relevant imagery for the given time window best represents
 240 the analyst annotation, we implement a greedy algorithm by running f_o on each x to generate a
 241 pseudo-label, y^* . The output of f_o , y^* , give predictions on if smoke is in the image, and if there is
 242 smoke, where the smoke is in that image and the density of that smoke. y^* serve as pseudo-labels
 243 for each density of smoke and are compared to the analyst annotations, y . To compare y^* and y , we
 244 calculate the IoU using the total set of pixels for y^* at that density of smoke and the entire set of
 245 pixels for y for a particular smoke density in each image as shown in equation 1. The image with the
 246 highest IoU score is chosen as the image, x_p , that best represents the analyst smoke annotation, y .
 247 Often used for pseudo-labeling, a confidence threshold value is defined to determine if a pseudo-label
 248 should to be included in a dataset [11]. We chose a confidence threshold that would include the
 249 sample, x_p , in \mathcal{X}_p if the maximum overall IoU (equation 1) between y^* and y over all densities was
 250 over 0.01.

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

251 Finally, we use \mathcal{X}_p to train an additional child model, f_c . We use the same dataset split method and
 252 model setup but change \mathcal{X}_M to \mathcal{X}_p to train the child model. For training both f_c and f_p we train each
 253 model over 10 epochs using the Adam optimizer on a single Nvidia A100 GPU allocating 10GB of
 254 memory over 80 hours of allotted training time.

255 Results

256 To interpret the performance of f_o , we report the IoU metrics in table 4 that were computed by
 257 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
 258 pixels that f_o predicts as that density of smoke and the entire set of pixels labeled by the analyst
 259 as a particular smoke density over all imagery contained in the testing dataset. Additionally, we
 260 compute the overall IoU for all densities by first computing the number of pixels that intersect their
 261 corresponding density and divide that by the total number of pixels that make up the union of model
 262 predicted and analyst labeled smoke in the testing dataset.

263 An illustration of a pseudo-label picked image better representing the analyst annotation when
 264 compared to the Mie-derived image selection is evident in Figure 7, where the heavy density smoke
 265 IoU increases from 0.01 to 0.59. The analyst annotation for these densities cover 5 hours of imagery,
 266 the Mie-derived selection optimizes for the image closest to sunrise while the pseudo-label image
 267 selection chooses the image with the highest overlap between the pseudo-label and the analyst
 268 annotation.

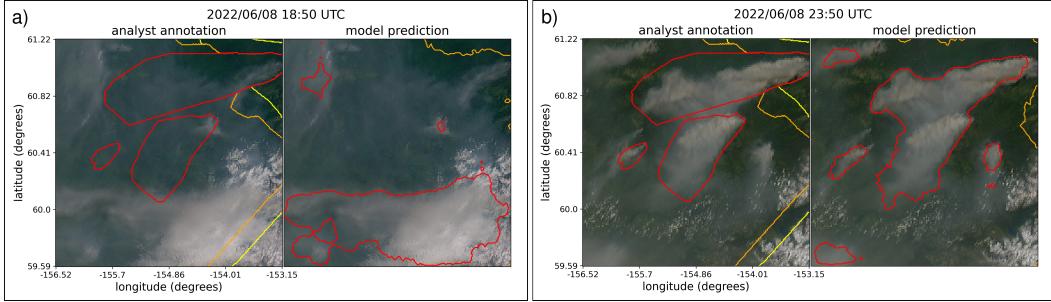


Figure 7: GOES-West imagery showing smoke on June 8th, 2022 in Alaska where, at this geolocation, daylight was between 12:43-7:53 UTC. The HMS smoke annotations displayed span from 18:50 to 23:50 UTC. a) shows the imagery that was selected using the Mie-derived data selection process b) shows the image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

269 3 Limitations

270 One of the concerns that comes with using pseudo-labeling methods is that you can perpetuate biases
 271 from the parent model into subsequent child models. Due to the increase in detectable forward
 272 scattered light off smoke particular matter, we expect the model to have a bias towards producing a
 273 higher success rate for smoke detection at larger solar zenith angles. Another concern is the possibility
 274 of a data leak between the adjacent days every 10 days for validation and testing set. Finally, the
 275 original HMS dataset is not split by type of fire and includes a large portion of small, controlled burns.
 276 This can be a limitation to consider if the dataset is being used to detect large wildfires. All these
 277 limitations are discussed and analyzed further in the Appendix.

278 4 Conclusion

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

295 To improve the dataset’s capacity to accurately represent wildfire smoke plumes, we train a parent
 296 machine learning model, f_o , using the Mie-derived dataset, \mathcal{X}_M , and run it on the relevant satellite
 297 images within the time-frame. The image with the maximum IoU score between the model’s smoke
 298 predictions, or pseudo-label, and the analyst smoke annotations are used to create the pseudo-label
 299 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
 300 testing sets from \mathcal{X}_M and \mathcal{X}_p . The results reported in table 4 suggest that \mathcal{X}_p was able to train a better
 301 performing model, f_c , that gave higher IoU metrics on both dataset’s testing sets in comparison to
 the original parent model, f_o .

303 The result of this study is a representative dataset, SmokeViz¹, that can be used to train machine
304 learning models for various wildfire smoke applications. A future goal is to produce a robust
305 and reliable machine learning based approach for detecting wildfires using satellite imagery. That
306 information can be used for wildfire monitoring and as data provided to public health officials for air
307 quality assessments. On a broader scale, we show how pseudo-labeling can be used to optimize a
308 dataset when the resolution for the data and corresponding labels do not match. This could be useful
309 in similar applications involving time-series/video data with a singular label where the data can be
310 compressed while still remaining representative of the label.

311 5 Acknowledgments and Disclosure of Funding

312 This research was supported in part by NOAA cooperative agreement NA22OAR4320151, for the
313 Cooperative Institute for Earth System Research and Data Science (CIESRDS). We thank Wilfrid
314 Schroeder and the Hazard Mapping Systems team for giving guidance on how they created their
315 smoke plume dataset. This work utilized the Alpine high performance computing resource at the
316 University of Colorado Boulder. Alpine is jointly funded by the University of Colorado Boulder, the
317 University of Colorado Anschutz, Colorado State University, and the National Science Foundation
318 (award 2201538). The statements, findings, conclusions, and recommendations are those of the
319 author(s) and do not necessarily reflect the views of NOAA or the U.S. Department of Commerce.

320 References

- 321 [1] J. E. Aldy, M. Auffhammer, M. Cropper, A. Fraas, and R. Morgenstern. Looking back at 50
322 years of the clean air act. *Journal of Economic Literature*, 60(1):179–232, 2022.
- 323 [2] R. S. Allison, J. M. Johnston, G. Craig, and S. Jennings. Airborne optical and thermal remote
324 sensing for wildfire detection and monitoring. *Sensors*, 16(8):1310, 2016.
- 325 [3] N. Andela, J. Kaiser, G. Van der Werf, and M. Wooster. New fire diurnal cycle characterizations
326 to improve fire radiative energy assessments made from modis observations. *Atmospheric
327 Chemistry and Physics*, 15(15):8831–8846, 2015.
- 328 [4] R. Ba, C. Chen, J. Yuan, W. Song, and S. Lo. Smokenet: Satellite smoke scene detection using
329 convolutional neural network with spatial and channel-wise attention. *Remote Sensing*, 11(14):
330 1702, 2019.
- 331 [5] M. Bah, M. Gunshor, and T. Schmit. Generation of goes-16 true color imagery without a green
332 band. *Earth and Space Science*, 5(9):549–558, 2018.
- 333 [6] J. Buckman, A. Roy, C. Raffel, and I. Goodfellow. Thermometer encoding: One hot way to
334 resist adversarial examples. In *International conference on learning representations*, 2018.
- 335 [7] M. Burke, A. Driscoll, S. Heft-Neal, J. Xue, J. Burney, and M. Wara. The changing risk and
336 burden of wildfire in the united states. *Proceedings of the National Academy of Sciences*, 118
337 (2):e2011048118, 2021.
- 338 [8] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic
339 image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs.
340 *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.
- 341 [9] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical
342 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages
343 248–255. Ieee, 2009.
- 344 [10] L. Deng. The mnist database of handwritten digit images for machine learning research [best of
345 the web]. *IEEE signal processing magazine*, 29(6):141–142, 2012.
- 346 [11] R. E. Ferreira, Y. J. Lee, and J. R. Dórea. Using pseudo-labeling to improve performance of
347 deep neural networks for animal identification. *Scientific Reports*, 13(1):13875, 2023.

¹<https://noaa-gsl-experimental-pds.s3.amazonaws.com/index.html#SmokeViz/>

- 348 [12] E. Gakidou, A. Afshin, A. A. Abajobir, K. H. Abate, C. Abbafati, K. M. Abbas, F. Abd-Allah,
 349 A. M. Abdulle, S. F. Abera, V. Aboyans, et al. Global, regional, and national comparative risk
 350 assessment of 84 behavioural, environmental and occupational, and metabolic risks or clusters
 351 of risks, 1990–2016: a systematic analysis for the global burden of disease study 2016. *The
 352 Lancet*, 390(10100):1345–1422, 2017.
- 353 [13] S. J. Goodman, T. J. Schmit, J. Daniels, and R. J. Redmon. *The GOES-R series: a new
 354 generation of geostationary environmental satellites*. Elsevier, 2019.
- 355 [14] E. J. Hyer, J. S. Reid, E. M. Prins, J. P. Hoffman, C. C. Schmidt, J. I. Miettinen, and L. Giglio.
 356 Patterns of fire activity over indonesia and malaysia from polar and geostationary satellite
 357 observations. *Atmospheric research*, 122:504–519, 2013.
- 358 [15] P. Iakubovskii. Segmentation models pytorch. https://github.com/qubvel/segmentation_models.pytorch, 2019.
- 360 [16] J. E. Keeley and A. D. Syphard. Large california wildfires: 2020 fires in historical context. *Fire
 361 Ecology*, 17:1–11, 2021.
- 362 [17] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- 363 [18] A. Larsen, I. Hanigan, B. J. Reich, Y. Qin, M. Cope, G. Morgan, and A. G. Rappold. A deep
 364 learning approach to identify smoke plumes in satellite imagery in near-real time for health risk
 365 communication. *Journal of exposure science & environmental epidemiology*, 31(1):170–176,
 366 2021.
- 367 [19] D.-H. Lee. Pseudo-label : The simple and efficient semi-supervised learning method for deep
 368 neural networks. *ICML 2013 Workshop : Challenges in Representation Learning (WREPL)*, 07
 369 2013.
- 370 [20] Y. Lee, C. D. Kummerow, and I. Ebert-Uphoff. Applying machine learning methods to detect
 371 convection using geostationary operational environmental satellite-16 (goes-16) advanced
 372 baseline imager (abi) data. *Atmospheric Measurement Techniques*, 14(4):2699–2716, 2021.
- 373 [21] J. L. McCarty, C. O. Justice, and S. Korontzi. Agricultural burning in the southeastern united
 374 states detected by modis. *Remote Sensing of Environment*, 108(2):151–162, 2007.
- 375 [22] D. McNamara, G. Stephens, M. Ruminski, and T. Kasheta. The hazard mapping system (hms) -
 376 noaa's multi-sensor fire and smoke detection program using environmental satellites. *Conference
 377 on Satellite Meteorology and Oceanography*, 01 2004.
- 378 [23] NOAA. Hazard mapping system fire and smoke product. URL <https://www.ospo.noaa.gov/Products/land/hms.html#about>.
- 380 [24] T. C. Phan and T. T. Nguyen. Remote sensing meets deep learning: exploiting spatio-temporal-
 381 spectral satellite images for early wildfire detection. 2019.
- 382 [25] T. Randriambelo, S. Baldy, M. Bessafi, M. Petit, and M. Despinoy. An improved detection
 383 and characterization of active fires and smoke plumes in south-eastern africa and madagascar.
 384 *International Journal of Remote Sensing*, 19(14):2623–2638, 1998.
- 385 [26] M. Raspaud, D. Hoese, A. Dybbroe, P. Lahtinen, A. Devasthale, M. Itkin, U. Hamann, L. Ø.
 386 Rasmussen, E. S. Nielsen, T. Leppelt, et al. Pytroll: An open-source, community-driven python
 387 framework to process earth observation satellite data. *Bulletin of the American Meteorological
 388 Society*, 99(7):1329–1336, 2018.
- 389 [27] A. Royer, P. Vincent, and F. Bonn. Evaluation and correction of viewing angle effects on
 390 satellite measurements of bidirectional reflectance. *Photogrammetric engineering and remote
 391 sensing*, 51(12):1899–1914, 1985.
- 392 [28] W. Schroeder, M. Ruminski, I. Csizar, L. Giglio, E. Prins, C. Schmidt, and J. Morisette.
 393 Validation analyses of an operational fire monitoring product: The hazard mapping system.
 394 *International Journal of Remote Sensing*, 29(20):6059–6066, 2008.

- 395 [29] A. Sharif Razavian, H. Azizpour, J. Sullivan, and S. Carlsson. Cnn features off-the-shelf: an
 396 astounding baseline for recognition. In *Proceedings of the IEEE conference on computer vision*
 397 and pattern recognition workshops, pages 806–813, 2014.
- 398 [30] C. Sun, A. Shrivastava, S. Singh, and A. Gupta. Revisiting unreasonable effectiveness of data in
 399 deep learning era. In *Proceedings of the IEEE international conference on computer vision*,
 400 pages 843–852, 2017.
- 401 [31] M. Tan and Q. Le. Efficientnetv2: Smaller models and faster training. In *International*
 402 *conference on machine learning*, pages 10096–10106. PMLR, 2021.
- 403 [32] Z. Wang, P. Yang, H. Liang, C. Zheng, J. Yin, Y. Tian, and W. Cui. Semantic segmentation and
 404 analysis on sensitive parameters of forest fire smoke using smoke-unet and landsat-8 imagery.
 405 *Remote Sensing*, 14(1):45, 2022.
- 406 [33] J. Wen and M. Burke. Wildfire smoke plume segmentation using geostationary satellite imagery.
 407 *ArXiv*, abs/2109.01637, 2021. URL <https://api.semanticscholar.org/CorpusID:237416777>.
- 409 [34] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural
 410 networks? *Advances in neural information processing systems*, 27, 2014.
- 411 [35] T. X.-P. Zhao, S. Ackerman, and W. Guo. Dust and smoke detection for multi-channel imagers.
 412 *Remote Sensing*, 2(10):2347–2368, 2010. ISSN 2072-4292. doi: 10.3390/rs2102347. URL
 413 <https://www.mdpi.com/2072-4292/2/10/2347>.

414 A Appendix

415 A.1 Original Data and Software Licenses

416 The HMS Smoke product does not have a license attached to it. For GOES imagery, NOAA states
 417 "There are no restrictions on the use of this data" and does not provide a license. Pytroll is distributed
 418 under the GNU General Public License v3.0 license while Segmentation Models Pytorch is distributed
 419 under the MIT License.

420 A.2 Statistical Visualizations for SmokeViz Dataset

421 Figures 8, 9, 10, 11 provide some statistical analysis on \mathcal{X}_p . As seen in figure 8, we see the highest
 422 number of samples for the year 2020 that showed a high volume of available annotations that year
 423 likely due to the large number of wildfires [16] during 2020. The peak for number of samples shown
 424 in figure 9 is March and April, coming right before the typical wildfire season that usually goes from
 425 late Spring through Fall. This may be due to the increase in prescribed agricultural burns before
 426 plants emerge from winter dormancy [21]. The HMS analysts do not have a way of distinguishing
 427 between planned or uncontrolled fire, so many of the annotations represent small agricultural burns
 428 along with wildfires.

429 As shown in figure 10, the states with the highest number of samples are California, Georgia and
 430 Florida. The high frequency in fires in the Southeast may be due to the aforementioned prescribed
 431 agricultural burns. Analysts are looking not only at the United States, but also Canada and Mexico,
 432 figure 11 shows a breakdown of the number of samples that originate from each country.

433 A.3 Model Performance Analysis

434 In order to get a better understanding of the dataset, we use the deep learning models to analyze
 435 certain data characteristics. Figure 12 shows variations in overall IoU values running f_o on the \mathcal{X}_p
 436 test set data. The highest IoU are during the typical wildfire season and outside the typical window
 437 for prescribed agricultural burns.

438 We report on how many samples come from each satellite in table 5, along with the \mathcal{X}_p test set
 439 IoU in comparison to the HMS analyst annotations. While GOES-EAST provides over triple the
 440 number of training samples, f_c performs better on GOES-WEST samples out of the test set. The

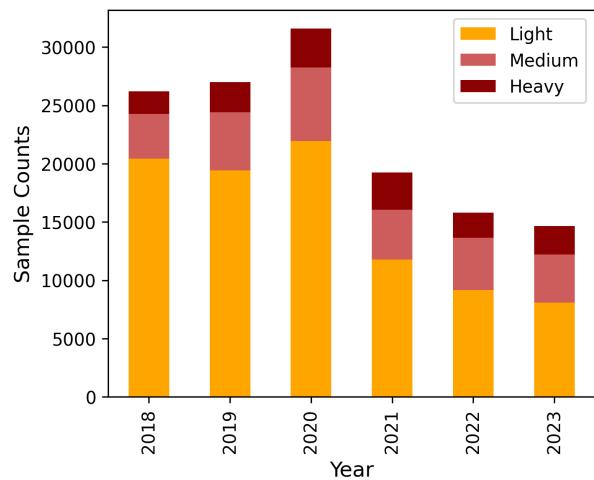


Figure 8: Sample count per year

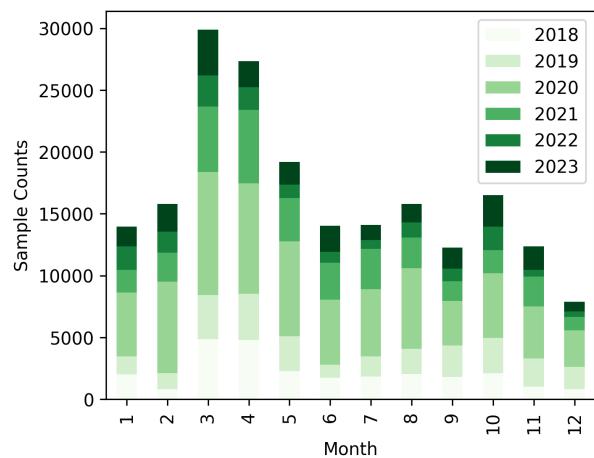


Figure 9: Sample count per month.

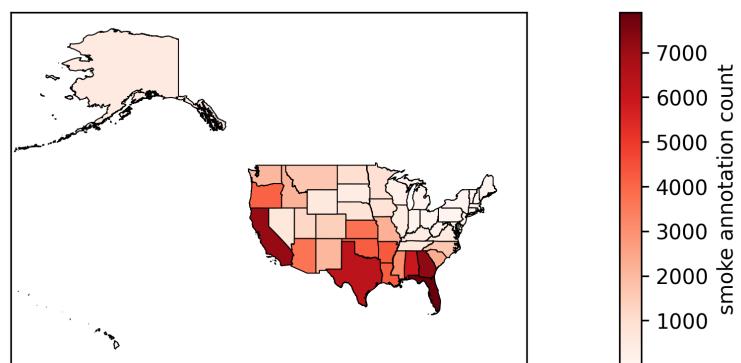


Figure 10: Sample count per US state.

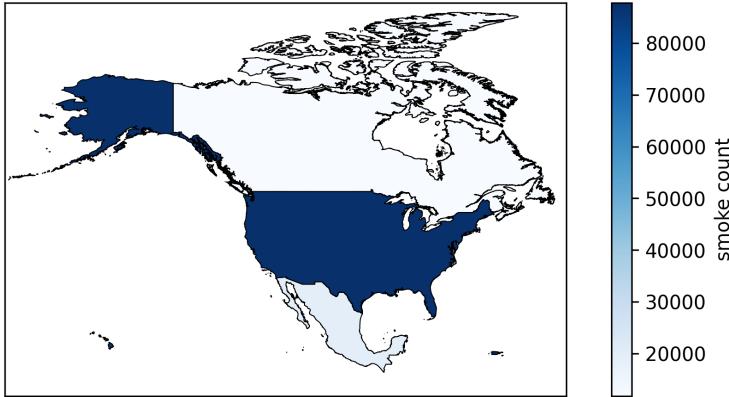


Figure 11: Sample count per North American country.

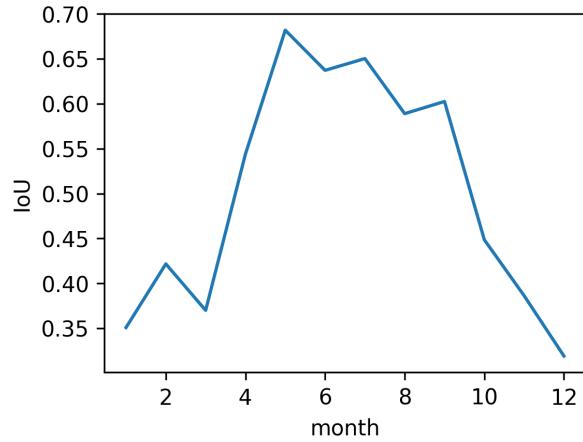


Figure 12: IoU between f_c predictions and analyst annotations per month for \mathcal{X}_p test set.

441 signal observed by a single satellite vary diurnally and annually in the amount of atmospheric noise
 442 and solar radiation. In turn, if provided with enough samples, this could create a more robust and
 443 generalizable model to the extent of being able to perform well on two different sensors with varying
 444 calibrations and line of sights.

445 As mentioned in the limitations, there may have been a bias introduced towards correctly classifying
 446 imagery close to sunrise or sunset. This bias may not only be introduced by our Mie-derived dataset
 447 that was used to train f_o , but also in the original HMS annotations. The configuration of the sun,
 448 smoke and satellite give the highest signal-to-noise ratio at the times near the sunrise and sunset,
 449 making smoke more easily observable. In contrast, the diurnal variations of wildfires cause the
 450 fire radiative power to be highest around solar noon [3]. Table 6 shows how the IoU between f_c
 451 predictions and analyst annotations for the test data from either \mathcal{X}_M or \mathcal{X}_p are not significantly
 452 affected by being within 2 hours to sunrise/sunset. The main difference we see from table 6 is the
 453 split of closer to daylight boundaries is shifted towards midday between \mathcal{X}_M to \mathcal{X}_p . This is because,
 454 for \mathcal{X}_p , we are choosing the imagery with the best overlap to the analyst product rather than the image
 455 from \mathcal{X}_M that optimized for highest possible signal-to-noise ratio if given constant signal.

456 In order to observe geographical regional variations we create quadrants, Northwest (NW), Southwest
 457 (SW), Northeast (NE) and Southeast (SE) in relation to the midpoint (40, -100) and show the
 458 sample distribution and model performance for each region in table 7. The table shows the worst f_c
 459 performance in the SE quadrant despite representing this largest fraction of the training data. This
 460 is likely due to the large number of aforementioned prescribed burns in that area. If the goal of
 461 the dataset is to be used to train a model to detect and monitor large wildfires, a weakness in the

Table 5: Sample count along with variations in f_c performance depending on which GOES satellite data is used.

Satellite	Test IoU	\mathcal{X}_p Test Samples	\mathcal{X}_p Samples
GOES-WEST	0.645	1827	30640
GOES-EAST	0.483	5647	119040

Table 6: Variations in f_c performance depending on temporal proximity to sunrise or sunset.

Time difference	\mathcal{X}_M Test Set IoU	\mathcal{X}_p Test Set IoU	\mathcal{X}_M Test Samples	\mathcal{X}_p Test Samples
<2 hours	0.412	0.546	3923 (63%)	3436 (46%)
>2 hours	0.411	0.538	2280 (37%)	4038 (54%)

462 dataset would be that it likely consists of a lot more small, controlled agricultural burns that aren't
463 representative of the intended task.

464 A weakness in the dataset split for 2022 validation and testing sets is that there are adjacent days
465 between the rotating 10 day splits. This is a weakness because wildfires often last more than one day,
466 smoke from the same fires are likely to leak between the datasets. The choice to split the dataset
467 every 10 days was a trade off between being able to keep another day for training and keeping the
468 validation and test set completely independent. Another consideration for the choice was that we
469 expect the diurnal variations in smoke characteristics to vary largely enough at either ends of the
470 nocturnal stagnations in fire activity [14]. The scope of this paper was to use the deep learning models
471 as a way of optimizing the dataset and comparing the datasets against each other. While the data leak
472 is not likely to have high consequences for this particular application (as suggested in table 8), we
473 encourage users of SmokeViz to split validation and test sets so that they are completely independent,
474 especially as new years of data are added.

475 A.4 Machine Learning Reproducibility

476 All relevant code is accessible at <https://github.com/reykoki/SmokeViz>. The models pre-
477 sented in this paper are not optimized for performance, but are intended to create sufficient pseudo-
478 labels to develop the SmokeViz dataset and then compare the performance of SmokeViz against the
479 original dataset. We did not perform any experimentation for deciding on architecture or hyperpa-
480 rameters shown in table 9, but did make educated decisions. We chose DeepLabV3+ because smoke
481 varies in scale and the DeepLabV3+ backbone uses a atrous spatial pyramid pooling module that
482 allows for varying scales of the same type of object. We use the Adam optimizer that will adapt the
483 learning rate during training and is suited for problems with large amounts of data. Batch size was
484 chosen due to the necessity to run the model on limited resources.

485 A.5 Datasheet for SmokeViz

486 Questions from the <https://arxiv.org/abs/1803.09010> paper, v7.

487 A.5.1 Motivation

488 The questions in this section are primarily intended to encourage dataset creators to clearly articulate
489 their reasons for creating the dataset and to promote transparency about funding interests.

Table 7: Along with sample count we show variations in f_c performance depending on quadrant.

Quadrant	\mathcal{X}_p Test IoU	\mathcal{X}_p Test Samples	\mathcal{X}_p Samples
NW	0.5932	1425	23335
SW	0.6094	1131	26577
NE	0.4726	252	8392
SE	0.4706	4666	76130

Table 8: Comparison of the IoU and loss between the full \mathcal{X}_p test set and the \mathcal{X}_p test set with adjacent days between the validation and test set removed.

\mathcal{X}_p Test Set	Overall IoU	Testing Loss
full test set	0.539	0.870
adjacent days removed	0.547	0.895

Table 9: Hyperparameters used to create f_o and f_c .

parameter	value
epochs	10
learning rate	1e-2
batch size	32
optimizer	Adam

490 **For what purpose was the dataset created?**

491 SmokeViz was created to serve as a large labeled dataset to be used in creating wildfire smoke plume
492 related machine learning models. Applications include wildfire smoke detection or smoke dispersion
493 modeling.

494 **Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g.,
495 company, institution, organization)?**

496 SmokeViz was created a group of researchers that at the time of the dataset creation were affiliated
497 with The National Oceanic and Atmospheric Administration and The University of Colorado, Boulder,
498 and The Cooperative Institute for Research in Environmental Sciences that connects CU, Boulder to
499 NOAA.

500 **Who funded the creation of the dataset?**

501 This work was funded by the National Oceanic and Atmospheric Administration and The Cooperative
502 Institute for Research in Environmental Sciences.

503 **Any other comments?**

504 None.

505 **A.5.2 Composition**

506 Most of these questions are intended to provide dataset consumers with the information they need to
507 make informed decisions about using the dataset for specific tasks. The answers to some of these
508 questions reveal information about compliance with the EU’s General Data Protection Regulation
509 (GDPR) or comparable regulations in other jurisdictions.

510 **What do the instances that comprise the dataset represent (e.g., documents, photos, people,
511 countries)?**

512 Each instance is a 256x256x3 RGB image from GOES imagery with an accompanying 256x256x3
513 binary masks corresponding to density of smoke. There are 3 densities of smoke - Light, Medium
514 and Heavy.

515 **How many instances are there in total (of each type, if appropriate)?**

516 There are 134500 samples, 90810 for light, 28023 for medium and 15667 for Heavy density smoke.

517 **Does the dataset contain all possible instances or is it a sample (not necessarily random) of
518 instances from a larger set?**

519 It is intended to contain all smoke data from 2018 through 2023 but we cut out imagery if it is too
520 bright or too dim based on photon count.

521 **What data does each instance consist of?**

522 The data is processed to correct for Rayleigh scattering, solar zenith angle and projected so each pixel
523 is representative of the same area of land. The algorithm is referenced in the SmokeViz paper.

524 **Is there a label or target associated with each instance?**

525 Yes, there are no instances that do not contain smoke.

526 **Is any information missing from individual instances?**

527 We have seen imagery where smoke is labeled but there's adjacent smoke plumes that were unlabeled.
528 With human labels comes human errors.

529 **Are relationships between individual instances made explicit (e.g., users' movie ratings, social
530 network links)?**

531 Some instances can overlap in geographic location, there can be multiple smoke plumes in one
532 instance, but the index of the HMS smoke annotation is listed and can be mapped back to the original
533 dataset for geolocation information.

534 **Are there recommended data splits (e.g., training, development/validation, testing)?**

535 We recommend using full years of data for training, validation and testing, but split testing and
536 validation every 10 days for 2022 in order to keep more data in the training set.

537 **Are there any errors, sources of noise, or redundancies in the dataset?**

538 The HMS smoke annotations that are used as truth are a source of noise as explained in the SmokeViz
539 paper. These include approximations of smoke polygons mismatching actual location and time
540 windows being too large that smoke moves during the time window. There is also noise caused by
541 atmospheric interactions with light. Redundancies occur when there are more than one smoke plume and
542 annotation in one image.

543 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,
544 websites, tweets, other datasets)?**

545 The dataset is self-contained.

546 **Does the dataset contain data that might be considered confidential (e.g., data that is pro-
547 tected by legal privilege or by doctor-patient confidentiality, data that includes the content of
548 individuals' non-public communications)?**

549 No.

550 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening,
551 or might otherwise cause anxiety?**

552 No.

553 **Does the dataset relate to people?**

554 No, not directly, wildfires do affect people, but these images are at 1km resolution.

555 **Does the dataset identify any subpopulations (e.g., by age, gender)?**

556 No.

557 **Is it possible to identify individuals (i.e., one or more natural persons), either directly or
558 indirectly (i.e., in combination with other data) from the dataset?**

559 No.

560 **Does the dataset contain data that might be considered sensitive in any way (e.g., data that
561 reveals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or
562 union memberships, or locations; financial or health data; biometric or genetic data; forms of
563 government identification, such as social security numbers; criminal history)?**

564 No.

565 **Any other comments?**

566 No.

567 **A.5.3 Collection process**

568 The answers to questions here may provide information that allow others to reconstruct the dataset
569 without access to it.

570 **How was the data associated with each instance acquired?**

571 The labeled from HMS smoke product is not validated or verified but is used as verification for
572 numerical smoke dispersion modeling. The GOES imagery is collected by the ABI sensor and is
573 corrected for any anomalies and also converted from photon count to radiance values.

574 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or
575 sensor, manual human curation, software program, software API)?**

576 Original low temporal resolution annotations were manual human analyst curated. To create the high
577 temporal resolution annotations, we use pseudo-labeling discussed in the SmokeViz paper.

578 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,
579 probabilistic with specific sampling probabilities)?**

580 The HMS smoke analysts are only looking for smoke during the daytime.

581 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and
582 how were they compensated (e.g., how much were crowdworkers paid)?**

583 The NOAA employed analysts are compensated as salaried federal employees.

584 **Over what timeframe was the data collected?**

585 2018-2023

586 **Were any ethical review processes conducted (e.g., by an institutional review board)?**

587 No.

588 **A.5.4 Preprocessing/cleaning/labeling**

589 The questions in this section are intended to provide dataset consumers with the information they
590 need to determine whether the “raw” data has been processed in ways that are compatible with their
591 chosen tasks. For example, text that has been converted into a “bag-of-words” is not suitable for tasks
592 involving word order.

593 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,
594 tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing
595 of missing values)?**

596 The data was processed according to the GOES True Color paper referenced in the SmokeViz methods
597 section.

598 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support
599 unanticipated future uses)?**

600 The raw data is available from the NOAA AWS webpage. <https://registry.opendata.aws/noaa-goes/>
601 The HMS smoke annotations are available here: <https://www.ospo.noaa.gov/products/land/hms.html>

602 **Is the software used to preprocess/clean/label the instances available?**

603 Yes, Pytroll implements the algorithm discussed in the GOES True Color paper referenced in the
604 SmokeViz paper.

605 **Any other comments?** None.

606 **A.5.5 Uses**

607 These questions are intended to encourage dataset creators to reflect on the tasks for which the dataset
608 should and should not be used. By explicitly highlighting these tasks, dataset creators can help dataset
609 consumers to make informed decisions, thereby avoiding potential risks or harms.

610 **Has the dataset been used for any tasks already?**

- 611 Not yet.
- 612 **Is there a repository that links to any or all papers or systems that use the dataset?**
- 613 No.
- 614 **What (other) tasks could the dataset be used for?** Smoke dispersion modeling, automated wildfire smoke detection.
- 616 **Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?**
- 618 No.
- 619 **Are there tasks for which the dataset should not be used?**
- 620 No. Any other comments? None
- 621 **A.5.6 Distribution**
- 622 **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?**
- 624 No.
- 625 **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?**
- 626 Amazon Web Services hosted by NOAA.
- 627 **When will the dataset be distributed?**
- 628 It is currently available.
- 629 **Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)?**
- 631 No. Have any third parties imposed IP-based or other restrictions on the data associated with the instances?
- 632
- 633 No.
- 634 **Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?**
- 636 No.
- 637 **Any other comments?**
- 638 None.
- 639 Maintenance
- 640 These questions are intended to encourage dataset creators to plan for dataset maintenance and communicate this plan with dataset consumers.
- 642 **Who is supporting/hosting/maintaining the dataset?**
- 643 NOAA.
- 644 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**
- 645 rey.koki@noaa.gov
- 646 **Is there an erratum?**
- 647 No.
- 648 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?**
- 649 yes
- 650 **If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were individuals in question told that their data would be retained for a fixed period of time and then deleted)?**

- 653 Not applicable.
- 654 **Will older versions of the dataset continue to be supported/hosted/maintained?**
- 655 No, if it needs to be updated, it is too large to keep multiple versions.
- 656 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?**
- 658 We encourage anyone that would like to contribute to SmokeViz to reach out to Rey Koki at
659 rey.koki@noaa.gov
- 660 **Any other comments?**
- 661 None

662 **NeurIPS Paper Checklist**

663 **1. Claims**

664 Question: Do the main claims made in the abstract and introduction accurately reflect the
665 paper's contributions and scope?

666 Answer: [Yes]

667 Justification: The claims of using pseudolabels to create a more robust dataset is reflected in
668 the paper's contributions.

669 Guidelines:

- 670 • The answer NA means that the abstract and introduction do not include the claims
671 made in the paper.
- 672 • The abstract and/or introduction should clearly state the claims made, including the
673 contributions made in the paper and important assumptions and limitations. A No or
674 NA answer to this question will not be perceived well by the reviewers.
- 675 • The claims made should match theoretical and experimental results, and reflect how
676 much the results can be expected to generalize to other settings.
- 677 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
678 are not attained by the paper.

679 **2. Limitations**

680 Question: Does the paper discuss the limitations of the work performed by the authors?

681 Answer: [Yes]

682 Justification: We address limitations of the dataset.

683 Guidelines:

- 684 • The answer NA means that the paper has no limitation while the answer No means that
685 the paper has limitations, but those are not discussed in the paper.
- 686 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 687 • The paper should point out any strong assumptions and how robust the results are to
688 violations of these assumptions (e.g., independence assumptions, noiseless settings,
689 model well-specification, asymptotic approximations only holding locally). The authors
690 should reflect on how these assumptions might be violated in practice and what the
691 implications would be.
- 692 • The authors should reflect on the scope of the claims made, e.g., if the approach was
693 only tested on a few datasets or with a few runs. In general, empirical results often
694 depend on implicit assumptions, which should be articulated.
- 695 • The authors should reflect on the factors that influence the performance of the approach.
696 For example, a facial recognition algorithm may perform poorly when image resolution
697 is low or images are taken in low lighting. Or a speech-to-text system might not be
698 used reliably to provide closed captions for online lectures because it fails to handle
699 technical jargon.
- 700 • The authors should discuss the computational efficiency of the proposed algorithms
701 and how they scale with dataset size.
- 702 • If applicable, the authors should discuss possible limitations of their approach to
703 address problems of privacy and fairness.
- 704 • While the authors might fear that complete honesty about limitations might be used by
705 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
706 limitations that aren't acknowledged in the paper. The authors should use their best
707 judgment and recognize that individual actions in favor of transparency play an impor-
708 tant role in developing norms that preserve the integrity of the community. Reviewers
709 will be specifically instructed to not penalize honesty concerning limitations.

710 **3. Theory Assumptions and Proofs**

711 Question: For each theoretical result, does the paper provide the full set of assumptions and
712 a complete (and correct) proof?

713 Answer: [NA]

714 Justification: No theoretical results are presented.

715 Guidelines:

- 716 • The answer NA means that the paper does not include theoretical results.
- 717 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 718 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 719 • The proofs can either appear in the main paper or the supplemental material, but if
- 720 they appear in the supplemental material, the authors are encouraged to provide a short
- 721 proof sketch to provide intuition.
- 722 • Inversely, any informal proof provided in the core of the paper should be complemented
- 723 by formal proofs provided in appendix or supplemental material.
- 724 • Theorems and Lemmas that the proof relies upon should be properly referenced.

725 **4. Experimental Result Reproducibility**

726 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
727 perimental results of the paper to the extent that it affects the main claims and/or conclusions
728 of the paper (regardless of whether the code and data are provided or not)?

729 Answer: [Yes]

730 Justification: We provide the code to create the datasets along with the final dataset hosted
731 on AWS by NOAA.

732 Guidelines:

- 733 • The answer NA means that the paper does not include experiments.
- 734 • If the paper includes experiments, a No answer to this question will not be perceived
735 well by the reviewers: Making the paper reproducible is important, regardless of
736 whether the code and data are provided or not.
- 737 • If the contribution is a dataset and/or model, the authors should describe the steps taken
738 to make their results reproducible or verifiable.
- 739 • Depending on the contribution, reproducibility can be accomplished in various ways.
740 For example, if the contribution is a novel architecture, describing the architecture fully
741 might suffice, or if the contribution is a specific model and empirical evaluation, it may
742 be necessary to either make it possible for others to replicate the model with the same
743 dataset, or provide access to the model. In general, releasing code and data is often
744 one good way to accomplish this, but reproducibility can also be provided via detailed
745 instructions for how to replicate the results, access to a hosted model (e.g., in the case
746 of a large language model), releasing of a model checkpoint, or other means that are
747 appropriate to the research performed.
- 748 • While NeurIPS does not require releasing code, the conference does require all submis-
749 sions to provide some reasonable avenue for reproducibility, which may depend on the
750 nature of the contribution. For example
 - 751 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
752 to reproduce that algorithm.
 - 753 (b) If the contribution is primarily a new model architecture, the paper should describe
754 the architecture clearly and fully.
 - 755 (c) If the contribution is a new model (e.g., a large language model), then there should
756 either be a way to access this model for reproducing the results or a way to reproduce
757 the model (e.g., with an open-source dataset or instructions for how to construct
758 the dataset).
 - 759 (d) We recognize that reproducibility may be tricky in some cases, in which case
760 authors are welcome to describe the particular way they provide for reproducibility.
761 In the case of closed-source models, it may be that access to the model is limited in
762 some way (e.g., to registered users), but it should be possible for other researchers
763 to have some path to reproducing or verifying the results.

764 **5. Open access to data and code**

765 Question: Does the paper provide open access to the data and code, with sufficient instruc-
766 tions to faithfully reproduce the main experimental results, as described in supplemental
767 material?

769 Answer: [Yes]

770 Justification: Pseudo-labeled derived dataset is released along with code to recreate it.

771 Guidelines:

- 772 • The answer NA means that paper does not include experiments requiring code.
- 773 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 774 • While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- 775 • The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 776 • The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 777 • The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- 778 • At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- 779 • Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

791 6. Experimental Setting/Details

792 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
793 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
794 results?

795 Answer: [Yes]

796 Justification: Dataset splits, hyperparameters, optimizer are specified.

797 Guidelines:

- 798 • The answer NA means that the paper does not include experiments.
- 799 • The experimental setting should be presented in the core of the paper to a level of detail
800 that is necessary to appreciate the results and make sense of them.
- 801 • The full details can be provided either with the code, in appendix, or as supplemental
802 material.

803 7. Experiment Statistical Significance

804 Question: Does the paper report error bars suitably and correctly defined or other appropriate
805 information about the statistical significance of the experiments?

806 Answer: [No]

807 Justification: The results are represented in by the intersection over union values, there are
808 no error bars.

809 Guidelines:

- 810 • The answer NA means that the paper does not include experiments.
- 811 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
812 dence intervals, or statistical significance tests, at least for the experiments that support
813 the main claims of the paper.
- 814 • The factors of variability that the error bars are capturing should be clearly stated (for
815 example, train/test split, initialization, random drawing of some parameter, or overall
816 run with given experimental conditions).
- 817 • The method for calculating the error bars should be explained (closed form formula,
818 call to a library function, bootstrap, etc.)
- 819 • The assumptions made should be given (e.g., Normally distributed errors).

- 820 • It should be clear whether the error bar is the standard deviation or the standard error
 821 of the mean.
 822 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 823 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 824 of Normality of errors is not verified.
 825 • For asymmetric distributions, the authors should be careful not to show in tables or
 826 figures symmetric error bars that would yield results that are out of range (e.g. negative
 827 error rates).
 828 • If error bars are reported in tables or plots, The authors should explain in the text how
 829 they were calculated and reference the corresponding figures or tables in the text.

830 **8. Experiments Compute Resources**

831 Question: For each experiment, does the paper provide sufficient information on the com-
 832 puter resources (type of compute workers, memory, time of execution) needed to reproduce
 833 the experiments?

834 Answer: [Yes]

835 Justification: We mention the A100 GPU, 10GB of memory and 80 hours of run time.

836 Guidelines:

- 837 • The answer NA means that the paper does not include experiments.
- 838 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
 839 or cloud provider, including relevant memory and storage.
- 840 • The paper should provide the amount of compute required for each of the individual
 841 experimental runs as well as estimate the total compute.
- 842 • The paper should disclose whether the full research project required more compute
 843 than the experiments reported in the paper (e.g., preliminary or failed experiments that
 844 didn't make it into the paper).

845 **9. Code Of Ethics**

846 Question: Does the research conducted in the paper conform, in every respect, with the
 847 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

848 Answer: [Yes]

849 Justification: There are no conflicts between the research and the NeurIPS Code of Ethics.

850 Guidelines:

- 851 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 852 • If the authors answer No, they should explain the special circumstances that require a
 853 deviation from the Code of Ethics.
- 854 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
 855 eration due to laws or regulations in their jurisdiction).

856 **10. Broader Impacts**

857 Question: Does the paper discuss both potential positive societal impacts and negative
 858 societal impacts of the work performed?

859 Answer: [Yes]

860 Justification: There are no negative, but there are positive that are mentioned in the paper
 861 such as better tools for public health decision making.

862 Guidelines:

- 863 • The answer NA means that there is no societal impact of the work performed.
- 864 • If the authors answer NA or No, they should explain why their work has no societal
 865 impact or why the paper does not address societal impact.
- 866 • Examples of negative societal impacts include potential malicious or unintended uses
 867 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
 868 (e.g., deployment of technologies that could make decisions that unfairly impact specific
 869 groups), privacy considerations, and security considerations.

- 870 • The conference expects that many papers will be foundational research and not tied
 871 to particular applications, let alone deployments. However, if there is a direct path to
 872 any negative applications, the authors should point it out. For example, it is legitimate
 873 to point out that an improvement in the quality of generative models could be used to
 874 generate deepfakes for disinformation. On the other hand, it is not needed to point out
 875 that a generic algorithm for optimizing neural networks could enable people to train
 876 models that generate Deepfakes faster.
- 877 • The authors should consider possible harms that could arise when the technology is
 878 being used as intended and functioning correctly, harms that could arise when the
 879 technology is being used as intended but gives incorrect results, and harms following
 880 from (intentional or unintentional) misuse of the technology.
- 881 • If there are negative societal impacts, the authors could also discuss possible mitigation
 882 strategies (e.g., gated release of models, providing defenses in addition to attacks,
 883 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
 884 feedback over time, improving the efficiency and accessibility of ML).

885 11. Safeguards

886 Question: Does the paper describe safeguards that have been put in place for responsible
 887 release of data or models that have a high risk for misuse (e.g., pretrained language models,
 888 image generators, or scraped datasets)?

889 Answer: [NA]

890 Justification: There are no risks for misuse.

891 Guidelines:

- 892 • The answer NA means that the paper poses no such risks.
- 893 • Released models that have a high risk for misuse or dual-use should be released with
 894 necessary safeguards to allow for controlled use of the model, for example by requiring
 895 that users adhere to usage guidelines or restrictions to access the model or implementing
 896 safety filters.
- 897 • Datasets that have been scraped from the Internet could pose safety risks. The authors
 898 should describe how they avoided releasing unsafe images.
- 899 • We recognize that providing effective safeguards is challenging, and many papers do
 900 not require this, but we encourage authors to take this into account and make a best
 901 faith effort.

902 12. Licenses for existing assets

903 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
 904 the paper, properly credited and are the license and terms of use explicitly mentioned and
 905 properly respected?

906 Answer: [Yes]

907 Justification: The raw NOAA datasets used to create SmokeViz do not have licenses while
 908 the python packages used do, we list these in the appendix.

909 Guidelines:

- 910 • The answer NA means that the paper does not use existing assets.
- 911 • The authors should cite the original paper that produced the code package or dataset.
- 912 • The authors should state which version of the asset is used and, if possible, include a
 913 URL.
- 914 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 915 • For scraped data from a particular source (e.g., website), the copyright and terms of
 916 service of that source should be provided.
- 917 • If assets are released, the license, copyright information, and terms of use in the
 918 package should be provided. For popular datasets, paperswithcode.com/datasets
 919 has curated licenses for some datasets. Their licensing guide can help determine the
 920 license of a dataset.
- 921 • For existing datasets that are re-packaged, both the original license and the license of
 922 the derived asset (if it has changed) should be provided.

- 923 • If this information is not available online, the authors are encouraged to reach out to
924 the asset's creators.

925 **13. New Assets**

926 Question: Are new assets introduced in the paper well documented and is the documentation
927 provided alongside the assets?

928 Answer: [Yes]

929 Justification: The dataset, supporting code and user-friendly Notebooks to play with the
930 dataset/model all support the assets accessibility.

931 Guidelines:

- 932 • The answer NA means that the paper does not release new assets.
933 • Researchers should communicate the details of the dataset/code/model as part of their
934 submissions via structured templates. This includes details about training, license,
935 limitations, etc.
936 • The paper should discuss whether and how consent was obtained from people whose
937 asset is used.
938 • At submission time, remember to anonymize your assets (if applicable). You can either
939 create an anonymized URL or include an anonymized zip file.

940 **14. Crowdsourcing and Research with Human Subjects**

941 Question: For crowdsourcing experiments and research with human subjects, does the paper
942 include the full text of instructions given to participants and screenshots, if applicable, as
943 well as details about compensation (if any)?

944 Answer: [NA]

945 Justification: The paper does not involve crowdsourcing nor research with human subjects.

946 Guidelines:

- 947 • The answer NA means that the paper does not involve crowdsourcing nor research with
948 human subjects.
949 • Including this information in the supplemental material is fine, but if the main contribu-
950 tion of the paper involves human subjects, then as much detail as possible should be
951 included in the main paper.
952 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
953 or other labor should be paid at least the minimum wage in the country of the data
954 collector.

955 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human
956 Subjects**

957 Question: Does the paper describe potential risks incurred by study participants, whether
958 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
959 approvals (or an equivalent approval/review based on the requirements of your country or
960 institution) were obtained?

961 Answer: [NA]

962 Justification: The paper does not involve crowdsourcing nor research with human subjects.

963 Guidelines:

- 964 • The answer NA means that the paper does not involve crowdsourcing nor research with
965 human subjects.
966 • Depending on the country in which research is conducted, IRB approval (or equivalent)
967 may be required for any human subjects research. If you obtained IRB approval, you
968 should clearly state this in the paper.
969 • We recognize that the procedures for this may vary significantly between institutions
970 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
971 guidelines for their institution.
972 • For initial submissions, do not include any information that would break anonymity (if
973 applicable), such as the institution conducting the review.