

---

# Datasheet for FIRE-D: NASA-centric Remote Sensing of Wildfire

---

**Yuzhou Chen**  
Temple University

**Nicholas LaHaye**  
NASA's JPL at Caltech

**Jae Won Choi**  
UT Dallas

**Zhiwei Zhen**  
UT Dallas

**Philip E. Davis**  
University of Utah

**Huikyo Lee**  
NASA's JPL at Caltech

**Manish Parashar**  
University of Utah

**Yulia R. Gel**  
UT Dallas

## 1 A Datasheet

2 This Datasheets for Dataset follows the template from [2, 1].

### 3 A.1 Motivation

**For what purpose was the dataset created?** Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.

4  
5 **A:** Satellite imagery is created for a variety of purposes, each serving different applications and  
6 academic research. A number of satellite-based instruments do not provide an operational per-pixel  
7 wildfire identification dataset. Recognizing the need, SIT-FUSE used these images to mark the  
8 regions where active fires are burning. By using these wildfire label data, machine learning models  
9 can learn to recognize patterns, features, and characteristics associated with wildfires, enabling them  
10 to autonomously detect and monitor such events in real-time or at scale.

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

11  
12 **A:** Dr. Nicholas LaHaye, NASA's Jet Propulsion Lab (JPL) has pre-processed the satellite imagery  
13 from various sources and created the wildfire label data, as a part of the NASA AIST funding efforts  
14 on open benchmarks.

**Who funded the creation of the dataset?** If there is an associated grant, please provide the name of the grantor and the grant name and number.

15  
16 **A:** This project has been supported by NASA Advanced Information Systems Technology (AIST)  
17 grants 21-AIST21\_2-0020 and 21-AIST21\_2-0059.

**Any other comments?**

18  
19 **A:** None.

## 20 A.2 Composition

21 **What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?** Are there multiple types of instances (e.g., movies, users, and ratings; people and interactions between them; nodes and edges)? Please provide a description.

21

22 **A:** The satellite imagery data include high-resolution images of the Earth's surface from space at  
23 different electromagnetic frequencies.

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

24

25 **A:** There are 157,469,891 samples labeled -1, 297,971,233 samples labeled 0, and 9,276 samples  
26 labeled 1.

27 **Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?** If the dataset is a sample, then what is the larger set? Is the sample representative of the larger set (e.g., geographic coverage)? If so, please describe how this representativeness was validated/verified. If it is not representative of the larger set, please describe why not (e.g., to cover a more diverse range of instances, because instances were withheld or unavailable).

27

28 **A:** The temporal coverage of satellite observations is longer than that of those prepared for FIRE-D.  
29 Additionally, while satellite-based instruments observe large regions of the Earth, FIRE-D focuses  
30 solely on the contiguous United States (CONUS). The dataset over the CONUS does not represent  
31 the entire scope of observations but is still useful to support a variety of wildfire studies using ML  
32 models. This is because almost all ML models designed to predict active fires are built for specific  
33 geographical regions.

34 **What data does each instance consist of?** "Raw" data (e.g., unprocessed text or images) or features? In either case, please provide a description.

34

35 **A:** High-resolution images.

36 **Is there a label or target associated with each instance?** If so, please provide a description.

36

37 **A:** We use JPL's Segmentation, Instance Tracking, and data Fusion Using multi-SEnsor imagery  
38 (SIT-FUSE), to generate a set of labels to mark active wildfires – previously validated on numerous  
39 pre-existing label sets for wildfires. For this task, each wildfire label set is generated at the input  
40 data's native resolution and contains 3 values:  $-1 = no\_input\_data$ ,  $0 = no\_fire$ , and  $1 = fire$ .

41 **Is any information missing from individual instances?** If so, please provide a description, explaining why this information is missing (e.g., because it was unavailable). This does not include intentionally removed information, but might include, e.g., redacted text.

41

42 **A:** No information is missing.

43 **Are relationships between individual instances made explicit (e.g., users' movie ratings, social network links)?** If so, please describe how these relationships are made explicit.

43

44 **A:** The location and time at which each image was independently taken are specified.

**Are there recommended data splits (e.g., training, development/validation, testing)?** If so, please provide a description of these splits, explaining the rationale behind them.

45

46 **A:** In this paper, we conduct unsupervised learning tasks and we do not split data into training,  
47 development/validation, testing. The dataset can be used to evaluate sensitivity of model performance  
48 with respect to different splits.

**Are there any errors, sources of noise, or redundancies in the dataset?** If so, please provide a description.

49

50 **A:** Satellite image has noise caused by atmospheric interference, sensor artifacts, cloud cover, or  
51 other environmental conditions. This kind of noise in satellite image can affect the accuracy of the  
52 wildfire label, but this is beyond the scope of this work.

**Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g., websites, tweets, other datasets)?** If it links to or relies on external resources, a) are there guarantees that they will exist, and remain constant, over time; b) are there official archival versions of the complete dataset (i.e., including the external resources as they existed at the time the dataset was created); c) are there any restrictions (e.g., licenses, fees) associated with any of the external resources that might apply to a dataset consumer? Please provide descriptions of all external resources and any restrictions associated with them, as well as links or other access points, as appropriate.

53

54 **A:** The dataset is self-contained.

**Does the dataset contain data that might be considered confidential (e.g., data that is protected by legal privilege or by doctor– patient confidentiality, data that includes the content of individuals’ non-public communications)?** If so, please provide a description.

55

56 **A:** The dataset contains no confidential information.

**Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening, or might otherwise cause anxiety?** If so, please describe why.

57

58 **A:** The dataset contains no information that that, if viewed directly or indirectly, might be offensive,  
59 insulting, threatening, or cause anxiety.

**Does the dataset identify any subpopulations (e.g., by age, gender)?** If so, please describe how these subpopulations are identified and provide a description of their respective distributions within the dataset.

60

61 **A:** The dataset does not identify any subpopulations and their respective socio-demographic charac-  
62 teristics.

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

63

64 **A:** It is not possible to identify individuals, either directly or indirectly from the dataset.

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

65

66 **A:** The dataset does not contain any sensitive information.

**Any other comments?**

67

68 **A:** None.

### 69 **A.3 Collection Process**

**How was the data associated with each instance acquired?** Was the data directly observable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or language)? If the data was reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If so, please describe how.

70

71 **A:** Generating satellite images from raw observations involves pre-processing the data to correct  
72 errors and artifacts, and geo-referencing the data.

**What mechanisms or procedures were used to collect the data (e.g., hardware apparatuses or sensors, manual human curation, software programs, software APIs)? How were these mechanisms or procedures validated?**

73

74 **A:** The pre-processed and geo-referenced observations are reconstructed into an image format,  
75 forming a spatial grid where each pixel represents a location on Earth. Additional enhancement  
76 techniques may be applied to improve visual quality or highlight specific features. The procedure  
77 follows an algorithm theoretical basis document (ATBD) for each satellite mission, and the ATBD  
78 includes the validation of the procedures.

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

79

80 **A:** The image datasets are spatial and temporal subsets of the entire archive, with FIRE-D specifically  
81 focusing on CONUS.

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

82

83 **A:** Dr. Nicholas LaHaye, NASA's Jet Propulsion Lab (JPL) has collected the satellite imagery from  
84 various sources and created the wildfire label data, as part of his commitment to the NASA-funded  
85 AIST project on open benchmarks.

**Over what timeframe was the data collected?** Does this timeframe match the creation timeframe of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please describe the timeframe in which the data associated with the instances was created.

86

87 This dataset contains 34 LandSat-8 scenes from September of 2022. Within the scenes there are  
88 455,450,400 total samples.

**Were any ethical review processes conducted (e.g., by an institutional review board)?** If so, please provide a description of these review processes, including the outcomes, as well as a link or other access point to any supporting documentation.

89

90 **A:** The ethical review process is not applicable to the collected data on Earth science processes.

**Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (e.g., websites)?**

91

92 **A:** Not applicable.

**Were the individuals in question notified about the data collection?** If so, please describe (or show with screenshots or other information) how notice was provided, and provide a link or other access point to, or otherwise reproduce, the exact language of the notification itself.

93

94 **A:** Not applicable.

**Did the individuals in question consent to the collection and use of their data?** If so, please describe (or show with screenshots or other information) how consent was requested and provided, and provide a link or other access point to, or otherwise reproduce, the exact language to which the individuals consented.

95

96 **A:** Not applicable.

**If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses?** If so, please provide a description, as well as a link or other access point to the mechanism (if appropriate).

97

98 **A:** Not applicable.

**Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data protection impact analysis) been conducted?** If so, please provide a description of this analysis, including the outcomes, as well as a link or other access point to any supporting documentation.

99

100 **A:** The satellite image datasets are publicly available, so the dataset does not require its own impact  
101 assessment.

**Any other comments?**

102

103 **A:** None.

#### 104 **A.4 Preprocessing/cleaning/labeling**

**Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)?** If so, please provide a description. If not, you may skip the remaining questions in this section.

105

106 **A:** No pre-processing or cleaning is performed, but the data represent a certain time period over  
107 CONUS.

**Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)?** If so, please provide a link or other access point to the “raw” data.

108

109 **A:** The raw data are publicly available from:

110 <https://search.earthdata.nasa.gov/search?q=HLS%20Daily%20Global>, and  
111 <https://scihub.copernicus.eu/dhus/#/home>.

**Is the software that was used to preprocess/clean/label the data available?** If so, please provide a link or other access point.

112

113 **A:** The software is publicly available here:

114 [https://github.com/nlahaye/SIT\\_FUSE/](https://github.com/nlahaye/SIT_FUSE/).

**Any other comments?**

115

116 **A:** None.

## 117 **A.5 Uses**

**Has the dataset been used for any tasks already?** If so, please provide a description.

118

119 **A:** Images from Earth observation satellites are used to generate scientific level 2 products. Level 2  
120 products typically involve the processing and analysis of satellite image to derive specific parameters  
121 or measurements. These products are often more refined and standardized, catering to scientific  
122 research, analysis, and modeling in atmospheric sciences.

**Is there a repository that links to any or all papers or systems that use the dataset?** If so, please provide a link or other access point.

123

124 **A:** GeoQuery is available at <https://github.com/186philip-davis/GeoQuery>.

**What (other) tasks could the dataset be used for?**

125

126 **A:** The wildfire labels can be used as input for physics-based numerical models, which can simulate  
127 transport of wildfire-induced smoke plumes and their impact.

**Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?** For example, is there anything that a dataset consumer might need to know to avoid uses that could result in unfair treatment of individuals or groups (e.g., stereotyping, quality of service issues) or other risks or harms (e.g., legal risks, financial harms)? If so, please provide a description. Is there anything a dataset consumer could do to mitigate these risks or harms?

128

129 **A:** Given that there has not been a major update in generating images from raw observations, the  
130 satellite imagery and wildfire labels can be used in their current form.

**Are there tasks for which the dataset should not be used?** If so, please provide a description.

A: We do not anticipate any restrictions on the tasks that the dataset should not be used.

**Any other comments?**

A: None.

## A.6 Distribution

**Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?** If so, please provide a description.

A: The dataset is publicly available.

**How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the dataset have a digital object identifier (DOI)?**

A: The wildfire labels dataset is available at [https://osf.io/v4uz9/?view\\_only=a6f58def39d344dba81952968ea234cd](https://osf.io/v4uz9/?view_only=a6f58def39d344dba81952968ea234cd)

**When will the dataset be distributed?**

A: This dataset is already available.

**Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)?** If so, please describe this license and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions.

A: The access repository and dataset area both provided under the MIT license, and can be accessed without cost.

**Have any third parties imposed IP-based or other restrictions on the data associated with the instances?** If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these restrictions.

A: We are not aware of any third-party restrictions or costs imposed.

**Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?** If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any supporting documentation.

A: The dataset does not require export controls or other regulatory restrictions.

**Any other comments?**

A: None.

152 **A.7 Maintenance**

153 **Who will be supporting/hosting/maintaining the dataset?**

154 **A:** The dataset is hosted by the Open Science Foundation (`osf.io`). The dataset and access repo  
155 will be supported by Philip Davis at the University of Utah, as part of NASA AIST initiative on open  
156 benchmarks.

157 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

158 **A:** philip.davis@sci.utah.edu.

159 **Is there an erratum?** If so, please provide a link or other access point.

160 **A:** There is no erratum.

161 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete instances)?** If so, please describe how often, by whom, and how updates will be communicated to dataset consumers (e.g., mailing list, GitHub)?

162 **A:** Errors in the dataset will be corrected on an as-necessary basis. Updates will be communicated on  
163 the OSF repository.

164 **If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were the individuals in question told that their data would be retained for a fixed period of time and then deleted)?** If so, please describe these limits and explain how they will be enforced.

165 **A:** The dataset does not relate to people.

166 **Will older versions of the dataset continue to be supported/hosted/maintained?** If so, please describe how. If not, please describe how its obsolescence will be communicated to dataset consumers.

167 **A:** Older versions will not be hosted or maintained. Obsolescence will be communicated on the OSF  
168 repository, if necessary.

169 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so?** If so, please provide a description. Will these contributions be validated/verified? If so, please describe how. If not, why not? Is there a process for communicating/distributing these contributions to dataset consumers? If so, please provide a description.

170 **A:** Users may fork and modify the OSF and/or GeoQeury repository using common facilities for  
171 doing so.

172 **Any other comments?**

173 **A:** None.



174 **References**

- 175 [1] Julien Cornebise, Ivan Oršolić, and Freddie Kalaitzis. Open high-resolution satellite imagery:  
176 The worldstrat dataset—with application to super-resolution. *Advances in Neural Information*  
177 *Processing Systems*, 35:25979–25991, 2022.
- 178 [2] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach,  
179 Hal Daumé Iii, and Kate Crawford. Datasheets for datasets. *Communications of the ACM*,  
180 64(12):86–92, 2021.