Multi-Source Weak Supervision Fusion for Disaster Scene Recognition in Videos

Maria Presa-Reyes*, Yudong Tao[†], Rui Ma[†], Shu-Ching Chen*, Mei-Ling Shyu[†]

*Knight Foundation School of Computing and Information Sciences

Florida International University, Miami, Florida 33199

Emails: {mpres029,chens}@cs.fiu.edu

[†]Department of Electrical and Computer Engineering

University of Miami, Coral Gables, FL 33124

Emails: {yxt128,rxm1351,shyu}@miami.edu

Abstract-Images or video recordings assist emergency responders in quickly inspecting the damage after a disaster event. New techniques are needed to help responders organize and find important information at the right time. However, most existing methods do not meet public safety standards due to a lack of training data. We propose a multi-source weak supervision fusion technique to train on a highly imbalanced dataset annotated with noisy labels. Using a Confident Learning technique, we reduce the noise effect while boosting the class labels' quality. We combine the predictive power from models trained on large-scale visual datasets using Differential Evolution. This research demonstrates a fully-automatic approach with great potential to reduce the required time and resources while delivering exceptional results. In the TRECVID2021 Disaster Scene Description and Indexing (DSDI) Challenge, our technique achieved the top score among all the submitted runs, independent of the training data utilized.

Keywords-damage assessment; deep learning; convolutional neural networks

I. Introduction

Image and video recognition algorithms have advanced rapidly and with better precision, and are expected to become a critical component of incident and disaster responses [1]. Using advanced technologies and deep learning methodologies such as Convolutional Neural Networks (CNNs), it is possible to deploy a drone ahead of the search team to swiftly identify the most damaged areas that should be prioritized during a disaster. The automated content-based analysis and classification of the observed disaster-related features in recorded videos will allow better curation and retrieval of critical information for situational awareness. Due to insufficient training data and standards, most of the existing methods do not fulfill public safety demands [2].

Civil Air Patrol (CAP) has the technical capability to function even when severe weather disrupts power, internet, phones, and airplane takeoffs, making it a critical and cost-effective tool for the Federal Emergency Management Agency (FEMA) to survey the impacted region swiftly and efficiently. CAP offers aerial pictures of flooded areas, collapsed dams, and other natural disaster-related events. To this end, several large-scale disaster imagery datasets, including the Incidents Dataset [3], LADI (Low Altitude

Disaster Imagery) [4], xBD [5], etc., have been recently released to stimulate the development of new research and technologies in this field. Given the volume of data being collected, it is also critical to develop sophisticated tools and systems for curating all of the information.

It is challenging to analyze the images taken by lowaltitude planes since they have a low height perspective, an oblique angle, and many disaster-related parts that image recognition systems do not usually take into account. We propose a weakly-supervised learning technique that incorporates data from a range of sources, many of which are of low quality or have been trained on subjects significantly different from the target classification task. The proposed fully-automatic solution would significantly decrease the time and expense associated with the classification jobs while delivering superior outcomes.

The main contributions of this paper are summarized as follows.

- We propose a new semi-supervised training technique that is robust to noisy, limited, and erroneous annotations and class labels from multiple sources.
- For the multi-source weak supervision fusion framework, a unique approach for recognizing and merging the relevant predictions from various pre-trained networks is proposed.
- The proposed approach is evaluated on the LADI dataset and achieved the top score among all the submitted runs in the TRECVID2021 [6] Disaster Scene Description and Indexing (DSDI) Challenge, independent of the training data utilized.

This paper is organized as follows. Section II examines approaches that use deep learning techniques to analyze low-altitude images. Section III introduces our proposed weakly-supervised framework, including confident learning and multi-source weak-supervision fusion. In Section IV, the effectiveness of our proposed framework is shown through the quantitative experimental results. Finally, Section V summarizes this paper and recommends future research.

II. RELATED WORK

Most current solutions rely on high-quality annotations to build reliable models that can sufficiently automate image processing and concept detection. Non-experts are likely to have only seen low-altitude photos on rare occasions. Consequently, it will be too costly to get enough highquality annotations to build a good training dataset. Numerous researchers have developed a variety of deep learning algorithms that are less reliant on the quality of the training data. The weakly-supervised tags and visual information are used to train semantic-aware hash functions [7]. Previously, deep canonical correlation analysis (DCCA) [8] was used to combine visual and text tag data. Many previously reported techniques rely on sparse line reconstruction, sparse coding, and dictionary learning to recover textual tags, which costs time and space and is not suited for large-scale applications. Research into automated disaster scene descriptions from images has grown in popularity. Newly-released disaster datasets such as xBD [5] and the Incidents Dataset [3] feature a top-down and a ground-level view of the damages. However, LADI [4] is unique in the low-altitude and oblique views found in its images. More recent studies explore an ensemble learning approach to tackle the class-imbalance and noisy-label issues [9], [10]. The incorporation of spatiotemporal information to increase the model's contextual awareness has also been investigated [11]. Our proposed framework aims to improve the quality of noisy labels in the LADI training data through a Confident Learning (CL) [12] strategy. Furthermore, a novel multi-source information fusion method is proposed to improve the performance of the target features that are underrepresented in LADI.

III. PROPOSED FRAMEWORK

Figure 1 illustrates the full flow of our proposed framework. CL is used to improve the quality of the noisy labels in the crowdsourced annotated training set, which is the first step in our multi-source architecture for combining weak supervision from different sources. Given the scores of numerous semantic concepts obtained from different machine annotators, several semantically related predictions are used to improve the performance of a target feature. The text that describes the target feature is turned into high-dimensional vectors, which are then used to look for semantic similarity and pick relevant concepts from other networks. We optimize a weighted average that incorporates all of the models' relevant predictions into a single scalar that serves to rank the video clip using Differential Evolution (DE).

A. Denoising with Confident Learning

According to the LADI researchers [4], annotations are organized as Human Intelligence Task (HIT) which asks the human worker whether any of the target features in each of the five categories (i.e., *damage*, *environment*, *infrastructure*, *water*, and *vehicle*) are correct. Each HIT is allocated to up

to five workers (asking just one category at a time) in order to reach the agreement on the label quality. Namely, for an image i and a target feature F_C that belongs to a specific category C (i.e., $F_C \in C$), the initial soft score S_{i,F_C} is calculated as follows.

$$S_{i,F_C} = \frac{\#Positive\ Votes_{i,F_C}}{Total\ Votes_{i,C}} \tag{1}$$

To calculate S_{i,F_C} , we assume that a particular image must have at least one vote from an annotator who was assigned a specific category C. Then, we employ cross-validation confident intervals [12] to derive out-of-sample prediction probabilities to further improve the label quality.

B. Multi-Source Weak Supervision Fusion

1) Machine Annotators: This study employs five CNN network configurations (i.e., ResNet50, DenseNet161, YOLOv4, ViT-B/16, and InceptionV3) pre-trained on five open-source datasets (i.e., Places365, Incidents Dataset, MS COCO, ImageNet21k, and LADI+Others). ResNet50 [13] and DenseNet161 are pre-trained on Places365 dataset which contains 1.8 million training images taken from 365 scene categories [14]. Another ResNet50 network is also pretrained on the Incidents Dataset containing 446,844 manually annotated images covering 43 incidents across various scenes [3]. YOLOv4 (You Only Look Once) [15] pre-trained on Microsoft Common Objects in Context (MS COCO) is one of the leading deep learning-based object detection frameworks. The ViT-B/16 [16] model pre-trained on the ImageNet21K dataset is proven to be a key component in our proposed framework. Last but not least, an InceptionV3 model trained on LADI plus other sources by Presa-Reyes et al. [11] have also been incorporated.

2) Multi-Source Concept Fusion: Given the predicted scores X^p of many machine annotators' semantic concepts (i.e., target classes), many of these related concepts may help identify a target feature F. A Universal Sentence Encoder based on the Deep Averaging Network (DAN) [17] converts the text describing the target feature into high-dimensional vectors T that are then utilized to obtain the semantic similarity among different concepts using the cosine distance θ of the vectors. To fuse multi-source concepts, the high-dimensional vectors of the target feature F and the semantic concept P are first matched, and the weighted average score of those closely correlated concepts are fused, i.e.,

$$S_F(k, w_F) = \sum_{p \in O} w_F^p \cdot X_k^p \tag{2}$$

where $O = \{P | \theta(T_F, T_P) > \vartheta\}$, $w_F^p \in \mathbb{Q}$ is the set of optimized weights representing the contributing power of each pre-trained model's predicted score X_k^p for a key frame k, and w_F is the vector with all w_F^p . Moreover, the values for w_F^p bounds and the ϑ threshold are empirically decided based on the validation performance. Furthermore,

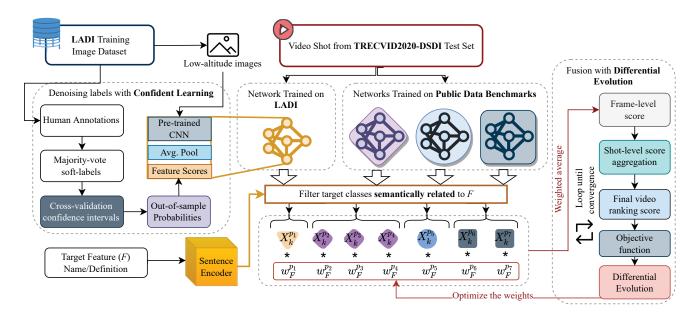


Figure 1: The proposed weakly-supervised deep learning framework implements a confident learning approach to denoise crowdsourced annotations along with a multi-modality fusion framework to search and combine relevant target features predicted by multiple networks.

there exist multiple key frames inside a given video shot v. Therefore, the average score over all the key frames in v is computed as the shot-level feature score, which can be formally written as

$$S_V(V, w_F) = \frac{1}{||V||} \sum_{k \in V} S_F(k, w_F)$$
 (3)

Then, for a given dataset of video shots $\mathcal V$ and a target feature F, the top-N shots with F can be defined as an ordered sequence $V_F = [V_1, V_2, \ldots, V_N]$, where $V_i \in \mathcal V$ and $\forall i > j, S_V(V_j, w_F) > S_V(V_i, w_F)$.

3) Weight Optimization based on Differential Evolution: The remaining problem is to determine the optimal weights w_F for each target feature F. Differential Evolution (DE) [18] is a kind of evolutionary optimization technique that works with a population of candidate solutions. It uses genetic operators like mutation and recombination to repeatedly enhance the population. The objective function G determines each candidate's fitness. If $G(s_1) < G(s_2)$, candidate s_1 is judged to be superior to candidate s_2 . The objective function seeks to improve the average precision for a specific target feature (i.e., to minimize $1 - AP^N$) by measuring the performance of a collection of retrieved results using the precision and recall metrics. Assuming the solution contains N video shots ordered by the final aggregated confidence scores, our objective is to minimize the error as shown in Equation 4. Semantically relevant predictions are then combined into a single scalar which is used to score the video clip.

$$\hat{w}_F = \arg\min_{w_F} G(w_F) = \arg\min_{w_F} \left[1 - AP^N(V_F) \right] \quad (4)$$

IV. EXPERIMENT RESULTS

A. Experimental Setup

This section discusses the LADI dataset to evaluate the proposed methodology compared to other currently available approaches to determine how effective it is. We fine-tune the weights of the CNN network trained on ImageNet to train the feature score models using only the LADI dataset via transfer learning. Using DE, the feature scores that have been predicted are then integrated with the predictions made by CNN networks pre-trained on relevant open dataset benchmarks.

1) Dataset: We test our methods using the LADI dataset, which comprises images acquired by CAP from a low-flying aircraft and maintained by FEMA. The LADI training dataset consists of images captured from an airplane, and the LADI test dataset consists of brief video clips captured from a UAV. The DSDI track's test dataset in 2021 comprises 2,802 video shots with a maximum duration of 60 seconds per shot, focusing on the devastation wrought by an earthquake tragedy. The test set supplied in TRECVID2020-DSDI [20] is used as validation during the DE processing in our case. The Mean Average Precision (MAP) metric is used to examine and compare the performance of different approaches.

Table I: Performance comparison among our proposed technique and competing methods.

Method	Training	Precision@k			Recall@k			F1@k			MAP
	Data	k=10	k=100	k=1000	k=10	k=100	k=1000	k=10	k=100	k=1000	IVIAI
BUPT_MCPRL [6]	L	0.271	0.225	0.228	0.271	0.232	0.405	0.271	0.227	0.244	0.159
VCL_CERTH [19]	L+	0.510	0.367	0.245	0.511	0.415	0.378	0.510	0.377	0.255	0.282
Presa-Reyes et al. [11]	O	0.413	0.392	0.285	0.413	0.448	0.682	0.413	0.404	0.316	0.298
Ours-CL-BA	L	0.394	0.346	0.279	0.394	0.383	0.648	0.394	0.351	0.307	0.254
Ours-CL-ZS	O	0.471	0.409	0.296	0.471	0.522	0.789	0.471	0.425	0.332	0.339
Ours-CL-DE (proposed)	L	0.384	0.351	0.286	0.384	0.395	0.683	0.384	0.360	0.315	0.268
	O	0.481	0.425	0.310	0.481	0.502	0.793	0.481	0.439	0.345	0.359

L+ LADI-based (L) training data plus additional human annotations (i.e., instance and segmentation).

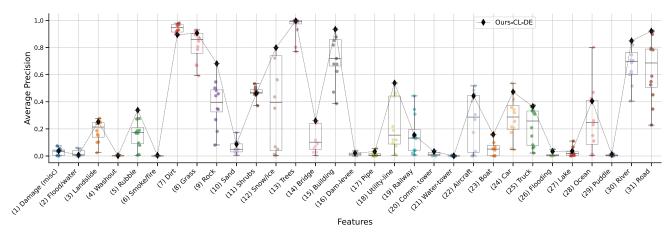


Figure 2: The boxplot shows the distribution for a feature's precision score compared across all submissions to TRECVID2021-DSDI, independent of which training dataset was used to train each technique. The placement of our proposed method's performance is demonstrated using a black diamond.

- 2) Competing Methods: To determine the effectiveness of the proposed technique, we compare it to other competing methods, such as BUPT_MCPR [6] and VCL_CERTH [19]. Both competing methods are trained solely on the LADI-dataset. In particular, the problem was approached by VCL_CERTH as a panoptic segmentation problem with additional instance and semantic segmentation annotations for 300 LADI images. The method proposed by Presa-Reyes et al. [11] trained on LADI plus other datasets was also included. Two baseline fusion techniques using the average of the best performing model (Ours-CL-BA) and aggregated predictive scores after z-score normalization (Ours-CL-ZS) are also explored to compare against our proposed DE fusion.
- 3) Feature Score Model and Fusion: Two feature score models, EfficientNet-B5 [21] and ResNet50 [13], are trained on the LADI's confident labels generated by the CL-based approach. Using transfer learning, we fine-tune the network's weights on ImageNet. The network's final classification head is replaced with a fully-connected layer followed by a sigmoid activation for multi-class soft-label classification. With a starting learning rate of $(\eta = 1e 4)$, we use the Adam solver to optimize our model. We use a weighted average

ensemble where the weights are optimized through DE to get superior performance by combining human and machine-generated annotations. For the DE search, we employ the DE/best/1/bin technique which generates new candidate solutions by randomly picking solutions from the population, subtracting one from the other, and adding a scaled version of the difference to the population's best candidate solution.

B. Results and Discussion

The proposed framework is compared to competing methods mainly categorized as LADI-based (L) the LADI + Others (O) track submission—where "Others" in our proposed approach refers to the inclusion of models pre-trained on open-source data benchmarks. Table I summarizes the performance comparison across different methods. The excellent results obtained by the panoptic segmentation approach proposed by VCL_CERTH on the LADI-based (L) track underline the necessity to integrate additional information about the images other than the noisy labels.

Our proposed technique achieves impressive results on the LADI + Others (O) track, particularly compared to other competing methods. The high recall rate illustrates our classification model's ability to detect and recover the majority of positive examples within a relevant target feature. By

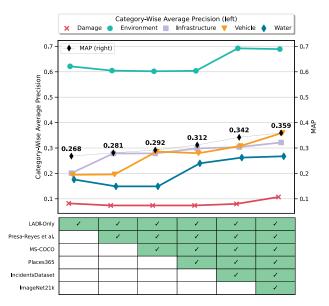


Figure 3: Ablation research demonstrating how the performance of the proposed Ours-CL-DE improves with the inclusion of each machine annotator's dataset.

comparing the baseline methods Ours-CL-BA and Ours-CL-ZS, we demonstrate our proposed Ours-CL-DE approach can find better weights when aggregating the predictions of different models. Furthermore, compared to the method that simply trains on the LADI-based (L) data, the proposed method introduced on LADI + Others (O) improves the MAP score by roughly 34%, indicating the effectiveness of our strategy of fusing the weak supervision from multiple sources.

In Figure 2, the average precision at the target feature level shows that our suggested approach has obtained the greatest performance for the target features such as debris, rock, snow/ice, building, utility-line, boat, river, and road. Figure 3 depicts the performance contribution of each additional dataset used to train the machine annotators previously described in Section III-B1. Starting from our proposed CL technique trained on LADI only, each additional dataset is added to the ensemble as depicted by a checkmark in the figure. The ResNet50 pre-trained on the Incidents Dataset contributed a performance boost for the environment category, detecting concepts such as 'snow covered' and 'field' and improving on features snow/ice and grass. The damage features, on the other hand, did not improve as expected given the damage concepts from the Incidents Dataset, necessitating further investigation. The environment features achieve better performance because they are simpler to discern from long distances and show lower inter-class variation than other categories. YOLOv4 network pre-trained on MS COCO contributes a performance boost for the vehicle categories, detecting concepts for 'aeroplane,' 'boat,'

'car,' and 'truck.'

We employ a weighted average ensemble that achieves better performance thanks to the integration of human and machine-generated annotations. Since it is clear that the relationship between the relevant features is not linear to their semantic similarity, our proposed technique has been proven to be a viable approach to identify the best predictions based on the performance of each machine annotator. Because our proposed technique outperforms existing methods with minimal training, they are an excellent means of leveraging and transferring information from the methods that have already been presented in previous research into any emerging topic.

V. CONCLUSION AND FUTURE WORK

Due to a lack of appropriate training data, most presentday picture recognition algorithms fail to meet public safety requirements. As part of our multi-source weak supervision fusion architecture, we apply the CL technique to enhance the quality of noisy labels in the crowdsourced annotated training set. Semantic similarity is used to identify relevant concepts predicted by other networks. We use DE to rank the video clip based on a weighted average of all relevant model predictions. Combining many classifiers pre-trained on wellknown data benchmarks improves the overall performance, but only the best and most relevant predicted score towards a particular target feature should be used. Overall, the study shows how this framework has great potential to save a significant amount of time and resources while still achieving outstanding results in the disaster scene description task. Although this work focuses on disaster scene description, the proposed methods have been developed with extendability in mind. Our approaches are effective for leveraging and transferring knowledge from past study into any new topic. As a potential future work, we will explore more advanced techniques of incorporating other multi-modality sources using our proposed technique, such as spatio-temporal data.

ACKNOWLEDGMENT

For Shu-Ching Chen, this research is partially supported by NSF CNS-1952089 and CNS-2125165.

REFERENCES

- [1] M.-F. R. Lee and T.-W. Chien, "Artificial intelligence and internet of things for robotic disaster response," 2020 International Conference on Advanced Robotics and Intelligent Systems (ARIS), pp. 1–6, 2020.
- [2] J. Evans, "Artificial intelligence and public standards: Report," Feb 2020. [Online]. Available: https://www.gov.uk/government/publications/artificial-intelligence-and-public-standards-report
- [3] E. Weber, N. Marzo, D. P. Papadopoulos, A. Biswas, A. Lapedriza, F. Ofli, M. Imran, and A. Torralba, "Detecting natural disasters, damage, and incidents in the wild," in *The European Conference on Computer Vision*, August 2020.

- [4] J. Liu, D. Strohschein, S. Samsi, and A. Weinert, "Large scale organization and inference of an imagery dataset for public safety," in *IEEE High Performance Extreme Computing Conference*, Sep. 2019, pp. 1–6.
- [5] R. Gupta, B. Goodman, N. N. Patel, R. Hosfelt, S. Sajeev, E. T. Heim, J. Doshi, K. Lucas, H. Choset, and M. E. Gaston, "Creating xbd: A dataset for assessing building damage from satellite imagery," *ArXiv*, vol. abs/1911.09296, 2019.
- [6] G. Awad, A. A. Butt, K. Curtis, J. Fiscus, A. Godil, Y. Lee, A. Delgado, J. Zhang, E. Godard, B. Chocot, L. Diduch, J. Liu, Y. Graham, G. J. F. Jones, and G. Quénot, "Evaluating multiple video understanding and retrieval tasks at trecvid 2021," in *TRECVID*. NIST, USA, 2021.
- [7] J. Tang and Z. Li, "Weakly supervised multimodal hashing for scalable social image retrieval," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 28, pp. 2730–2741, 2018.
- [8] G. Andrew, R. Arora, J. Bilmes, and K. Livescu, "Deep canonical correlation analysis," in *International conference* on machine learning. PMLR, 2013, pp. 1247–1255.
- [9] Y. Li, H. Wang, S. Sun, and B. Buckles, "Integrating multiple deep learning models to classify disaster scene videos," in 2020 IEEE High Performance Extreme Computing Conference, 2020.
- [10] S. Okazaki, Q. Kong, M. Klinkigt, and T. Yoshinaga, "Hitachi at TRECVID DSDI 2020," in TRECVID. NIST, USA, 2020.
- [11] M. Presa-Reyes, Y. Tao, S.-C. Chen, and M.-L. Shyu, "Deep learning with weak supervision for disaster scene description in low-altitude imagery," *IEEE Transactions on Geoscience and Remote Sensing*, pp. 1–1, 2021.
- [12] C. Northcutt, L. Jiang, and I. Chuang, "Confident learning: Estimating uncertainty in dataset labels," *Journal of Artificial Intelligence Research*, vol. 70, pp. 1373–1411, 2021.
- [13] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [14] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [15] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," arXiv preprint arXiv:2004.10934, 2020.
- [16] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint* arXiv:2010.11929, 2020.

- [17] D. Cer, Y. Yang, S.-y. Kong, N. Hua, N. Limtiaco, R. St. John, N. Constant, M. Guajardo-Cespedes, S. Yuan, C. Tar, B. Strope, and R. Kurzweil, "Universal sentence encoder for English," in *Proceedings of the 2018 Conference* on Empirical Methods in Natural Language Processing: System Demonstrations. Brussels, Belgium: Association for Computational Linguistics, Nov. 2018, pp. 169–174. [Online]. Available: https://aclanthology.org/D18-2029
- [18] R. Storn and K. Price, "Differential evolution—a simple and efficient heuristic for global optimization over continuous spaces," *Journal of global optimization*, vol. 11, no. 4, pp. 341–359, 1997.
- [19] E. Christakis, Z. Batzos, K. Konstantoudakis, K. Christaki, P. Boutis, D. Sainidis, D. Tsiakmakis, G. Almpanis, T. Dimou, and P. Daras, "Low altitude image analysis using panoptic segmentation," in *TRECVID*. NIST, USA, 2021.
- [20] G. Awad, A. A. Butt, K. Curtis, Y. Lee, J. Fiscus, A. Godil, A. Delgado, J. Zhang, E. Godard, L. Diduch, J. Liu, A. F. Smeaton, Y. Graham, G. J. F. Jones, W. Kraaij, and G. Quénot, "TRECVID 2020: comprehensive campaign for evaluating video retrieval tasks across multiple application domains," in TRECVID. NIST, USA, 2020.
- [21] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *International Conference* on Machine Learning. PMLR, 2019, pp. 6105–6114.