
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 [1]. 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 United States [2]. This is particularly concerning, as ambient $PM_{2.5}$ is a leading environmental risk
26 factor for adverse health outcomes and premature mortality [3]. These trends/risks highlight the
27 urgent need 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 [4] 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 pre-training for Earth observation tasks involving dynamic atmospheric phenomena.

53 The contributions presented in this paper include **SmokeViz**, the largest satellite-based dataset for
54 wildfire smoke segmentation, with over 180,000 samples from GOES imagery, our proposed **PLDR**,
55 a physics-guided semi-supervised method for aligning coarse human annotations with temporally
56 optimal satellite imagery and benchmark segmentation baselines with standardized training splits to
57 support reproducibility and future studies.

58 2 Related Work

59 2.1 Smoke Detection and Labeling Methods

60 Multi-channel thresholding remains a widely used method for distinguishing smoke from similar
61 atmospheric signatures such as dust or clouds using channel-specific radiance values [5]. These
62 thresholds are typically derived from labeled historical data and are fine-tuned to specific regions and
63 fuel types, limiting their generalizability [6]. In contrast, the SmokeViz dataset spans a wide range of
64 biogeographies across North America and can serve as a source of refined analyst-labeled examples
65 for developing more generalizable thresholding techniques.

66 Large parameterized numerical models are used for forecasting smoke dispersion, but not for smoke
67 detection itself. Systems such as HRRR-Smoke and RRFS [7, 8] rely on computationally intensive
68 forecasts requiring nearly 200 dynamic meteorological inputs. A key limitation of these models is
69 the absence of a real-time smoke analysis product for data assimilation, resulting in delayed model
70 spin-up and compounded forecast errors. Model predictions from SmokeViz could help fill this gap,
71 offering a real-time, satellite-driven alternative to support data assimilation for operational smoke
72 dispersion forecasting.

73 Manual smoke labeling is performed by trained analysts through visual inspection of satellite imagery.
74 NOAA’s Hazard Mapping System (HMS) provides a analyst-labeled wildfire smoke dataset [9, 10].
75 HMS analysts examine GOES imagery sequences to track smoke plume movement and annotate the
76 approximate spatial extent and qualitative density of smoke (light, medium, heavy), as illustrated
77 in Figure 2.1. Annotations are issued on a rolling basis and span time windows ranging from
78 instantaneous to over 20 hours [11]. While HMS provides high-quality expert annotations, its
79 operational format introduces challenges for supervised learning: annotations are temporally coarse,
80 vary in length, and lack one-to-one correspondence with satellite frames. SmokeViz refines HMS
81 annotations into temporally resolved, frame-aligned labels, enabling real-time, continuous predictions
82 of smoke extent and density.

83 2.2 Deep Learning Datasets and Models for Wildfire Smoke

84 Recent efforts have applied deep learning to wildfire smoke detection using a variety of satellite
85 sources and label strategies. SmokeNet [12] employs a convolutional neural network (CNN) to classify
86 MODIS image scenes as containing smoke or not, using student-provided labels. SatlasPretrain [13]
87 includes a small set of Sentinel-2 images labeled for smoke as part of a larger multi-label pre-training

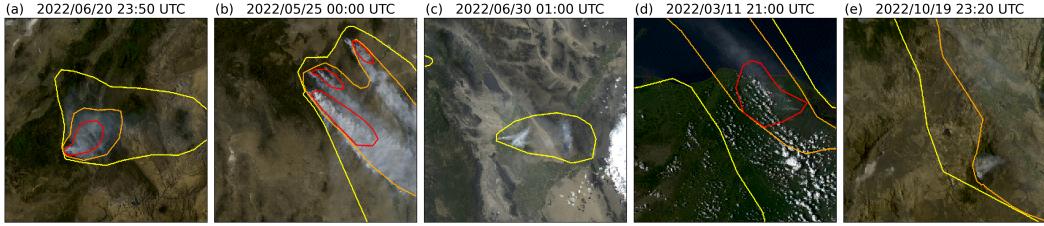


Figure 1: HMS smoke annotations overlaid on GOES imagery. Yellow, orange, and red contours indicate light, medium, and heavy smoke density, respectively. (a) and (b) show canonical smoke plumes; (c)–(e) illustrate density label variation across scenes.

dataset. While scene classification methods can provide wildfire detection information, they do not capture spatial characteristics of smoke plumes that segmentation would be more appropriate to capture.

Several datasets have been developed for smoke segmentation, but they are limited in scope. Wen et al. [14] trained a CNN on GOES-East imagery over California and Nevada using HMS annotations from the 2018 wildfire season. Larsen et al. [15] used Himawari-8 data to detect smoke at the pixel level for a single fire event, using a threshold-based algorithm as ground truth. Table 1 compares these datasets in terms of scale, source, and labeling. SmokeViz stands out by offering over 180,000 samples with analyst-generated, frame-aligned labels covering multiple fire seasons, regions, and biogeographies. Not only do we use geostationary satellites with persistent observations, but we choose either GOES-East or GOES-West based on which satellite has optimal observational conditions of the event. It is, 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.
[12]	1016	MODIS (P)	5	students	SC	no
[13]	125	Sentinel-2 (P)	3	crowd sourced	SC	yes
[14]	4095	GOES-East (G)	5	HMS analysts	SS	no
[15]	975	Himiwari-8 (G)	7	algorithm	SS	no
SmokeViz	183,672	GOES-East+West (G)	3	HMS analysts	SS	yes

In addition to its relevance for wildfire applications, SmokeViz contributes a challenging benchmark for general-purpose remote sensing vision tasks. Unlike many existing datasets that avoid cloudy scenes [16, 17] or focus on sharply bounded features such as cropland [17], infrastructure [18], or oceanic clouds [19, 20], smoke has amorphous, fading boundaries in both space and time. Incorporating smoke segmentation into large-scale pre-training corpora, such as SatlasPretrain [13], could enhance generalizable models for Earth observation.

2.3 Pseudo-labeling and Semi-Supervised Learning

Semi-supervised learning techniques such as pseudo-labeling have been widely used to expand training data by leveraging unlabeled samples [21]. Typically, a parent model is trained on labeled data and then used to generate pseudo-labels for an unlabeled dataset, which are in turn used to train subsequent models in an iterative process.

In contrast, we propose a non-iterative variation focused not on data expansion, but dataset data-to-label precision. Our method, **pseudo-label dimension reduction (PLDR)**, generates intermediary pseudo-labels (IPLs) for each satellite frame within the HMS annotation window. Rather than using these labels for training, we use them to identify the satellite image with the greatest alignment to the analyst annotation. This enables the construction of SmokeViz, a temporally disambiguated, one-to-one image-to-label dataset. The resulting dataset methodically pairs the analyst-generated

117 smoke plume labels with selected GOES imagery, enabling high-resolution, temporally accurate
118 segmentation model training.

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

130 3 Methods

131 3.1 Datasets

132 We use imagery from the latest GOES satellites—GOES-16 (East), GOES-17, and GOES-18 (West),
133 each equipped with the Advanced Baseline Imager (ABI), which captures 16 spectral bands from
134 visible to infrared wavelengths every 10 minutes. We process bands 1-3 using PyTroll [22] to generate
135 1km true-color composites [23], matching the imagery reviewed by HMS analysts. These bands
136 correspond to the shortest wavelengths available on ABI and yield the highest signal-to-noise ratio
137 (SNR).

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

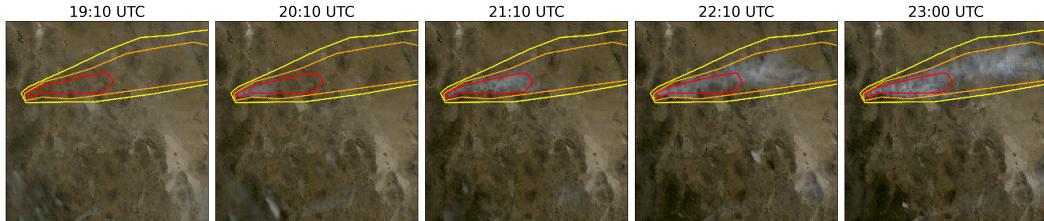


Figure 2: True color GOES-East imagery from May 5th, 2022, Southeast New Mexico (31.38°N , 107.87°W) during the start of the Foster Fire. The red, orange and yellow lines represent the heavy, medium and low density HMS smoke annotations that span 19:10–23:00 UTC.

146 We select either GOES-East or GOES-West based on the solar zenith angle (SZA) to optimize
147 for forward Mie scattering, which enhances smoke visibility in satellite imagery. Smoke particles
148 ($100\text{nm}-10\mu\text{m}$) scatter light predominantly via Mie scattering when $\lambda < d$, favoring short wavelengths
149 and forward angles (Figure 3). To generate the Mie-derived dataset, we evaluate the available satellite
150 platforms for each annotation time window and choose the satellite (East or West) that is expected
151 to observe the strongest forward scattering geometry based on sun-satellite alignment. This ensures
152 selection of the satellite view with the highest potential smoke SNR if smoke were present. Therefore,
153 we select (1) the satellite expected to yield the strongest Mie forward scattering (Figures 4(a) vs
154 4(b)) and (2) the three shortest wavelength ABI bands (C01-C03: 0.47, 0.64, and $0.865\mu\text{m}$) (Figures
155 4(c)-4(e)).

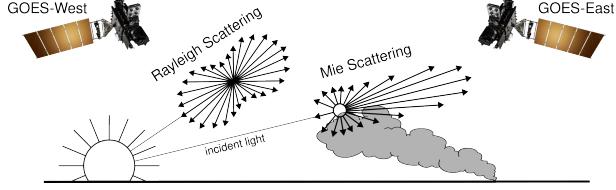


Figure 3: If the particle size is $< \frac{1}{10}$ the λ of the interacting light, then the primary scattering will be Rayleigh. Mie scattering is the predominant scattering mechanism when the particle size is larger than the λ of light. This schematic demonstrates that when the sun is setting in the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

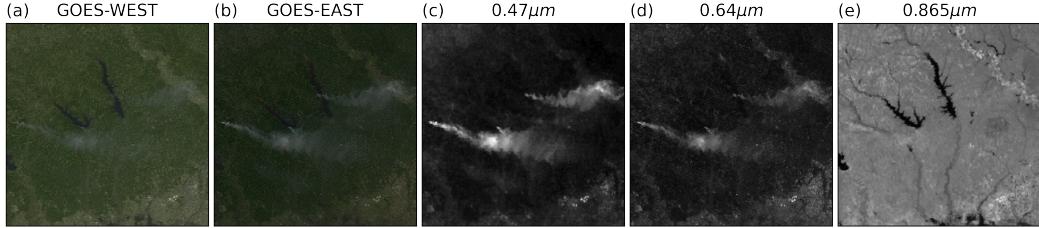


Figure 4: True color (a) GOES-WEST and (b) GOES-EAST imagery from March 23rd, 2022 centered at $(31.1^\circ, -93.8^\circ)$ in Texas, USA taken at 23:20 UTC. The GOES-EAST raw band imagery for (c) blue, (d) red and (e) veggie bands show variations in the SNR for smoke detection in relation to the λ of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

156 3.1.1 From Full Dataset \mathcal{D} to Mie-Derived Dataset \mathcal{D}_M

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

167 3.1.2 PLDR Dataset \mathcal{D}_p

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

174 To build the parent model f_o , that will create the intermediary pseudo-labels (IPLs), we implement
 175 Segmentation Models PyTorch [25] with EfficientNetV2 [26] as the encoder and PSPNet [27]
 176 as the decoder. Input images are $256 \times 256 \times 3$ true-color snapshots; the output is a $256 \times 256 \times 3$
 177 classification map predicting categorical smoke density. We use thermometer encoding (Table 2) and
 178 apply binary cross-entropy loss across density levels. Thermometer encoding is chosen over one-hot
 179 encoding because it captures the ordinal structure of smoke density categories (none < light < medium
 180 < heavy). In thermometer encoding, each higher class includes all lower class activations (e.g., heavy
 181 = [1 1 1]), allowing the model to learn not just class distinctions, but the relative severity of smoke.
 182 We use a confidence threshold of $IoU > 0.01$ [28] to exclude samples with negligible overlap.

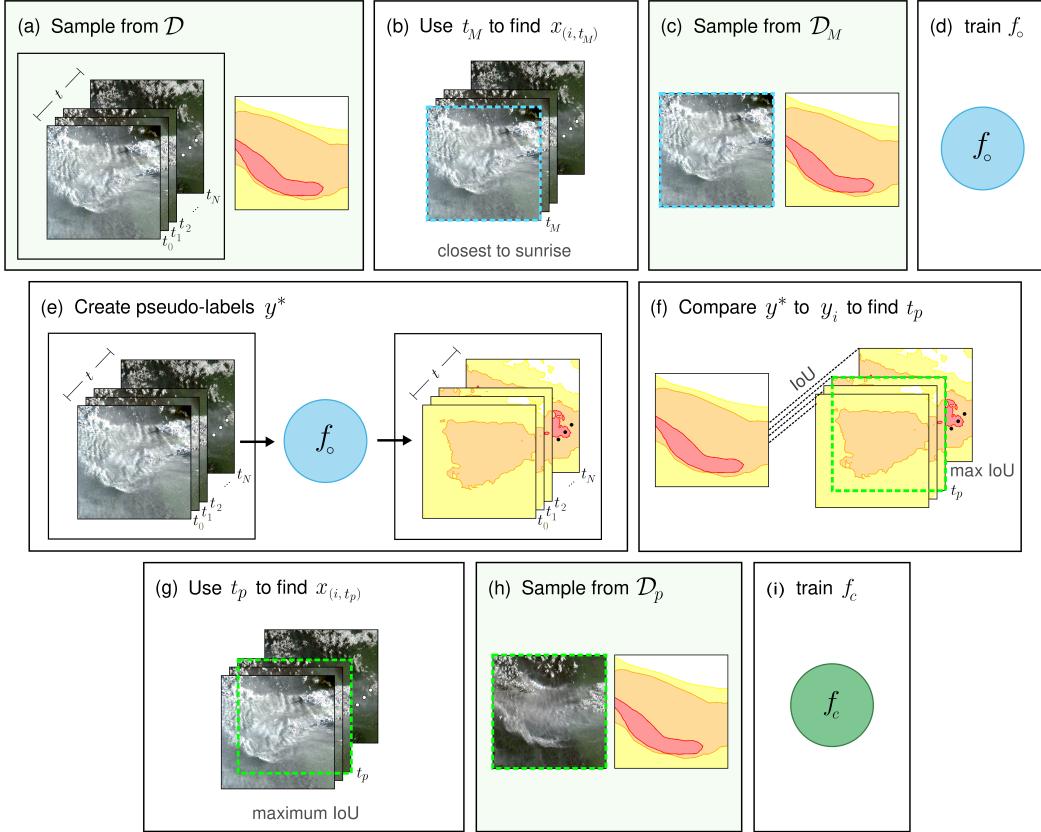


Figure 5: PLDR applied to create the SmokeViz dataset. Green boxes indicate dataset stages. (a) For original dataset \mathcal{D} - analyst annotation y_i corresponds to N satellite images across time window t so that $([x_{(i,t_0)}, \dots, x_{(i,t_N)}], y_i) \in \mathcal{D}$; (b) use Mie scattering to 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. 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	2018-21, 24
validation	20,056	19,223	2023
testing	21,541	20,224	2022

183 Figures 6 and 7 give statistical information on SmokeViz as well as highlight the possible influence
 184 of agricultural burns on the dataset distribution and possible model performance. Figure 6 shows
 185 that sample counts in SmokeViz peak in March and April, corresponding to agricultural burning
 186 rather than wildfire activity. During these months, IoU performance is relatively low in comparison

187 to the scores observed from May through September which align with peak wildfire activity. Figure 7
 188 further supports this trend, showing that the southeastern (SE) quadrant, where agricultural burns are
 189 prevalent, contributes 55% of all samples but exhibits relatively low IoU performance. These patterns
 190 suggest that agricultural burns, which are typically smaller in spatial extent and less visually distinct
 191 than large wildfires, present a greater challenge for accurate detection and segmentation by the model.



Figure 6: Monthly distribution of samples in the full dataset \mathcal{D}_p (left), and monthly IoU scores between f_c predictions and analyst annotations on the \mathcal{D}_p test set (right).

Figure 7: Sample percent contribution by region in the full \mathcal{D}_p dataset and IoU performance of f_c on the \mathcal{D}_p test set across quadrants centered at (40°N, -105°W).

192 3.2 Benchmark Models

193 We benchmark the SmokeViz dataset \mathcal{D}_p using PSPNet [27] and DeepLabV3+ [29] with Efficient-
 194 NetV2 [26], DPT [30] with ViT [31], Segformer [32] and UperNet [33] with EfficientViT [34]. Each
 195 model is trained for 100 epochs while limited to 24 hours using a batch size of 16 and the Adam
 196 optimizer on 8 16GB Nvidia P100 GPUs. These architectures are selected for their relatively low
 197 memory requirements and effectiveness in segmenting multi-scale objects such as smoke plumes.

198 4 Results

199 We evaluate the performance of the parent (f_o) and child (f_c) models using Intersection over Union
 200 (IoU), precision and recall metrics on the test sets of both \mathcal{D}_M and SmokeViz (\mathcal{D}_p), as shown in
 201 Table 4. For each smoke density class, IoU is calculated as the pixel-level intersection between model
 202 predictions and HMS analyst labels, divided by their union, aggregated over all test samples.

Table 4: Segmentation metrics comparing f_o and f_c on the \mathcal{D}_M and \mathcal{D}_p test sets.

Metric	f_o		f_c	
	\mathcal{D}_M	\mathcal{D}_p	\mathcal{D}_M	\mathcal{D}_p
Heavy IoU	0.1510	0.2179	0.1649	0.2604
Medium IoU	0.2572	0.3417	0.2786	0.3965
Light IoU	0.3933	0.5054	0.4395	0.5873
Overall IoU	0.3483	0.4499	0.3898	0.5250
Precision	0.6990	0.7875	0.6942	0.7907
Recall	0.4098	0.5121	0.4706	0.6098

203 As shown in Table 4, in terms of IoU, f_c , that was trained on \mathcal{D}_p , consistently outperform f_o , that
 204 was trained on \mathcal{D}_M , across all smoke density categories. For both f_o and f_c , IoU improves when

205 evaluated on \mathcal{D}_p . The highest overall IoU = 0.5250, is achieved by f_c on \mathcal{D}_p , indicating that PLDR
206 improves image-label alignment and reduces training noise.

207 Precision remains relatively consistent between models, both f_o and f_c achieve precision ≈ 0.69
208 on \mathcal{D}_M , and precision ≈ 0.79 on \mathcal{D}_p . The improvement between datasets but not between models
209 suggest that \mathcal{D}_p 's test set may contain fewer samples where the image contains pixels of true smoke
210 that the label states are not smoke. In contrast, recall improves substantially with PLDR refinement:
211 for f_o , recall increases from recall = 0.4098 on \mathcal{D}_M to 0.5121 on \mathcal{D}_p ; for f_c , it improves from 0.4706
212 to 0.6098. These results demonstrate that training on \mathcal{D}_p rather than \mathcal{D}_M increases the model's ability
213 to detect existing wildfire plumes while maintaining a low false detection rate.

214 Figure 8 illustrates a case in which the PLDR-selected frame better represents the HMS annotation
215 than the Mie-derived selection. Here, the heavy smoke IoU improves from 0.01 to 0.59. While the
216 Mie-derived image is selected based on its proximity to sunrise, PLDR chooses the frame with the
217 highest overlap between the model-generated intermediary pseudo-label and the analyst annotation.
218 This example highlights PLDR's advantage in resolving temporal ambiguity.

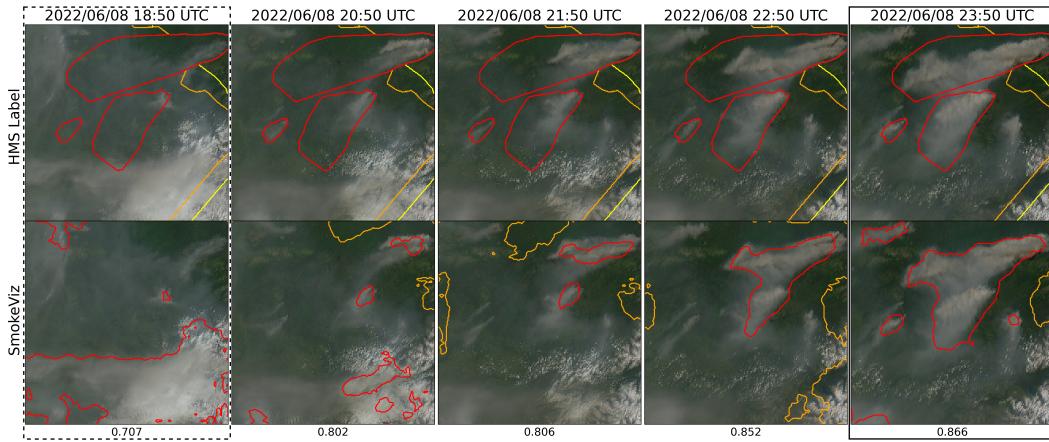


Figure 8: GOES-West imagery from June 8, 2022, over Alaska (61.06°N, 156.12°W). Daylight spanned 12:43-7:53 UTC. The single static HMS annotation (top row) spans 18:50-23:50 UTC is compared with f_o -generated per-frame smoke predictions (bottom row). The leftmost frame (dotted) represents the Mie-derived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

219 To further examine the performance of f_c , we can qualitatively compare its predictions against HMS
220 annotations for samples from \mathcal{D}_p in Figure 9. The model outputs capture more spatially detailed and
221 coherent smoke boundaries compared to the coarser, polygon-based analyst labels.

Table 5: Comparison of segmentation benchmark model IoU metrics on the SmokeViz dataset. Note that the first column is f_c .

encoder decoder	EfficientNet[26] PSPNet[27]	[26] DeepLabV3+[29]	EfficientViT[34] Segformer[32]	[34] UperNet [33]	ViT [31] DPT[30]
Heavy	0.2604	0.2766	0.2351	0.2374	0.2450
Medium	0.3965	0.3854	0.3707	0.3444	0.3745
Light	0.5873	0.5912	0.5259	0.5541	0.5803
Overall	0.5250	0.5246	0.4732	0.4886	0.5136

222 To benchmark performance across segmentation architectures, we evaluate several encoder-decoder
223 models trained on \mathcal{D}_p . Table 5 reports IoU scores by smoke density and overall. While DeepLabV3+
224 achieves the highest IoU for heavy smoke, PSPNet yields the best overall performance. Results across
225 models are relatively consistent, highlighting the robustness of the SmokeViz dataset for training
226 diverse architectures.

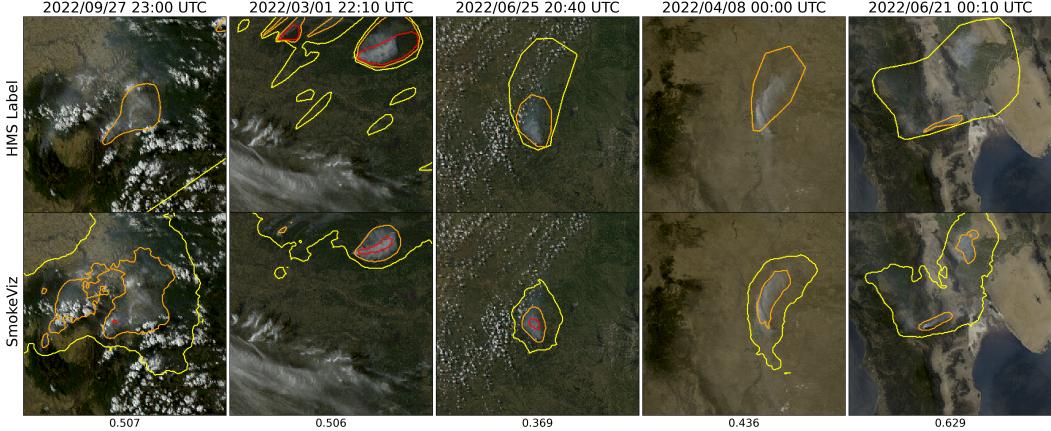


Figure 9: Examples of HMS annotations (top row) vs f_c output (bottom row) on \mathcal{D}_p samples. The overall IoU score is reported at the bottom of each column.

227 5 Limitations

228 More discussion and analysis on the two primary limitations can be found in the Supplementary Materials.
 229 First, pseudo-labeling methods may propagate biases from the parent model into downstream
 230 models. In our case, the increased detectability of forward-scattered light from smoke particulates
 231 may bias the model toward higher performance at larger solar zenith angles. Second, the HMS anno-
 232 tations do not distinguish between fire types and include a large number of controlled agricultural
 233 burns, which may limit the dataset’s applicability for targeting large-scale wildfires.
 234 Several additional limitations remain important directions for future work such as evaluating the
 235 model’s ability to distinguish smoke from dust and investigating uncertainty in the analyst annotations.

236 6 Conclusion

237 In this study, we present **SmokeViz**, a refined satellite imagery dataset for semantic segmentation of
 238 wildfire smoke plumes. Starting from the original NOAA HMS annotations of coarse, many-to-one
 239 approximations of smoke boundaries, we transform the dataset into a one-to-one mapping between
 240 satellite frames and smoke annotations. While the Mie-derived dataset selection process maximized
 241 the potential for detecting smoke if present, it did not account for whether smoke was actually visible
 242 in the selected image, leading to a high incidence of label-image mismatch and associated training
 243 noise. To address this, we introduce **pseudo-label dimension reduction (PLDR)**, a physics-guided,
 244 semi-supervised method that uses a parent model trained on the Mie-derived dataset (\mathcal{D}_M) to generate
 245 pseudo-labels across each annotation’s time window. We then select the image with the highest spatial
 246 overlap between the intermediary pseudo-label and the HMS annotation to construct a refined dataset
 247 (\mathcal{D}_p). A child model trained on \mathcal{D}_p achieves higher segmentation performance than the original parent
 248 model, as measured by IoU on both test sets, demonstrating the value of pseudo-label-based temporal
 249 alignment.

250 SmokeViz serves as a robust and representative dataset for training models to detect wildfire smoke in
 251 GOES imagery at the frame level. In addition to supporting real-time smoke segmentation, this dataset
 252 has potential applications in early wildfire detection, air quality monitoring, and as a smoke analysis
 253 product for data assimilation into dispersion models. It also provides a challenging benchmark for
 254 remote sensing models tasked with segmenting diffuse, low-contrast features like smoke. More
 255 generally, this work illustrates how PLDR can be used to resolve resolution mismatches between data
 256 and labels, especially in settings with time-series or video data paired with coarse annotations. The
 257 dataset is publicly available at <https://noaa-gsl-experimental-pds.s3.amazonaws.com/index.html#SmokeViz/> with code available at <https://github.com/anonymous-smokeviz/SmokeViz>.
 259

260 **References**

- 261 [1] Joseph E Aldy, Maximilian Auffhammer, Maureen Cropper, Arthur Fraas, and Richard Mor-
262 genstern. Looking back at 50 years of the clean air act. *Journal of Economic Literature*, 60(1):
263 179–232, 2022.
- 264 [2] Marshall Burke, Anne Driscoll, Sam Heft-Neal, Jiani Xue, Jennifer Burney, and Michael Wara.
265 The changing risk and burden of wildfire in the united states. *Proceedings of the National
266 Academy of Sciences*, 118(2):e2011048118, 2021.
- 267 [3] Emmanuela Gakidou, Ashkan Afshin, Amanuel Alemu Abajobir, Kalkidan Hassen Abate,
268 Cristiana Abbafati, Kaja M Abbas, Foad Abd-Allah, Abdishakur M Abdulle, Semaw Ferede
269 Abera, Victor Aboyans, et al. Global, regional, and national comparative risk assessment
270 of 84 behavioural, environmental and occupational, and metabolic risks or clusters of risks,
271 1990–2016: a systematic analysis for the global burden of disease study 2016. *The Lancet*, 390
272 (10100):1345–1422, 2017.
- 273 [4] Steven J Goodman, Timothy J Schmit, Jaime Daniels, and Robert J Redmon. *The GOES-R
274 series: a new generation of geostationary environmental satellites*. Elsevier, 2019.
- 275 [5] Tom X.-P. Zhao, Steve Ackerman, and Wei Guo. Dust and smoke detection for multi-channel
276 imagers. *Remote Sensing*, 2(10):2347–2368, 2010. ISSN 2072-4292. doi: 10.3390/rs2102347.
277 URL <https://www.mdpi.com/2072-4292/2/10/2347>.
- 278 [6] T Randriambelo, S Baldy, M Bessafi, Michel Petit, and Marc Despinoy. An improved detection
279 and characterization of active fires and smoke plumes in south-eastern africa and madagascar.
280 *International Journal of Remote Sensing*, 19(14):2623–2638, 1998.
- 281 [7] Eric James, Ravan Ahmadov, and Georg A Grell. Realtime wildfire smoke prediction in the
282 united states: The hrrr-smoke model. In *EGU General Assembly Conference Abstracts*, page
283 19526, 2018.
- 284 [8] Ravan Ahmadov, Haiqin Li, Johana Romero-Alvarez, Jordan Schnell, Sudheer Bhimireddy,
285 Eric James, Ka Yee Wong, Ming Hu, Jacob Carley, Partha Bhattacharjee, et al. Forecasting
286 smoke and dust in noaa’s next-generation high-resolution coupled numerical weather prediction
287 model. Technical report, Copernicus Meetings, 2024.
- 288 [9] Donna McNamara, George Stephens, Mark Ruminski, and Tim Kasheta. The hazard mapping
289 system (hms) - noaa’s multi-sensor fire and smoke detection program using environmental
290 satellites. *Conference on Satellite Meteorology and Oceanography*, 01 2004.
- 291 [10] W Schroeder, M Ruminski, I Csiszar, L Giglio, E Prins, C Schmidt, and J Morisette. Validation
292 analyses of an operational fire monitoring product: The hazard mapping system. *International
293 Journal of Remote Sensing*, 29(20):6059–6066, 2008.
- 294 [11] NOAA. Hazard mapping system fire and smoke product, 2024. URL <https://www.ospo.noaa.gov/Products/land/hms.html#about>.
- 295 [12] Rui Ba, Chen Chen, Jing Yuan, Weiguo Song, and Siuming Lo. Smokenet: Satellite smoke
296 scene detection using convolutional neural network with spatial and channel-wise attention.
297 *Remote Sensing*, 11(14):1702, 2019.
- 298 [13] Favyen Bastani, Piper Wolters, Ritwik Gupta, Joe Ferdinando, and Aniruddha Kembhavi.
299 Satlaspretrain: A large-scale dataset for remote sensing image understanding. In *Proceedings of
300 the IEEE/CVF International Conference on Computer Vision*, pages 16772–16782, 2023.
- 301 [14] Jeff Wen and M. Burke. Wildfire smoke plume segmentation using geostationary sat-
302 elite imagery. *ArXiv*, abs/2109.01637, 2021. URL <https://api.semanticscholar.org/CorpusID:237416777>.
- 303 [15] Alexandra Larsen, Ivan Hanigan, Brian J Reich, Yi Qin, Martin Cope, Geoffrey Morgan, and
304 Ana G Rappold. A deep learning approach to identify smoke plumes in satellite imagery in
305 near-real time for health risk communication. *Journal of exposure science & environmental
epidemiology*, 31(1):170–176, 2021.

- 309 [16] Gencer Sumbul, Marcela Charfuelan, Begüm Demir, and Volker Markl. Bigearthnet: A large-
 310 scale benchmark archive for remote sensing image understanding. In *IGARSS 2019-2019 IEEE*
 311 *international geoscience and remote sensing symposium*, pages 5901–5904. IEEE, 2019.
- 312 [17] J Jakubik, S Roy, C Phillips, P Fraccaro, D Godwin, B Zadrozny, D Szwarcman, C Gomes,
 313 G Nyirjesy, B Edwards, et al. Foundation models for generalist geospatial artificial intelligence.
 314 arxiv 2023. *arXiv preprint arXiv:2310.18660*, 2023.
- 315 [18] Stefano Zorzi, Shabab Bazrafkan, Stefan Habenschuss, and Friedrich Fraundorfer. Polyworld:
 316 Polygonal building extraction with graph neural networks in satellite images. In *Proceedings*
 317 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1848–1857,
 318 2022.
- 319 [19] Asanobu Kitamoto, Jared Hwang, Bastien Vuillod, Lucas Gautier, Yingtao Tian, and Tarin
 320 Clanuwat. Digital typhoon: Long-term satellite image dataset for the spatio-temporal modeling
 321 of tropical cyclones. *Advances in Neural Information Processing Systems*, 36, 2024.
- 322 [20] Bjorn Stevens, Sandrine Bony, Hélène Brogniez, Laureline Hentgen, Cathy Hohenegger,
 323 Christoph Kiemle, Tristan S L’Ecuyer, Ann Kristin Naumann, Hauke Schulz, Pier A Siebesma,
 324 et al. Sugar, gravel, fish and flowers: Mesoscale cloud patterns in the trade winds. *Quarterly*
 325 *Journal of the Royal Meteorological Society*, 146(726):141–152, 2020.
- 326 [21] Dong-Hyun Lee. Pseudo-label : The simple and efficient semi-supervised learning method
 327 for deep neural networks. *ICML 2013 Workshop : Challenges in Representation Learning*
 328 (*WREPL*), 07 2013.
- 329 [22] Martin Raspaud, David Hoese, Adam Dybbroe, Panu Lahtinen, Abhay Devasthale, Mikhail
 330 Itkin, Ulrich Hamann, Lars Ørum Rasmussen, Esben Stigård Nielsen, Thomas Leppelt, et al.
 331 Pytroll: An open-source, community-driven python framework to process earth observation
 332 satellite data. *Bulletin of the American Meteorological Society*, 99(7):1329–1336, 2018.
- 333 [23] MK Bah, MM Gunshor, and TJ Schmit. Generation of goes-16 true color imagery without a
 334 green band. *Earth and Space Science*, 5(9):549–558, 2018.
- 335 [24] Alain Royer, Pierre Vincent, and Ferdinand Bonn. Evaluation and correction of viewing angle
 336 effects on satellite measurements of bidirectional reflectance. *Photogrammetric engineering*
 337 *and remote sensing*, 51(12):1899–1914, 1985.
- 338 [25] Pavel Iakubovskii. Segmentation models pytorch. https://github.com/qubvel/segmentation_models.pytorch, 2019.
- 340 [26] Mingxing Tan and Quoc Le. Efficientnetv2: Smaller models and faster training. In *International*
 341 *conference on machine learning*, pages 10096–10106. PMLR, 2021.
- 342 [27] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene
 343 parsing network. In *Proceedings of the IEEE conference on computer vision and pattern*
 344 *recognition*, pages 2881–2890, 2017.
- 345 [28] Rafael EP Ferreira, Yong Jae Lee, and João RR Dórea. Using pseudo-labeling to improve
 346 performance of deep neural networks for animal identification. *Scientific Reports*, 13(1):13875,
 347 2023.
- 348 [29] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille.
 349 Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and
 350 fully connected crfs. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):
 351 834–848, 2017.
- 352 [30] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction.
 353 In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 12179–
 354 12188, 2021.
- 355 [31] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,
 356 Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al.
 357 An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint*
 358 *arXiv:2010.11929*, 2020.

- 359 [32] Enze Xie, Wenhui Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo.
360 Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances*
361 *in neural information processing systems*, 34:12077–12090, 2021.
- 362 [33] Tete Xiao, Yingcheng Liu, Bolei Zhou, Yuning Jiang, and Jian Sun. Unified perceptual parsing
363 for scene understanding. In *Proceedings of the European conference on computer vision*
364 (*ECCV*), pages 418–434, 2018.
- 365 [34] Xinyu Liu, Houwen Peng, Ningxin Zheng, Yuqing Yang, Han Hu, and Yixuan Yuan. Efficientvit:
366 Memory efficient vision transformer with cascaded group attention. In *Proceedings of the*
367 *IEEE/CVF conference on computer vision and pattern recognition*, pages 14420–14430, 2023.

368 **NeurIPS Paper Checklist**

369 **1. Claims**

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

372 Answer: [Yes]

373 Justification: The claims of using intermediary pseudo-labels to create a more robust dataset
374 is reflected in the paper's contributions.

375 Guidelines:

- 376 • The answer NA means that the abstract and introduction do not include the claims
377 made in the paper.
- 378 • The abstract and/or introduction should clearly state the claims made, including the
379 contributions made in the paper and important assumptions and limitations. A No or
380 NA answer to this question will not be perceived well by the reviewers.
- 381 • The claims made should match theoretical and experimental results, and reflect how
382 much the results can be expected to generalize to other settings.
- 383 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
384 are not attained by the paper.

385 **2. Limitations**

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

387 Answer: [Yes]

388 Justification: We address limitations of the dataset.

389 Guidelines:

- 390 • The answer NA means that the paper has no limitation while the answer No means that
391 the paper has limitations, but those are not discussed in the paper.
- 392 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 393 • The paper should point out any strong assumptions and how robust the results are to
394 violations of these assumptions (e.g., independence assumptions, noiseless settings,
395 model well-specification, asymptotic approximations only holding locally). The authors
396 should reflect on how these assumptions might be violated in practice and what the
397 implications would be.
- 398 • The authors should reflect on the scope of the claims made, e.g., if the approach was
399 only tested on a few datasets or with a few runs. In general, empirical results often
400 depend on implicit assumptions, which should be articulated.
- 401 • The authors should reflect on the factors that influence the performance of the approach.
402 For example, a facial recognition algorithm may perform poorly when image resolution
403 is low or images are taken in low lighting. Or a speech-to-text system might not be
404 used reliably to provide closed captions for online lectures because it fails to handle
405 technical jargon.
- 406 • The authors should discuss the computational efficiency of the proposed algorithms
407 and how they scale with dataset size.
- 408 • If applicable, the authors should discuss possible limitations of their approach to
409 address problems of privacy and fairness.
- 410 • While the authors might fear that complete honesty about limitations might be used by
411 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
412 limitations that aren't acknowledged in the paper. The authors should use their best
413 judgment and recognize that individual actions in favor of transparency play an impor-
414 tant role in developing norms that preserve the integrity of the community. Reviewers
415 will be specifically instructed to not penalize honesty concerning limitations.

416 **3. Theory assumptions and proofs**

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

419 Answer: [NA]

420 Justification: No theoretical results are presented.

421 Guidelines:

- 422 • The answer NA means that the paper does not include theoretical results.
- 423 • All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- 425 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 426 • The proofs can either appear in the main paper or the supplemental material, but if
- 427 they appear in the supplemental material, the authors are encouraged to provide a short
- 428 proof sketch to provide intuition.
- 429 • Inversely, any informal proof provided in the core of the paper should be complemented
- 430 by formal proofs provided in appendix or supplemental material.
- 431 • Theorems and Lemmas that the proof relies upon should be properly referenced.

432 **4. Experimental result reproducibility**

433 Question: Does the paper fully disclose all the information needed to reproduce the main ex-

434 perimental results of the paper to the extent that it affects the main claims and/or conclusions

435 of the paper (regardless of whether the code and data are provided or not)?

436 Answer: [Yes]

437 Justification: We provide the code to create the datasets along with the final dataset hosted

438 on AWS by NOAA.

439 Guidelines:

- 440 • The answer NA means that the paper does not include experiments.
- 441 • If the paper includes experiments, a No answer to this question will not be perceived
- 442 well by the reviewers: Making the paper reproducible is important, regardless of
- 443 whether the code and data are provided or not.
- 444 • If the contribution is a dataset and/or model, the authors should describe the steps taken
- 445 to make their results reproducible or verifiable.
- 446 • Depending on the contribution, reproducibility can be accomplished in various ways.
- 447 For example, if the contribution is a novel architecture, describing the architecture fully
- 448 might suffice, or if the contribution is a specific model and empirical evaluation, it may
- 449 be necessary to either make it possible for others to replicate the model with the same
- 450 dataset, or provide access to the model. In general, releasing code and data is often
- 451 one good way to accomplish this, but reproducibility can also be provided via detailed
- 452 instructions for how to replicate the results, access to a hosted model (e.g., in the case
- 453 of a large language model), releasing of a model checkpoint, or other means that are
- 454 appropriate to the research performed.
- 455 • While NeurIPS does not require releasing code, the conference does require all submis-
- 456 sions to provide some reasonable avenue for reproducibility, which may depend on the
- 457 nature of the contribution. For example
 - 458 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
 - 459 to reproduce that algorithm.
 - 460 (b) If the contribution is primarily a new model architecture, the paper should describe
 - 461 the architecture clearly and fully.
 - 462 (c) If the contribution is a new model (e.g., a large language model), then there should
 - 463 either be a way to access this model for reproducing the results or a way to reproduce
 - 464 the model (e.g., with an open-source dataset or instructions for how to construct
 - 465 the dataset).
 - 466 (d) We recognize that reproducibility may be tricky in some cases, in which case
 - 467 authors are welcome to describe the particular way they provide for reproducibility.
 - 468 In the case of closed-source models, it may be that access to the model is limited in
 - 469 some way (e.g., to registered users), but it should be possible for other researchers
 - 470 to have some path to reproducing or verifying the results.

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

472 Question: Does the paper provide open access to the data and code, with sufficient instruc-

473 tions to faithfully reproduce the main experimental results, as described in supplemental

474 material?

475 Answer: [Yes]

476 Justification: The SmokeViz dataset is released along with code used to develop it.

477 Guidelines:

- 478 • The answer NA means that paper does not include experiments requiring code.
- 479 • Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 480 • 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).
- 481 • 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.
- 482 • 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.
- 483 • 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.
- 484 • At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- 485 • Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

497 **6. Experimental setting/details**

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

501 Answer: [Yes]

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

503 Guidelines:

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

509 **7. Experiment statistical significance**

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

512 Answer: [No]

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

515 Guidelines:

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

- 526 • It should be clear whether the error bar is the standard deviation or the standard error
 527 of the mean.
 528 • It is OK to report 1-sigma error bars, but one should state it. The authors should
 529 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
 530 of Normality of errors is not verified.
 531 • For asymmetric distributions, the authors should be careful not to show in tables or
 532 figures symmetric error bars that would yield results that are out of range (e.g. negative
 533 error rates).
 534 • If error bars are reported in tables or plots, The authors should explain in the text how
 535 they were calculated and reference the corresponding figures or tables in the text.

536 **8. Experiments compute resources**

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

540 Answer: [Yes]

541 Justification: We mention 8 16GB P100 GPUs and limit to 24 hours of run time.

542 Guidelines:

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

551 **9. Code of ethics**

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

554 Answer: [Yes]

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

556 Guidelines:

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

562 **10. Broader impacts**

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

565 Answer: [Yes]

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

568 Guidelines:

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

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

591 11. Safeguards

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

595 Answer: [NA]

596 Justification: There are no risks for misuse.

597 Guidelines:

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

608 12. Licenses for existing assets

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

612 Answer: [Yes]

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

615 Guidelines:

- 616 • The answer NA means that the paper does not use existing assets.
- 617 • The authors should cite the original paper that produced the code package or dataset.
- 618 • The authors should state which version of the asset is used and, if possible, include a
 619 URL.
- 620 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 621 • For scraped data from a particular source (e.g., website), the copyright and terms of
 622 service of that source should be provided.
- 623 • If assets are released, the license, copyright information, and terms of use in the
 624 package should be provided. For popular datasets, paperswithcode.com/datasets
 625 has curated licenses for some datasets. Their licensing guide can help determine the
 626 license of a dataset.
- 627 • For existing datasets that are re-packaged, both the original license and the license of
 628 the derived asset (if it has changed) should be provided.

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

631 **13. New assets**

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

634 Answer: [Yes]

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

637 Guidelines:

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

646 **14. Crowdsourcing and research with human subjects**

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

650 Answer: [NA]

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

652 Guidelines:

- 653 • The answer NA means that the paper does not involve crowdsourcing nor research with
654 human subjects.
655 • Including this information in the supplemental material is fine, but if the main contribu-
656 tion of the paper involves human subjects, then as much detail as possible should be
657 included in the main paper.
658 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
659 or other labor should be paid at least the minimum wage in the country of the data
660 collector.

661 **15. Institutional review board (IRB) approvals or equivalent for research with human
662 subjects**

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

667 Answer: [NA]

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

669 Guidelines:

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

680 **16. Declaration of LLM usage**

681 Question: Does the paper describe the usage of LLMs if it is an important, original, or
682 non-standard component of the core methods in this research? Note that if the LLM is used
683 only for writing, editing, or formatting purposes and does not impact the core methodology,
684 scientific rigorousness, or originality of the research, declaration is not required.

685 Answer: [NA]

686 Justification: There are no LLM components to this work.

687 Guidelines:

- 688 • The answer NA means that the core method development in this research does not
689 involve LLMs as any important, original, or non-standard components.
690 • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)
691 for what should or should not be described.