
SmokeViz: A Large-Scale Satellite Dataset for Wildfire Smoke Detection and Segmentation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The global rise in wildfire frequency and intensity over the past decade underscores
2 the need for improved fire monitoring techniques. To advance deep learning re-
3 search on wildfire detection and its associated human health impacts, we introduce
4 **SmokeViz**, a large-scale machine learning dataset of smoke plumes in satellite
5 imagery. The dataset is derived from expert annotations created by smoke analysts
6 at the National Oceanic and Atmospheric Administration, which provide coarse
7 temporal and spatial approximations of smoke presence. To enhance annotation
8 precision, we propose **pseudo-label dimension reduction (PLDR)**, a generalizable
9 method that applies pseudo-labeling to refine datasets with mismatching temporal
10 and/or spatial resolutions. Unlike typical pseudo-labeling applications that aim to
11 increase the number of labeled samples, PLDR maintains the original labels but
12 increases the dataset quality by solving for intermediary pseudo-labels (IPLs) that
13 align each annotation to the most representative input data. For SmokeViz, a parent
14 model produces IPLs to identify the single satellite image within each annotations
15 time window that best corresponds with the smoke plume. This refinement process
16 produces a succinct and relevant deep learning dataset consisting of over 180,000
17 manual annotations. The SmokeViz dataset is expected to be a valuable resource
18 to develop further wildfire-related machine learning models and is publicly avail-
19 able at <https://noaa-gsl-experimental-pds.s3.amazonaws.com/index.html#SmokeViz/>.
20

21

1 Introduction

22 Due in part to public policy, average fine particulate matter ($PM_{2.5}$) levels in the United States have
23 declined over recent decades [2]. However, from 2010 to 2020, the contribution of wildfire smoke to
24 $PM_{2.5}$ concentrations more than doubled, accounting for up to half of total $PM_{2.5}$ exposure in Western
25 U.S. [6]. This is particularly concerning, as ambient $PM_{2.5}$ is a leading environmental risk factor for
26 adverse health outcomes and premature mortality [10]. These trends/risks highlight the urgent need
27 for scalable and timely smoke monitoring systems to mitigate public health risks.

28 Satellite imagery offers the spatial coverage and temporal frequency needed for large-scale smoke
29 monitoring. In comparison to polar-orbiting satellites like Suomi or Sentinel, geostationary satellites
30 such as the GOES series [11] are especially well-suited to this task, providing persistent observation
31 over fixed regions—essential for capturing the dynamic behavior of wildfire smoke plumes. The
32 high temporal resolution and wide coverage of GOES imagery enable real-time tracking of smoke
33 concentration and movement, supporting air quality assessments and early warning systems.

34 Even with the advances in remote sensing, existing deep learning satellite datasets for wildfire smoke
35 detection face several limitations. They are often small in scale, restricted to specific regions or events,
36 and focus on scene-level classification rather than pixel-level segmentation. Most do not differentiate

37 between smoke density levels, are not publicly available, and lack standardized benchmarks for
38 semantic segmentation. While NOAA’s Hazard Mapping System (HMS) provides a large-scale,
39 expert-labeled dataset, its annotations span multi-hour time windows that vary in duration. This
40 creates a temporal mismatch between the labels and individual satellite frames, complicating their
41 direct use for supervised learning.

42 To address these challenges, we introduce **SmokeViz**, a large-scale satellite dataset for semantic
43 segmentation of wildfire smoke plumes. SmokeViz includes over 180,000 annotated samples derived
44 from GOES-East and GOES-West imagery, aligned with HMS analyst annotations. To resolve the
45 temporal ambiguity in the original labels, we propose a semi-supervised method called **pseudo-label**
46 **dimension reduction (PLDR)**, which uses intermediary pseudo-labels to select the satellite image
47 that best matches each smoke annotation. The resulting dataset provides one-to-one image-to-label
48 pairs with ordinal smoke density masks, suitable for supervised deep learning.

49 **SmokeViz** serves as a benchmark for wildfire smoke segmentation and as a resource for the broader
50 machine learning community working with geospatial, temporal, and remote sensing data. It supports
51 new directions in ordinal segmentation, semi-supervised learning with temporal uncertainty, and
52 pretraining for Earth observation tasks involving dynamic atmospheric phenomena.

53 **Our contributions are:**

- 54 • We introduce **SmokeViz**, the largest satellite-based dataset for wildfire smoke segmentation,
55 with over 180,000 samples from GOES imagery.
- 56 • We propose **PLDR**, a physics-guided semi-supervised method for aligning coarse human
57 annotations with temporally optimal satellite imagery.
- 58 • We provide benchmark segmentation baselines and standardized training splits to support
59 reproducibility and downstream research.

60 **2 Related Work**

61 **2.1 Smoke Detection and Labeling Methods**

62 Multi-channel thresholding remains a widely used method for distinguishing smoke from similar
63 atmospheric signatures such as dust or clouds using channel-specific radiance values [33]. These
64 thresholds are typically derived from labeled historical data and are fine-tuned to specific regions
65 and fuel types, limiting their generalizability [21]. In contrast, the SmokeViz dataset spans a wide
66 range of biogeographies across North America and can serve as a source of refined analyst-labeled
67 examples for developing more generalizable thresholding techniques.

68 Large parameterized numerical models are used for forecasting smoke dispersion, but not for smoke
69 detection itself. Systems such as HRRR-Smoke and RRFS [14, 1] rely on computationally intensive
70 forecasts requiring nearly 200 dynamic meteorological inputs. A key limitation of these models is
71 the absence of a real-time smoke analysis product for data assimilation, resulting in delayed model
72 spin-up and compounded forecast errors. Model predictions from SmokeViz could help fill this gap,
73 offering a real-time, satellite-driven alternative to support data assimilation for operational smoke
74 dispersion forecasting.

75 Manual smoke labeling is performed by trained analysts through visual inspection of satellite imagery.
76 NOAA’s Hazard Mapping System (HMS) provides a analyst-labeled wildfire smoke dataset [19, 25].
77 HMS analysts examine GOES imagery sequences to track smoke plume movement and annotate the
78 approximate spatial extent and qualitative density of smoke (light, medium, heavy), as illustrated
79 in Figure 2.1. Annotations are issued on a rolling basis and span time windows ranging from
80 instantaneous to over 20 hours [20]. While HMS provides high-quality expert annotations, its
81 operational format introduces challenges for supervised learning: annotations are temporally coarse,
82 vary in length, and lack one-to-one correspondence with satellite frames. SmokeViz refines HMS
83 annotations into temporally resolved, frame-aligned labels, enabling real-time, continuous predictions
84 of smoke extent and density.

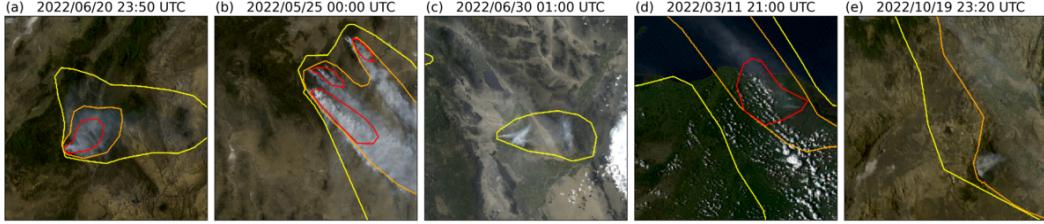


Figure 1: HMS smoke annotations overlaid on GOES imagery. Yellow, orange, and red contours indicate light, medium, and heavy smoke density, respectively. (a) and (b) show canonical smoke plumes; (c)–(e) illustrate density label variation across scenes.

85 2.2 Deep Learning Datasets and Models for Wildfire Smoke

86 Recent efforts have applied deep learning to wildfire smoke detection using a variety of satellite
 87 sources and label strategies. SmokeNet [3] employs a convolutional neural network (CNN) to classify
 88 MODIS image scenes as containing smoke or not, using student-provided labels. SatlasPretrain [5]
 89 includes a small set of Sentinel-2 images labeled for smoke as part of a larger multi-label pretraining
 90 dataset. While scene classification methods can provide wildfire detection information, they do
 91 not capture spatial characteristics of smoke plumes that segmentation would be more appropriate to
 92 capture.

93 Several datasets have been developed for smoke segmentation, but they are limited in scope. Wen et al.
 94 [29] trained a CNN on GOES-East imagery over California and Nevada using HMS annotations from
 95 the 2018 wildfire season. Larsen et al. [16] used Himawari-8 data to detect smoke at the pixel level for
 96 a single fire event, using a threshold-based algorithm as ground truth. Table 1 compares these datasets
 97 in terms of scale, source, and labeling. SmokeViz stands out by offering over 180,000 samples with
 98 analyst-generated, frame-aligned labels covering multiple fire seasons, regions, and biogeographies.
 99 Not only do we use geostationary satellites with persistent observations, but we choose either GOES-
 100 East or GOES-West based on which satellite has optimal observational conditions of the event. It is,
 101 to our knowledge, the largest and most diverse dataset for smoke plume segmentation.

Table 1: Comparison of satellite smoke plume datasets, detailing the number of smoke plume samples, satellite source (polar orbiting (P) or geostationary (G)), number of spectral bands, labeling method, classification type - scene classification (SC) or semantic segmentation (SS), and public availability.

reference	# samples	satellite	# bands	label	task	avail.
[3]	1016	MODIS (P)	5	students	SC	no
[5]	125	Sentinel-2 (P)	3	crowd sourced	SC	yes
[29]	4095	GOES-East (G)	5	HMS analysts	SS	no
[16]	975	Himawari-8 (G)	7	algorithm	SS	no
SmokeViz	183,672	GOES-East+West (G)	3	HMS analysts	SS	yes

102 In addition to its relevance for wildfire applications, SmokeViz contributes a challenging benchmark
 103 for general-purpose remote sensing vision tasks. Unlike many existing datasets that avoid cloudy
 104 scenes [27, 13] or focus on sharply bounded features such as cropland [13], infrastructure [34], or
 105 oceanic clouds [15, 26], smoke has amorphous, fading boundaries in both space and time. Incorpor-
 106 ating smoke segmentation into large-scale pretraining corpora, such as SatlasPretrain [5], could
 107 enhance generalizable models for Earth observation.

108 2.3 Pseudo-labeling and Semi-Supervised Learning

109 Semi-supervised learning techniques such as pseudo-labeling have been widely used to expand
 110 training data by leveraging unlabeled samples [17]. Typically, a parent model is trained on labeled
 111 data and then used to generate pseudo-labels for an unlabeled dataset, which are in turn used to train
 112 subsequent models in an iterative process.

113 In contrast, we propose a non-iterative variation focused not on data expansion, but dataset data-to-label
114 precision. Our method, **pseudo-label dimension reduction (PLDR)**, generates intermediary
115 pseudo-labels (IPLs) for each satellite frame within the HMS annotation window. Rather than using
116 these labels for training, we use them to identify the satellite image with the greatest alignment to
117 the analyst annotation. This enables the construction of SmokeViz, a temporally disambiguated,
118 one-to-one image-to-label dataset. The resulting dataset methodically pairs the analyst-generated
119 smoke plume labels with selected GOES imagery, enabling high-resolution, temporally accurate
120 segmentation model training.

121 Beyond wildfire smoke segmentation, PLDR offers a general framework for aligning coarse or weakly
122 matched datasets. This is particularly useful in domains such as remote sensing, medical imaging,
123 and video analysis, where annotations often span temporal or spatial intervals rather than individual
124 frames. In Earth observation specifically, atmospheric parameters are often combined from disparate
125 sources with inconsistent spatial and temporal resolutions, making it difficult to integrate them into
126 unified training datasets. By using intermediary pseudo-labels to identify the most representative
127 input sample, PLDR transforms many-to-one or one-to-many supervision into clean one-to-one
128 mappings. This enables more precise alignment between data and labels, facilitating integration
129 across heterogeneous sources without requiring new human annotations. As such, PLDR serves as a
130 practical preprocessing strategy for repurposing legacy datasets with temporal ambiguity into precise
131 training resources for modern deep learning models.

132 3 Methods

133 3.1 Datasets

134 We use imagery from the latest GOES satellites—GOES-16 (East), GOES-17, and GOES-18
135 (West)—each equipped with the Advanced Baseline Imager (ABI), which captures 16 spectral
136 bands from visible to infrared wavelengths every 10 minutes. We process bands 1–3 using PyTroll
137 [23] to generate 1 km true-color composites [4], matching the imagery reviewed by HMS ana-
138 lists. These bands correspond to the shortest wavelengths available on ABI and yield the highest
139 signal-to-noise ratio (SNR).

140 To approximate the dynamic movement of smoke, HMS analysts annotate plumes using multi-frame
141 satellite animations. These annotations span varying time windows, averaging three hours. Since the
142 HMS annotations are designed to reflect overall plume extent during a time window rather than at
143 any specific moment, smoke boundaries in individual frames may not align well with the annotation
144 (Figure 2). A naive modeling approach would use all frames within each time window as input, but
145 this introduces non-uniform sequence lengths and significantly increases memory and computational
146 demands and complicates the use of CNN architectures. Instead, we establish a one-to-one mapping
147 by identifying the single satellite frame that best matches each analyst annotation.

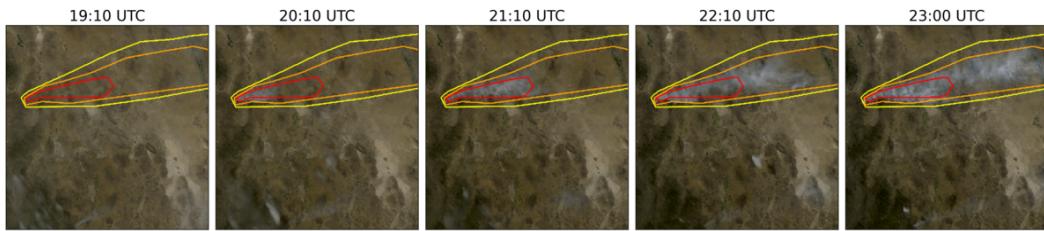


Figure 2: True color GOES-East imagery from May 5th, 2022, Southeast New Mexico (31.38°N , 107.87°W) during the start of the Foster Fire. The red, orange and yellow lines represent the heavy, medium and low density HMS smoke annotations that span 19:10–23:00 UTC.

148 We select either GOES-East or GOES-West based on the solar zenith angle (SZA) to optimize
149 for forward Mie scattering, which enhances smoke visibility in satellite imagery. Smoke particles
150 ($100\text{ nm}–10\text{ }\mu\text{m}$) scatter light predominantly via Mie scattering when $\lambda < d$, favoring short wave-
151 lengths and forward angles (figure 3). To generate the Mie-derived dataset, we evaluate the available
152 satellite platforms for each annotation time window and choose the satellite (East or West) that is
153 expected to observe the strongest forward scattering geometry based on sun-satellite alignment. This

154 ensures selection of the satellite view with the highest potential smoke SNR if smoke were present.
 155 Therefore, we select (1) the satellite expected to yield the strongest Mie forward scattering (Figures
 156 4(a) vs 4(b)) and (2) the three shortest wavelength ABI bands (C01–C03: 0.47, 0.64, and 0.865 μm)
 157 (Figures 4(c)-4(e)).

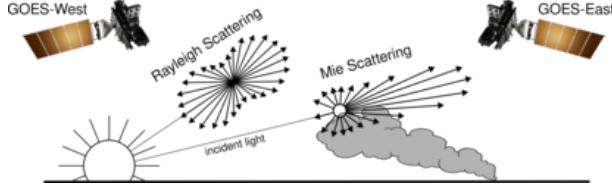


Figure 3: If the particle size is $< \frac{1}{10}$ the λ of the interacting light, then the primary scattering will be Rayleigh. Mie scattering is the predominant scattering mechanism when the particle size is larger than the λ of light. This schematic demonstrates that when the sun is setting in the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

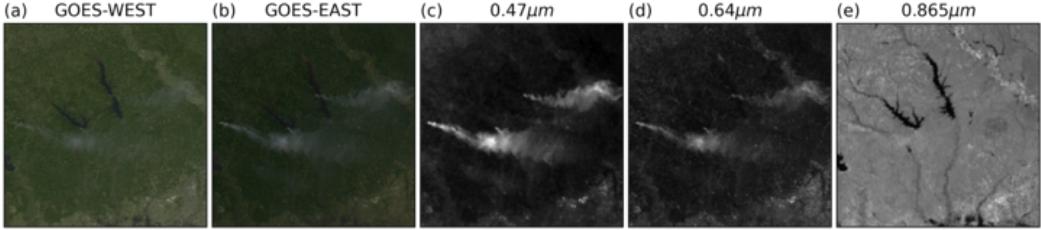


Figure 4: True color (a) GOES-WEST and (b) GOES-EAST imagery from March 23rd, 2022 centered at $(31.1^\circ, -93.8^\circ)$ in Texas, USA taken at 23:20 UTC. The GOES-EAST raw band imagery for (c) blue, (d) red and (e) veggie bands show variations in the SNR for smoke detection in relation to the λ of light being measured.

158 3.1.1 From Full Dataset \mathcal{D} to Mie-Derived Dataset \mathcal{D}_M

159 Let $\mathcal{D} = \{\mathcal{X}, \mathcal{Y}\}$ be the original dataset, where each label $y_i \in \mathcal{Y}$ corresponds to multiple satellite
 160 images $[x_{(i,t_0)}, \dots, x_{(i,t_N)}] \in \mathcal{X}$ over a given time window. Using Mie scattering principles, we select
 161 the image $x_{(i,t_M)}$ with the highest expected smoke SNR to form a one-to-one dataset $\mathcal{D}_M = \{\mathcal{X}_M, \mathcal{Y}\}$
 162 such that $\mathcal{X}_M \subset \mathcal{X}$ and $|\mathcal{X}_M| = |\mathcal{Y}|$. Based on forward scattering criteria, the trivial strategy would
 163 be to pull imagery from GOES-West right after sunrise and from GOES-East right before sunset
 164 when the SZA is closest to 90° . To avoid image artifacts caused by extreme SZA, we exclude scenes
 165 with $\text{SZA} > 88^\circ$ [24]. The resulting dataset \mathcal{D}_M (Table 3) contains over 200,000 samples where the
 166 satellite image is chosen based on which frame within the annotation time window would exhibit
 167 the strongest forward scattering geometry and thus the highest potential smoke SNR if smoke were
 168 present.

169 3.1.2 PLDR Dataset \mathcal{D}_p

170 The \mathcal{D}_M data selection process introduces a potential bias for resulting models to limit smoke
 171 identification to higher SZAs. Additionally, \mathcal{D}_M is limited to providing the timestamp for maximum
 172 possible smoke SNR, it does not give information to point to which image aligns best with the
 173 smoke label. To address these limitations, we propose using \mathcal{D}_M as a intermediary dataset in the
 174 PLDR workflow (Figure 5) that will predict the satellite image that best matches the analyst's smoke
 175 annotation to produce \mathcal{D}_p .

176 To build f_o , we implement Segmentation Models PyTorch [12] with EfficientNetV2 [28] as the
 177 encoder and PSPNet [32] as the decoder. Input images are 256×256 true-color snapshots; the output
 178 is a 256×256 classification map predicting categorical smoke density. We use thermometer encoding
 179 (Table 2) and apply binary cross-entropy loss across density levels. Thermometer encoding is chosen
 180 over one-hot encoding because it captures the ordinal structure of smoke density categories (none

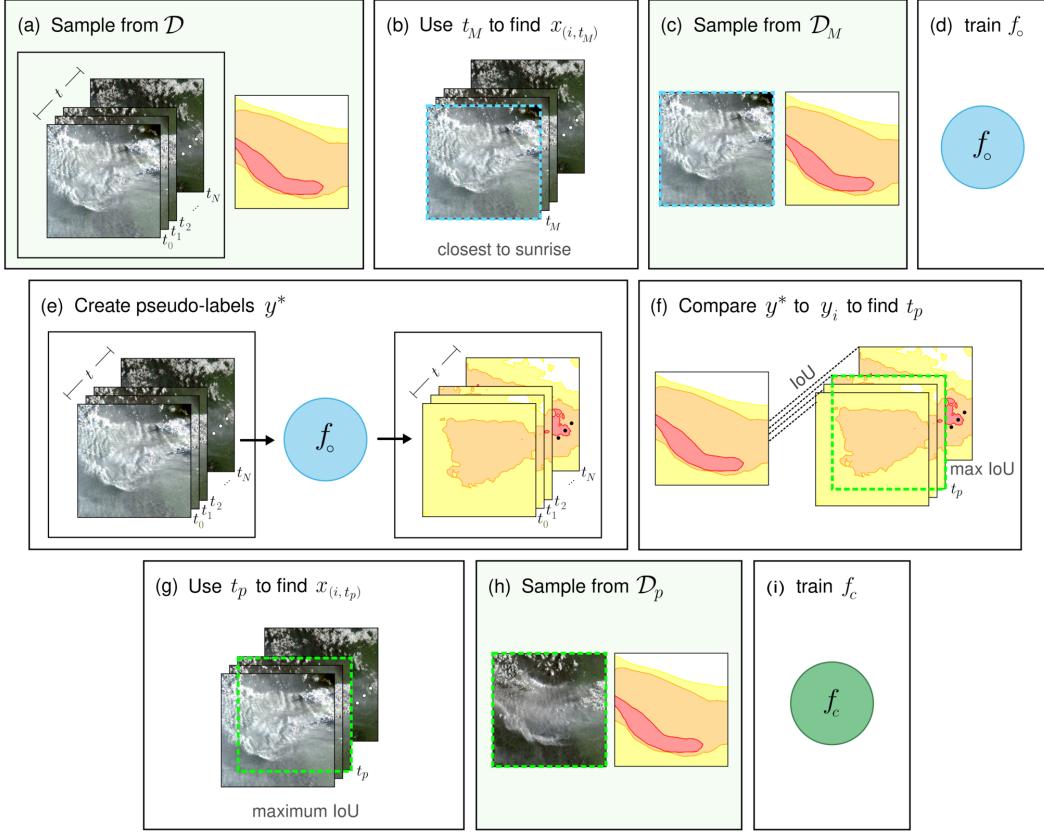


Figure 5: PLDR applied to create the SmokeViz dataset. Green boxes indicate dataset stages. (a) For original dataset \mathcal{D} - analyst annotation y_i corresponds to N satellite images across time window t so that $([x_{(i,t_0)}, \dots, x_{(i,t_N)}], y_i) \in \mathcal{D}$; (b) use Mie scattering to find the time, t_M , that corresponds with satellite image $x_{(i,t_M)}$ that would produce the highest possible SNR if smoke was present; (c) resulting \mathcal{D}_M is one-to-one $(x_{(i,t_M)}, y_i) \in \mathcal{D}_M$; (d) parent model f_o is trained on \mathcal{D}_M such that $f_o(x_{(i,t_M)}) = y_i$; (e) apply a greedy algorithm $f_o([x_{(i,t_0)}, \dots, x_{(i,t_N)}]) = [y_{(i,t_0)}^*, \dots, y_{(i,t_N)}^*]$ to create IPLs y^* for each candidate image; (f) compute the intersection over union (IoU) between y^* and y_i to identify the time t_p where the IPL and analyst annotation have the maximum IoU; (g) match t_p to its corresponding image $x_{(i,t_p)}$ that is predicted to best match the analyst annotation; (h) SmokeViz dataset \mathcal{D}_p created; (i) child model f_c is trained on \mathcal{D}_p such that $f_c(x_{(i,t_p)}) = y_i$ is used to detect and classify the density of wildfire smoke plumes in GOES imagery.

Table 2: A comparison of how smoke density would be represented by one-hot encoding commonly used for categorical data to thermometer encoding often used for ordinal data.

density	one-hot	thermometer
none	[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]

Table 3: Dataset split for \mathcal{D}_M and \mathcal{D}_p , samples for 2024 go up to November 1st **do 2018-2022 instead**. We use an entire year of data for both validation and testing sets to capture year-long wildfire trends.

dataset	\mathcal{D}_M	\mathcal{D}_p	years
training	165,609	144,225	2018-22
validation	20,056	19,223	2023
testing	21,541	20,224	2022

181 < light < medium < heavy). In thermometer encoding, each higher class includes all lower class
 182 activations (e.g., heavy = [1 1 1]), allowing the model to learn not just class distinctions, but the
 183 relative severity of smoke. We use a confidence threshold of IoU > 0.01 [9] to exclude samples with
 184 negligible overlap.

185 **3.2 Benchmark Models**

186 We benchmark the SmokeViz dataset \mathcal{D}_p using DeepLabV3+ [7] and PSPNet [32] with EfficientNetV2
 187 [28], DPT [22] with ViT [8], Segformer [31] and UperNet [30] with EfficientVit [18]. Each model is
 188 trained for 100 epochs using a batch size of 16 and the Adam optimizer on 8 16GB Nvidia P100 GPUs.
 189 These architectures are selected for their relatively low memory requirements and effectiveness in
 190 segmenting multi-scale objects such as smoke plumes.

191 **4 Results**

192 We evaluate the performance of f_o and f_c using Intersection over Union (IoU) metrics on the test
 193 sets of both \mathcal{D}_M and \mathcal{D}_p , as shown in Table 4. For each smoke density class, IoU is calculated as the
 194 pixel-level intersection between model predictions and HMS analyst labels, divided by their union,
 195 aggregated over all test samples. Overall IoU is computed by summing intersections across all density
 196 classes and dividing by the total union of predicted and labeled smoke pixels.

Table 4: IoU results per smoke density and overall, comparing f_o and f_c run on \mathcal{D}_M and \mathcal{D}_p test sets.

	f_o		f_c	
	\mathcal{D}_M	\mathcal{D}_p	\mathcal{D}_M	\mathcal{D}_p
heavy	0.278	0.368	0.218	0.411
medium	0.310	0.417	0.319	0.484
light	0.480	0.585	0.491	0.660
overall	0.430	0.533	0.438	0.607

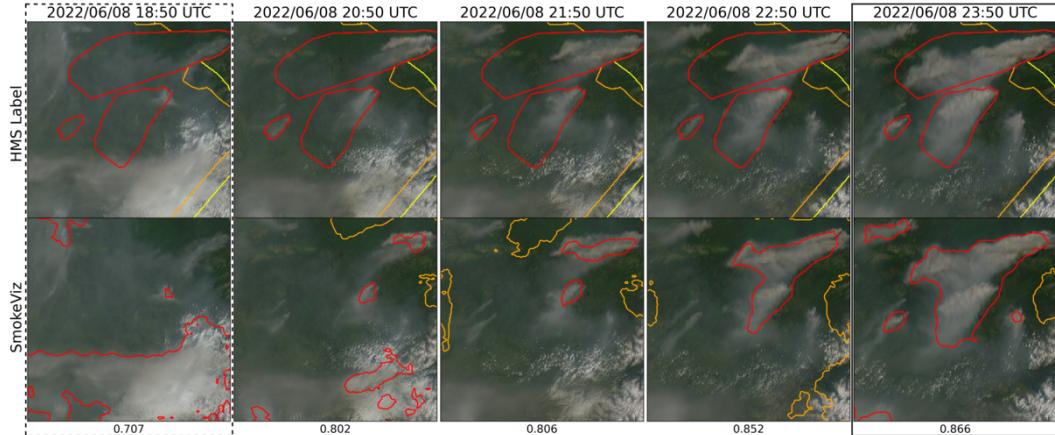


Figure 6: GOES-West imagery from June 8, 2022, over Alaska (61.06°N , 156.12°W). Daylight spanned 12:43–7:53 UTC. The HMS annotations (top row) span 18:50–23:50 UTC and are compared with f_o -generated smoke predictions (bottom row). The leftmost frame (dotted) represents the Mie-derived image; the rightmost frame (solid) was selected via PLDR and achieves higher IoU.

197 Figure 6 illustrates a case in which the PLDR-selected frame better represents the HMS annotation
 198 than the Mie-derived selection. Here, the heavy smoke IoU improves from 0.01 to 0.59. While the
 199 Mie-derived image is selected based on its proximity to sunrise, PLDR chooses the frame with the
 200 highest overlap between the model-generated intermediary pseudo-label and the analyst annotation.
 201 This example highlights PLDR’s advantage in resolving temporal ambiguity.

202 To further examine the performance of f_c , we can qualitatively compare its predictions against HMS
 203 annotations for samples from \mathcal{D}_p in Figure 7. The model outputs capture more spatially detailed and
 204 coherent smoke boundaries compared to the coarser, polygon-based analyst labels.

205 To benchmark performance across segmentation architectures, we evaluate several encoder-decoder
 206 models trained on \mathcal{D}_p . Table 5 reports IoU scores by smoke density and overall. While DeepLabV3+

Table 5: Comparison of segmentation benchmark model IoU metrics on \mathcal{D}_p .

encoder decoder	EfficientNetV2 [28] DeepLabV3+ [7]	[28] PSPNet [32]	ViT [8] DPT [22]	EfficientViT [18] Segformer [31]	[18] UperNet [30]
heavy	0.2894	0.3222	0.2091	0.2185	0.3099
medium	0.4091	0.4289	0.3946	0.3978	0.4042
light	0.4424	0.5045	0.5155	0.4331	0.4275
overall	0.4172	0.4677	0.4608	0.4055	0.4098

207 achieves the highest IoU for heavy smoke, PSPNet yields the best overall performance. Results across
 208 models are relatively consistent, highlighting the robustness of the SmokeViz dataset for training
 209 diverse architectures.

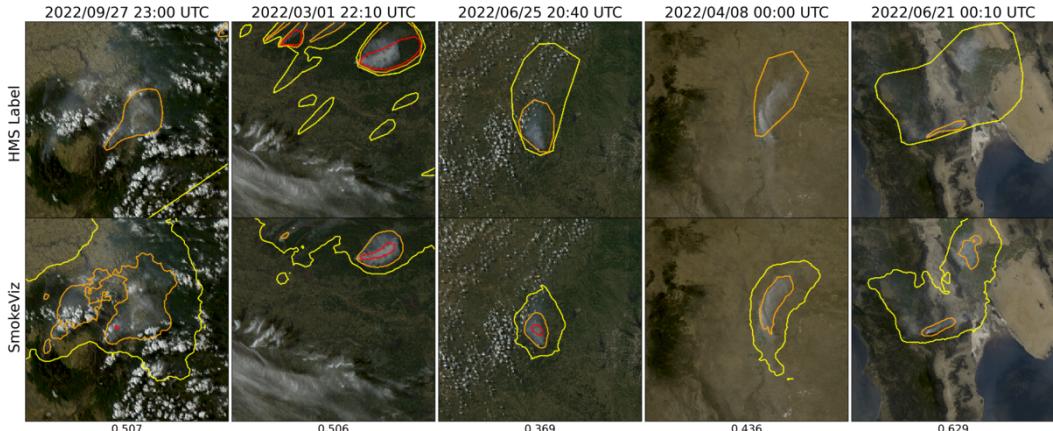


Figure 7: Examples of HMS annotations (top row) vs f_c output (bottom row) on \mathcal{D}_p samples. The overall IoU score is reported at the bottom of each column.

210 5 Limitations

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

219 6 Conclusion

220 In this study, we have refined an existing dataset originally curated by NOAA’s HMS team, trans-
 221 forming it from a many-to-one imagery-to-annotation format to a more succinct, one-to-one satellite
 222 image-to-annotation dataset. The initial HMS dataset provided a general approximation of where
 223 smoke had been present for a given time window, though it did not guarantee the actual existence
 224 of smoke in the labeled pixels during the given times. Our goal was to create a dataset that could
 225 be used, along with additional applications, to train a model to detect wildfire smoke in real-time
 226 on an image-by-image level. The Mie-derived dataset selection process determined that if smoke
 227 was present, what timestamp within the analyst time window would give the highest smoke
 228 signal-to-noise ratio. While optimizing for being able to detect smoke, if it is present, the Mie-dataset
 229 selection had no metric to determine if the smoke was effectively present in the selected image. Since
 230 many of the images within the HMS time-window either contained no smoke at all or the smoke was

231 not contained within the geospatial bounds of the annotations, the Mie-derived dataset contained
232 a large number of mislabeled samples. Discrepancies between data and labels can be detrimental
233 towards the model’s capacity to improve on feature representations in the target domain. During
234 model training, the penalization of accurate predictions can inadvertently introduce biases towards
235 misclassifying noise as meaningful signal.

236 To improve the dataset’s capacity to accurately represent wildfire smoke plumes, we train a parent
237 machine learning model, f_o , using the Mie-derived dataset, \mathcal{D}_M , and run it on the relevant satellite
238 images within the time-frame. The image with the maximum IoU score between the model’s smoke
239 predictions, or pseudo-label, and the analyst smoke annotations are used to create the pseudo-label
240 generated dataset, \mathcal{D}_p . We then train a child model, f_c , using \mathcal{D}_p and test f_o and f_c on both the 2022
241 testing sets from \mathcal{D}_M and \mathcal{D}_p . The results reported in table 4 suggest that \mathcal{D}_p was able to train a better
242 performing model, f_c , that gave higher IoU metrics on both dataset’s testing sets in comparison to
243 the original parent model, f_o .

244 The result of this study is a representative dataset, SmokeViz, that can be used to train machine
245 learning models for various wildfire smoke applications. A future goal is to produce a robust
246 and reliable machine learning based approach for detecting wildfires using satellite imagery. That
247 information can be used for wildfire detection and monitoring in along with a highly needed smoke
248 product for data assimilation into smoke dispersion models. Additionally, this dataset can be used as
249 a benchmark for how well remote sensing segmentation models can perform on dispersed edges such
250 as smoke. On a broader scale, we show how pseudo-labeling can be used to optimize a dataset when
251 the resolution for the data and corresponding labels do not match. This could be useful in similar
252 applications involving time-series/video data with a singular label where the data can be compressed
253 while still remaining representative of the label. All data is made publicly available at [aws download
254 link] and all code can be found at <https://github.com/anonymous-smokeviz/SmokeViz>.

255 References

- 256 [1] R. Ahmadov, H. Li, J. Romero-Alvarez, J. Schnell, S. Bhimireddy, E. James, K. Y. Wong,
257 M. Hu, J. Carley, P. Bhattacharjee, et al. Forecasting smoke and dust in noaa’s next-generation
258 high-resolution coupled numerical weather prediction model. Technical report, Copernicus
259 Meetings, 2024.
- 260 [2] J. E. Aldy, M. Auffhammer, M. Cropper, A. Fraas, and R. Morgenstern. Looking back at 50
261 years of the clean air act. *Journal of Economic Literature*, 60(1):179–232, 2022.
- 262 [3] R. Ba, C. Chen, J. Yuan, W. Song, and S. Lo. Smokenet: Satellite smoke scene detection using
263 convolutional neural network with spatial and channel-wise attention. *Remote Sensing*, 11(14):
264 1702, 2019.
- 265 [4] M. Bah, M. Gunshor, and T. Schmit. Generation of goes-16 true color imagery without a green
266 band. *Earth and Space Science*, 5(9):549–558, 2018.
- 267 [5] F. Bastani, P. Wolters, R. Gupta, J. Ferdinando, and A. Kembhavi. Satlaspretrain: A large-scale
268 dataset for remote sensing image understanding. In *Proceedings of the IEEE/CVF International
269 Conference on Computer Vision*, pages 16772–16782, 2023.
- 270 [6] M. Burke, A. Driscoll, S. Heft-Neal, J. Xue, J. Burney, and M. Wara. The changing risk and
271 burden of wildfire in the united states. *Proceedings of the National Academy of Sciences*, 118
272 (2):e2011048118, 2021.
- 273 [7] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic
274 image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs.
275 *IEEE transactions on pattern analysis and machine intelligence*, 40(4):834–848, 2017.
- 276 [8] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,
277 M. Minderer, G. Heigold, S. Gelly, et al. An image is worth 16x16 words: Transformers for
278 image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- 279 [9] R. E. Ferreira, Y. J. Lee, and J. R. Dórea. Using pseudo-labeling to improve performance of
280 deep neural networks for animal identification. *Scientific Reports*, 13(1):13875, 2023.

- 281 [10] E. Gakidou, A. Afshin, A. A. Abajobir, K. H. Abate, C. Abbafati, K. M. Abbas, F. Abd-Allah,
 282 A. M. Abdulle, S. F. Abera, V. Aboyans, et al. Global, regional, and national comparative risk
 283 assessment of 84 behavioural, environmental and occupational, and metabolic risks or clusters
 284 of risks, 1990–2016: a systematic analysis for the global burden of disease study 2016. *The
 285 Lancet*, 390(10100):1345–1422, 2017.
- 286 [11] S. J. Goodman, T. J. Schmit, J. Daniels, and R. J. Redmon. *The GOES-R series: a new
 287 generation of geostationary environmental satellites*. Elsevier, 2019.
- 288 [12] P. Iakubovskii. Segmentation models pytorch. https://github.com/qubvel/segmentation_models.pytorch, 2019.
- 290 [13] J. Jakubik, S. Roy, C. Phillips, P. Fraccaro, D. Godwin, B. Zdrozny, D. Szwarcman, C. Gomes,
 291 G. Nyirjesy, B. Edwards, et al. Foundation models for generalist geospatial artificial intelligence.
 292 arxiv 2023. *arXiv preprint arXiv:2310.18660*.
- 293 [14] E. James, R. Ahmadov, and G. A. Grell. Realtime wildfire smoke prediction in the united states:
 294 The hrrr-smoke model. In *EGU General Assembly Conference Abstracts*, page 19526, 2018.
- 295 [15] A. Kitamoto, J. Hwang, B. Vuillod, L. Gautier, Y. Tian, and T. Clanuwat. Digital typhoon: Long-
 296 term satellite image dataset for the spatio-temporal modeling of tropical cyclones. *Advances in
 297 Neural Information Processing Systems*, 36, 2024.
- 298 [16] A. Larsen, I. Hanigan, B. J. Reich, Y. Qin, M. Cope, G. Morgan, and A. G. Rappold. A deep
 299 learning approach to identify smoke plumes in satellite imagery in near-real time for health risk
 300 communication. *Journal of exposure science & environmental epidemiology*, 31(1):170–176,
 301 2021.
- 302 [17] D.-H. Lee. Pseudo-label : The simple and efficient semi-supervised learning method for deep
 303 neural networks. *ICML 2013 Workshop : Challenges in Representation Learning (WREPL)*, 07
 304 2013.
- 305 [18] X. Liu, H. Peng, N. Zheng, Y. Yang, H. Hu, and Y. Yuan. Efficientvit: Memory efficient vision
 306 transformer with cascaded group attention. In *Proceedings of the IEEE/CVF conference on
 307 computer vision and pattern recognition*, pages 14420–14430, 2023.
- 308 [19] D. McNamara, G. Stephens, M. Ruminski, and T. Kasheta. The hazard mapping system (hms) -
 309 noaa’s multi-sensor fire and smoke detection program using environmental satellites. *Conference
 310 on Satellite Meteorology and Oceanography*, 01 2004.
- 311 [20] NOAA. Hazard mapping system fire and smoke product. URL <https://www.ospo.noaa.gov/Products/land/hms.html#about>.
- 312 [21] T. Randriambelo, S. Baldy, M. Bessafi, M. Petit, and M. Despinoy. An improved detection
 313 and characterization of active fires and smoke plumes in south-eastern africa and madagascar.
 314 *International Journal of Remote Sensing*, 19(14):2623–2638, 1998.
- 315 [22] R. Ranftl, A. Bochkovskiy, and V. Koltun. Vision transformers for dense prediction. In
 316 *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12179–
 317 12188, 2021.
- 318 [23] M. Raspaud, D. Hoese, 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 [24] A. Royer, P. Vincent, and F. Bonn. Evaluation and correction of viewing angle effects on
 323 satellite measurements of bidirectional reflectance. *Photogrammetric engineering and remote
 324 sensing*, 51(12):1899–1914, 1985.
- 325 [25] W. Schroeder, M. Ruminski, I. Csizar, L. Giglio, E. Prins, C. Schmidt, and J. Morisette.
 326 Validation analyses of an operational fire monitoring product: The hazard mapping system.
 327 *International Journal of Remote Sensing*, 29(20):6059–6066, 2008.

- 329 [26] B. Stevens, S. Bony, H. Brogniez, L. Hentgen, C. Hohenegger, C. Kiemle, T. S. L’Ecuyer, A. K.
330 Naumann, H. Schulz, P. A. Siebesma, et al. Sugar, gravel, fish and flowers: Mesoscale cloud
331 patterns in the trade winds. *Quarterly Journal of the Royal Meteorological Society*, 146(726):
332 141–152, 2020.
- 333 [27] G. Sumbul, M. Charfuelan, B. Demir, and V. Markl. Bigearthnet: A large-scale benchmark
334 archive for remote sensing image understanding. In *IGARSS 2019-2019 IEEE international
335 geoscience and remote sensing symposium*, pages 5901–5904. IEEE, 2019.
- 336 [28] M. Tan and Q. Le. Efficientnetv2: Smaller models and faster training. In *International
337 conference on machine learning*, pages 10096–10106. PMLR, 2021.
- 338 [29] J. Wen and M. Burke. Wildfire smoke plume segmentation using geostationary satellite imagery.
339 *ArXiv*, abs/2109.01637, 2021. URL [https://api.semanticscholar.org/CorpusID:
340 237416777](https://api.semanticscholar.org/CorpusID:237416777).
- 341 [30] T. Xiao, Y. Liu, B. Zhou, Y. Jiang, and J. Sun. Unified perceptual parsing for scene understanding.
342 In *Proceedings of the European conference on computer vision (ECCV)*, pages 418–434, 2018.
- 343 [31] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo. Segformer: Simple and
344 efficient design for semantic segmentation with transformers. *Advances in neural information
345 processing systems*, 34:12077–12090, 2021.
- 346 [32] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *Proceedings of
347 the IEEE conference on computer vision and pattern recognition*, pages 2881–2890, 2017.
- 348 [33] T. X.-P. Zhao, S. Ackerman, and W. Guo. Dust and smoke detection for multi-channel imagers.
349 *Remote Sensing*, 2(10):2347–2368, 2010. ISSN 2072-4292. doi: 10.3390/rs2102347. URL
350 <https://www.mdpi.com/2072-4292/2/10/2347>.
- 351 [34] S. Zorzi, S. Bazrafkan, S. Habenschuss, and F. Fraundorfer. Polyworld: Polygonal building
352 extraction with graph neural networks in satellite images. In *Proceedings of the IEEE/CVF
353 Conference on Computer Vision and Pattern Recognition*, pages 1848–1857, 2022.