



# Scrapping The Web For Early Wildfire Detection: A New Annotated Dataset of Images and Videos of Smoke Plumes In-the-wild

Mateo Lostanlen<sup>\*1</sup> Nicolás Isla<sup>\*2,3</sup> Jose Guillen<sup>2</sup>  
Felix Veith<sup>1</sup> Cristian Buc<sup>3</sup> Valentin Barriere<sup>2,3</sup>  
<sup>1</sup>PyroNear <sup>2</sup>Universidad de Chile <sup>3</sup>Cenia <sup>\*</sup>Shared First Authorship



UNIVERSIDAD DE CHILE

## Introduction

Wildfires are a growing global challenge, intensified by climate change and human activities, leading to increased frequency and intensity of these disasters. Early wildfire detection (EWD) is critical for enabling rapid response and minimizing environmental, economic, and societal damage.

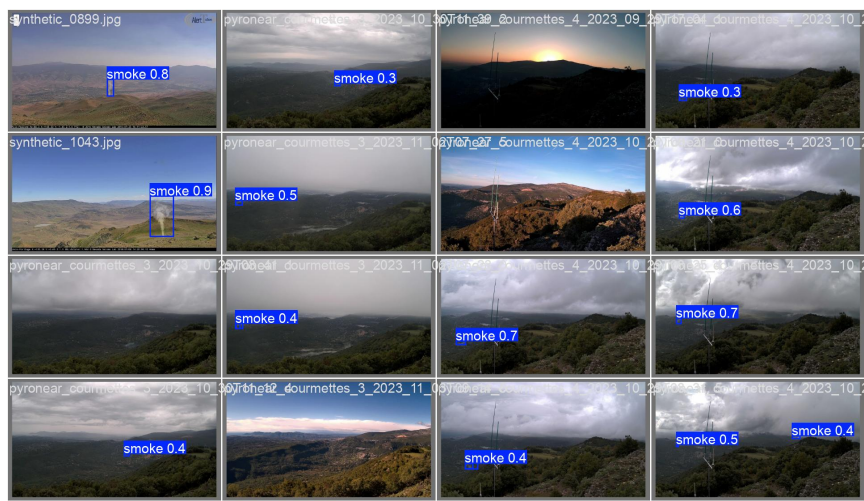


Figure 1. Examples of our dataset, containing real images and videos from France, Spain and United States, and synthetic images.

The **PyroNear<sub>2024</sub>** dataset was created to enhance the performance of smoke plume detection models, incorporating diverse and challenging data from real-world scenarios.

## Objectives

The following objectives outline the main goals of this work, aiming to address critical challenges and contribute to advancements in the field of early detection.

- To collect a diverse and representative dataset of videos for EWD using web scrapping and smoke plume detection models.
- To benchmark the image part of our dataset against other publicly available EWD image datasets, in a cross-dataset setting.
- To evaluate the image component of our dataset within a cross-dataset framework in comparison with other publicly available EWD image datasets.
- To take advantage of the video-based nature of our data in order to train temporal models improving early detection.

## Methodology

- **Data Collection:**
  - The **PyroNear<sub>2024</sub>** dataset includes 50,000 images from the AlertWildfire network, Google searches, and 200 synthetic images generated with Blender.
  - Collaborative annotation refined 120,000 initial labels into 24,000 high-quality labeled images.
- **Image Data Creation:**
  - The dataset targets early fire detection by subsampling initial wildfire frames.
- **Video Data Creation:**
  - Temporal sequences created by capturing images 15 minutes before and after detected smoke events.
- **Model Training and Evaluation:**
  - YOLOv8-small model trained for image-based detection with hyperparameter tuning on validation data and tested on unseen data to assess F1 score, precision, and recall.
  - For video detection, YOLOv8 extracted frames processed by a CNN-LSTM to capture temporal patterns for early fire detection.

## Summary of Datasets.

Dataset	# Wildfires	Total Images	Avg. BBox Area (%)
AiForMankind	31	2935/2584	1.341
Fuego	38	1661/1572	0.175
Nemo	62	2859/2570	7.337
SmokeFrames-2.4k	75	2410/976	14.167
SmokeFrames-50k	643	54576/36304	14.384
PyroNear <sub>2024</sub> – I	532/400	3292/2934	1.043
PyroNear <sub>2024</sub> – V	532/400	35406/19336	0.590

Table 1. Summary of Datasets. Columns marked with  $\bar{\square}$  indicate Total/Wildfire images.

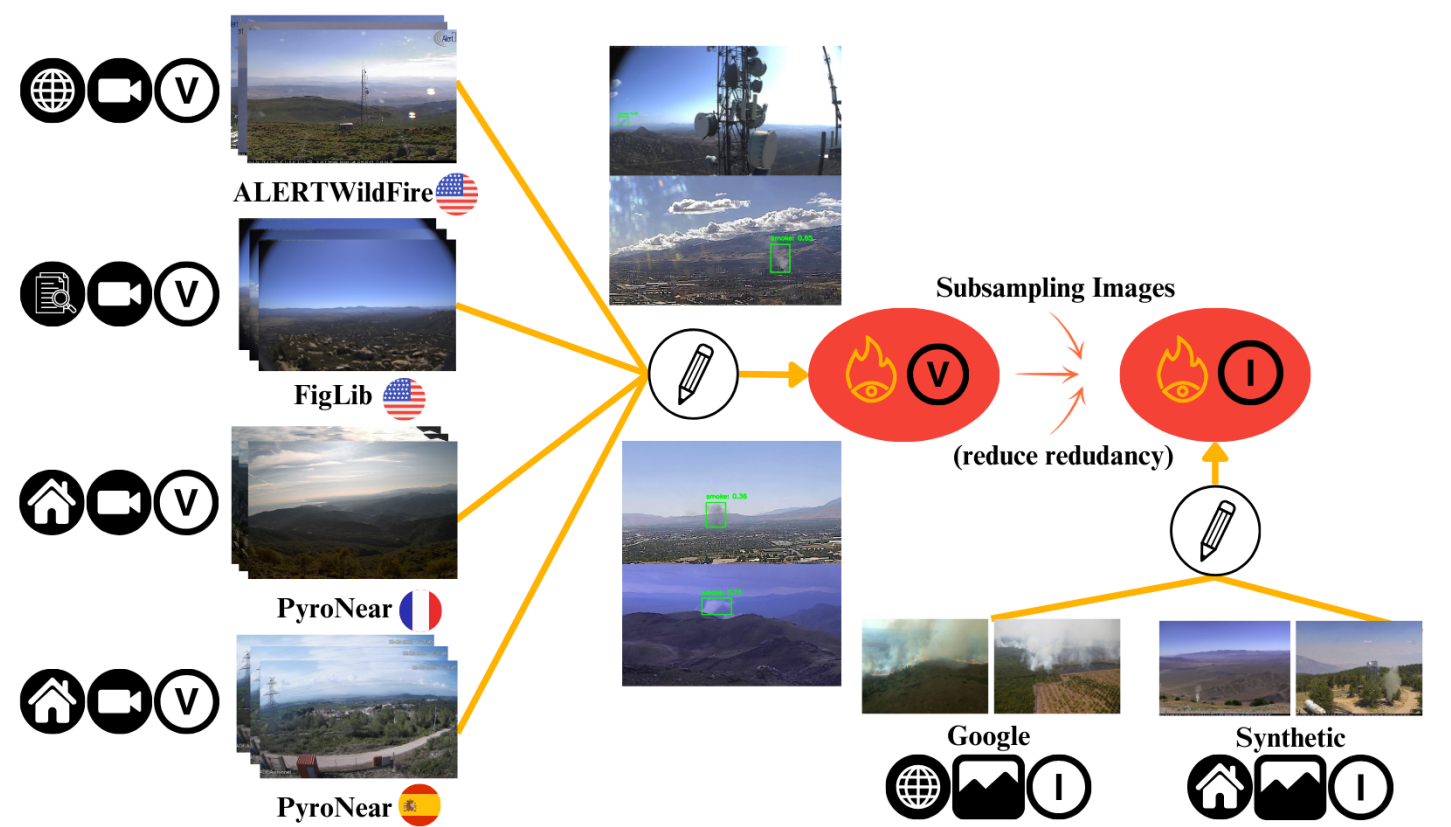


Figure 2. Summary of the whole process to create the video and the image dataset.

## Results Images

Train Dataset	Overall Precision	Overall Recall	Overall F1 Score
PyroNear-I	0.632	0.649	0.639
Nemo	0.731	0.586	0.602
AiForMankind	0.644	0.512	0.549
SmokeFrames-2.4k	0.892	0.417	0.491
SmokeFrames-50k	0.517	0.438	0.447
Fuego	0.506	0.449	0.462

Table 2. Summary of cross-dataset training results. The "Overall" row represents the average of test results across all combinations for each training dataset.

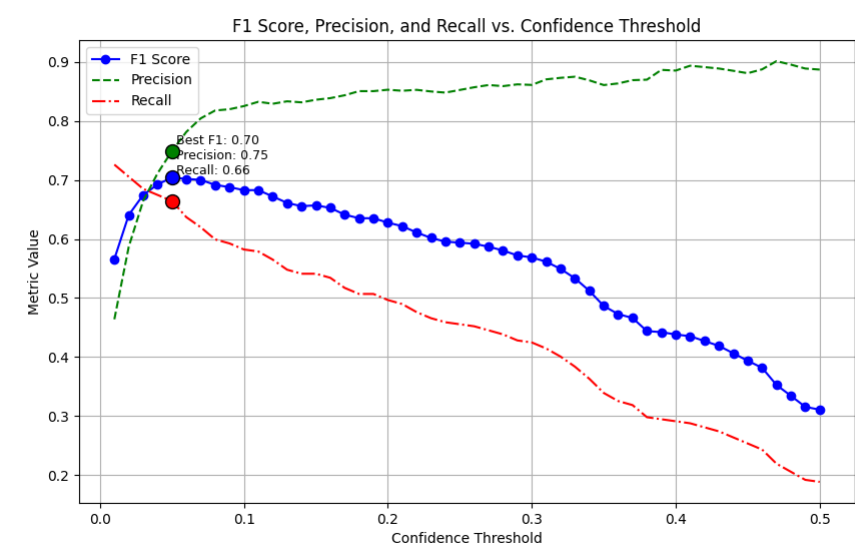


Figure 3. F1 on PyroNear<sub>2024</sub> – I with the selection of the best detection threshold  $\tau_d$

## Results Videos

Model	Precision	Recall	F1 Score	Time Elapsed
YOLOv8 (one frame)	0.920	0.529	0.672	1'46"
YOLOv8+CNN-LSTM	0.912	0.578	0.708	1'05"

Table 3. Performance comparison of image-based and sequential models on the PyroNear<sub>2024</sub> – V dataset using Precision, Recall and F1, as well as the time elapsed before detecting the fire.

## Conclusions

- **Introduction of PyroNear<sub>2024</sub> Dataset:**
  - A diverse dataset for early wildfire smoke plume detection.
  - Includes data from three countries, covering 532 wildfires.
- **Key Achievements:**
  - **Improved Detection:**
    - PyroNear<sub>2024</sub> – I combined with other datasets enhances performance and stability for smoke detection.
    - Small bounding box annotations tailored for real-world early detection scenarios.
  - **Sequential Data:**
    - Videos incorporated to train temporal models.
    - Improved recall.
    - Reduced detection times to 1:05 minutes.

