

---

# *Improving Deep Learning-Based Wildfire Smoke Plume Detection with a Multi-Model Ensemble Approach*

---

Anonymous Author(s)

Affiliation

Address

email

## **Abstract**

1 With increasing frequency and severity of wildfires, there is an urgent need for  
2 wildfire and smoke detection tools that can effectively and rapidly monitor smoke  
3 at a large scale. Recent advancements in computer vision have demonstrated the  
4 potential of deep learning models to automatically separate high-resolution images  
5 into labeled regions for high-accuracy feature detection. However, single-model  
6 approaches can struggle with generalization and accuracy in diverse conditions,  
7 which is necessary in the operational use of a smoke detection tool. To address  
8 these challenges, we propose using an ensemble of deep learning models to produce  
9 more accurate annotations of wildfire smoke plumes and their relative density (light,  
10 medium, heavy) in satellite imagery. Our preliminary results indicate that ensemble  
11 techniques can improve performance compared to using a single model, and that  
12 further investigation is needed to optimize the ensemble. This approach aims to  
13 provide a more reliable satellite-based tool for real-time monitoring of smoke.  
14 This will have numerous downstream impacts, such as aiding fire and hazard  
15 management efforts, improving modeling of wildfire behavior and air quality, and  
16 ultimately contributing to climate resilience and adaptation strategies.

## **1 Introduction**

18 In the last four decades, wildfire activity has increased drastically in the U.S. In fact, the conditions  
19 leading to wildfires have been shown to occur more frequently as a direct result of climate change [1].  
20 Wildfires produce particulate emissions, such as smoke and ash, which can be tracked to quantify the  
21 impact a fire has beyond its burn area. Satellite observations reveal that the number of days with smoke  
22 in the air have substantially increased in the U.S. during the last two decades [1]. Furthermore, the  
23 human impacts of smoke exposure include increased morbidity and mortality as well as downstream  
24 economic costs [2].

25 Determining causal links between wildfire activity and health impacts in regions distant from the  
26 source fire requires accurate large-scale monitoring of wildfire smoke and its movement. An intuitive  
27 approach is to utilize satellite-based methods that can monitor evolution of smoke plumes over large  
28 areas. However, such methods have yet to provide precise and high-frequency information on smoke  
29 density, or are confined to small case-study regions [3–5].

## 30 **2 Background**

31 The National Oceanic and Atmospheric Administration (NOAA) Geostationary Operational Environ-  
32 mental Satellites (GOES) provide high spatial and temporal resolution imagery of North America  
33 [6]. The NOAA Hazard Mapping System (HMS) Fire and Smoke Product currently relies on expert  
34 human analysts to annotate the presence of smoke over North America using GOES imagery [7].  
35 However, this product is limited by the availability of human analysts and their time, with annotations  
36 outputted only once to a few times a day and usually having a delay between smoke occurrence and  
37 the annotation. Thus, emergency responders may be delayed without real-time smoke conditions, or  
38 could miss early detections of fire considering that smoke can obfuscate fire visibility. To address  
39 these limitations, we leverage advancements in DL to automate the detection of smoke from GOES  
40 imagery, using the existing HMS dataset for training. By automating this task, we can enable more  
41 frequent detection of smoke plumes, which will aid in active wildfire monitoring and mitigating air  
42 quality impacts.

43 This proposal focuses on enhancing the capability of DL models to detect smoke using multi-model  
44 ensemble methods. Ensembles, which combine the predictions of multiple models, have been  
45 shown to often perform better than a single model in classification tasks [8]. Particularly, utilizing  
46 a diverse set of classifiers in an ensemble is important to achieve the improvement in performance  
47 [9]. Furthermore, combining the predictions of multiple independently-trained neural networks can  
48 improve generalization and detection accuracy [10–12]. To our knowledge, these ensemble methods  
49 have not yet been used for wildfire smoke detection.

## 50 **3 Proposed methods**

51 We use the SmokeViz dataset (Section 7.1) which consists of 183,672 smoke plume samples spanning  
52 2018–2024, each with three spectral channels of GOES imagery paired with human analyst annotations  
53 of light, medium, or heavy smoke. We use 2018–2021 and 2024 for training, 2023 for validation,  
54 and 2022 for the test set, ensuring the testing and validation data years are independent of the  
55 training data. During validation and testing, we quantify model performance with Intersection  
56 over Union (IoU) score, which measures pixel-wise alignment between the model prediction and  
57 the ground truth. The IoU metric supports pathways to climate impact since improved IoU score  
58 directly relates to more smoke being accurately detected (reducing false detections and increasing true  
59 detections) which equips hazard management to appropriately respond to current wildfire and smoke  
60 conditions. We utilize a variety of encoder-decoder architectures designed for semantic segmentation  
61 that include different features such as multi-scale fields-of-view and precise boundary detection  
62 [13–15]. Additionally, we selected the best-performing single architecture and trained it with 12  
63 seeds to generate different initial random weights. By varying the initialization, the model’s parameter  
64 space can be searched more fully; ideally, different minima of the loss function will be found. All the  
65 models were trained independently for 24 hours on 8 Nvidia P100 GPUs using the Adam optimizer  
66 and the binary cross entropy loss function.

67 The ensemble method we are using in this preliminary analysis is an unweighted average of  $N$   
68 model outputs (Figure 1), as used in [11]. We experimented with the ensemble size in supplementary  
69 section 7.2, and present results here with  $N = 8$ . We create an architecture-based ensemble and a  
70 initialization-based ensemble.

## 71 **4 Preliminary results**

72 Table 1 shows the IoU scores on the test set for individual models and ensembles. The ensemble of 8  
73 different architectures outperforms the individual models, with an improvement in all IoU metrics.  
74 The ensemble of 8 different initial weights (but the same architecture, PAN) also outperforms the  
75 individual models, with a similar improvement in the IoU scores. Future work is necessary to reveal  
76 the mechanisms behind the ensemble’s improvement in performance, as well as how to optimally  
77 select models that are included in the ensemble.

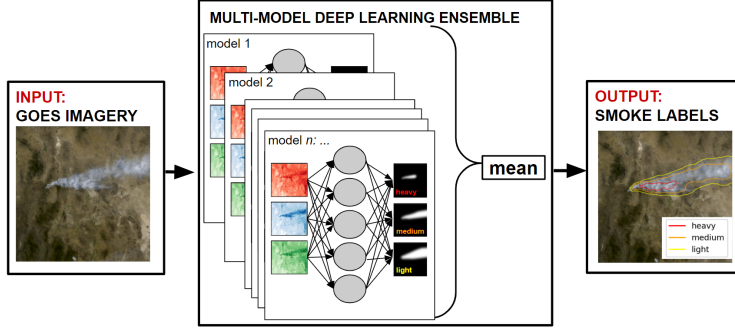


Figure 1: Multi-model ensemble framework. GOES imagery is inputted to  $N$  independently-trained models whose output is combined with an un-weighted average to produce the ensemble prediction of pixel-wise smoke labels.

Figure 2 shows an example of smoke plume detection from the testing dataset. The ensemble predictions have smoother boundaries than the individual model outputs, making the prediction more comparable to the human analyst-drawn annotations.

Table 1: Test IoU results across heavy, medium and light smoke density and over all densities with single models and ensemble schemes.

	Heavy	Medium	Light	Overall
Single Model: DLV3P [13]	0.347	0.441	0.666	0.599
Single Model: PAN [14]	0.349	0.478	0.664	0.604
Architecture Ensemble	0.400	0.507	0.692	0.635
Random Initial Weights Ensemble	0.409	0.512	0.684	0.631

80

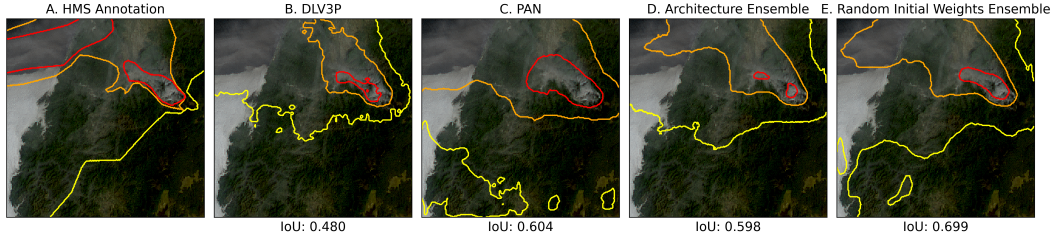


Figure 2: Example of smoke plume detection at  $(43.37^\circ, -123.25^\circ)$  on 2022/10/15 15:50 UTC. Red, orange, and yellow contours represent heavy, medium and light density smoke annotation/prediction, respectively. (A) displays the ground truth annotation; (B-C) show the predictions of two individual models; (D) shows the prediction of the architecture-based ensemble; (E) shows the prediction of an ensemble made with models initialized with different random weights.

## 5 Limitations and future work

This proposal explores two schemes for building ensembles of DL models that both improve on testing set IoU and smooth annotation boundaries. However, further investigation is required to give insight on how the ensemble reduces error and improves generalizability and what the optimal ensemble size and type are. One area to explore is "stacking", where an optimized meta-model is used to combine multi-model outputs, as used in [16–18]. Furthermore, future work will utilize the multi-model ensemble to quantify uncertainty in smoke annotations, enabling users like wildfire response teams and environmental agencies to assess the reliability of detections in real time.

## 6 Pathways to climate impact

We propose ensemble methods for improving DL-based smoke detection. Preliminary that can be used to aid fire and hazard management by automating the monitoring of smoke in real-time from

92 satellite imagery with smooth and accurate smoke annotations. This will improve prediction of  
93 wildfire movement and air quality impacts, ultimately supporting climate resilience and adaptation  
94 strategies.

## References

- [1] Marshall Burke, Anne Driscoll, Sam Heft-Neal, Jiani Xue, Jennifer Burney, and Michael Wara. The changing risk and burden of wildfire in the united states. *Proceedings of the National Academy of Sciences*, 118(2):e2011048118, 2021.
- [2] Wayne E. Cascio. Wildland fire smoke and human health. *The Science of the Total Environment*, 624:586–595, 05 2018. Epub 2017 Dec 27.
- [3] Jeff Wen and M. Burke. Wildfire smoke plume segmentation using geostationary satellite imagery. *ArXiv*, abs/2109.01637, 2021.
- [4] Jiayun Yao, Sean M. Raffuse, Michael Brauer, Grant J. Williamson, David M.J.S. Bowman, Fay H. Johnston, and Sarah B. Henderson. Predicting the minimum height of forest fire smoke within the atmosphere using machine learning and data from the calipso satellite. *Remote Sensing of Environment*, 206:98–106, 2018.
- [5] Alexandra Larsen, Ivan Hanigan, Brian J. Reich, Yi Qin, Martin Cope, Geoffrey Morgan, and Ana G. Rappold. A deep learning approach to identify smoke plumes in satellite imagery in near-real time for health risk communication. *Journal of Exposure Science & Environmental Epidemiology*, 31(1):170–176, 2021.
- [6] S. J. Goodman, T. J. Schmit, J. Daniels, and R. J. Redmon. *The GOES-R Series: A New Generation of Geostationary Environmental Satellites*. Elsevier, 2019.
- [7] Donna McNamara, George Stephens, Mark Ruminski, and Tim Kasheta. The hazard mapping system (hms) - noaa’s multi-sensor fire and smoke detection program using environmental satellites. *Conference on Satellite Meteorology and Oceanography*, 01 2004.
- [8] Thomas G. Dietterich. Ensemble methods in machine learning. *Multiple Classifier Systems*, pages 1–15, 2000.
- [9] Ludmila I. Kuncheva and Christopher J. Whitaker. Measures of diversity in classifier ensembles and their relationship with the ensemble accuracy. *Machine Learning*, 51(2):181–207, 2003.
- [10] L.K. Hansen and P. Salamon. Neural network ensembles. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 12(10):993–1001, 1990.
- [11] Aurélien Bibaut Cheng Ju and Mark van der Laan. The relative performance of ensemble methods with deep convolutional neural networks for image classification. *Journal of Applied Statistics*, 45(15):2800–2818, 2018. PMID: 31631918.
- [12] Giorgio Giacinto and Fabio Roli. Design of effective neural network ensembles for image classification purposes. *Image and Vision Computing*, 19(9):699–707, 2001.
- [13] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation, 2018.
- [14] Hanchao Li, Pengfei Xiong, Jie An, and Lingxue Wang. Pyramid attention network for semantic segmentation. *CoRR*, abs/1805.10180, 2018.
- [15] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++: A nested u-net architecture for medical image segmentation. *CoRR*, abs/1807.10165, 2018.
- [16] Gonçalo Falcão, Armando M. Fernandes, Nuno Garcia, Helena Aidos, and Pedro Tomás. Stacking deep learning models for early detection of wildfire smoke plumes. In *2023 31st European Signal Processing Conference (EUSIPCO)*, pages 1370–1374, 2023.

- [17] Binxu Zhai and Jianguo Chen. Development of a stacked ensemble model for forecasting and analyzing daily average pm2.5 concentrations in beijing, china. *Science of The Total Environment*, 635:644–658, 2018.
- [18] Jiayue Gu, Shuguang Liu, Zhengzheng Zhou, Sergey R. Chalov, and Qi Zhuang. A stacking ensemble learning model for monthly rainfall prediction in the taihu basin, china. *Water*, 14(3), 2022.

## 7 Supplementary Material

### 7.1 Data and Code Availability

The code for this work is available at <https://github.com/anonymous-ensemble-smoke/ensemble-AI-smoke-detection/tree/main>. The dataset can be accessed at <https://noaa-gsl-experimental-pds.s3.amazonaws.com/index.html#SmokeViz/>.

### 7.2 Ensemble Size Analysis

Figure 3 shows the IoU performance over all smoke densities as a function of ensemble size,  $N$ , for the two ensemble schemes. The ensemble with different initial weights generally improves as models are added to the ensemble. The ensemble of different architectures improves with more models up to 8 models, but then decreased in IoU with more models added to the ensemble. This decrease in performance could be due to the additional architectures not having enough variation in model bias to improve ensemble performance. Future work will aim to clarify exactly how different ensemble sizes behave and reduce error.

An additional example from the test data set is shown in Figure 4, where the individual model output has jagged boundaries and the ensemble outputs smooth over these edges. We see a peak in performance at  $N = 8$  in this sample where the  $N = 8$  ensemble has the highest IoU score, and the smoothing does not seem to improve in the  $N = 12$  ensemble output. This sample supports the proposed idea that ensemble DL can smooth over rough edges in semantic segmentation, and warrants further investigation for the optimal ensemble and how to use the multi-model approach to quantify uncertainty.

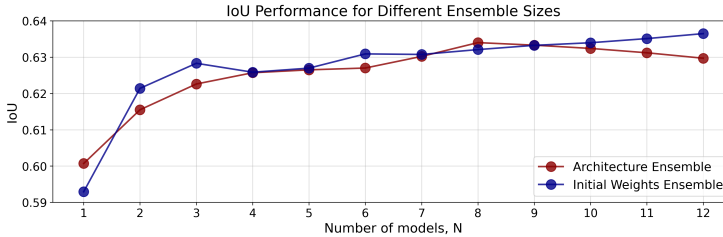


Figure 3: Overall IoU as a function of  $N$  for two ensemble design schemes: random initial weights (blue) and architecture-based (red).

### 7.3 Intersection Over Union (IoU) formula

Equation 1 provides a mathematical formula for IoU, where  $y_i$  represents the ground truth and  $y_i^*$  represents the model's prediction.

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

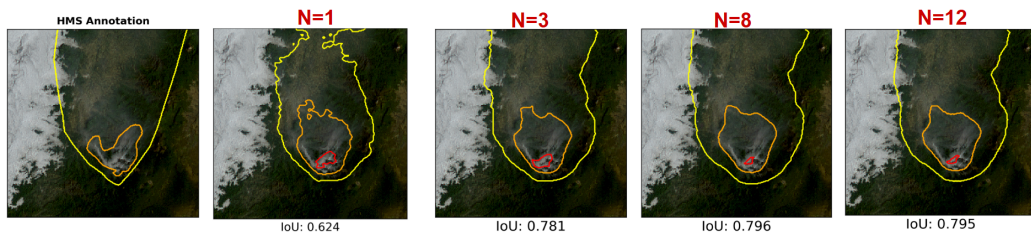


Figure 4: Example of smoke plume detection at (44.24, -122.74) on 2022/09/27 15:30 UTC. Red contours outline the heavy density smoke, orange contours outline the medium density smoke, and yellow contours outline the light density smoke annotations. The first panel displays the ground truth HMS annotation; the second panel is the individual model output of DLV3P; the following panels the prediction of an architecture-based ensemble as it increases in size,  $N$ .