
Assessing Post-Disaster Damage from Satellite Imagery using Semi-Supervised Learning Techniques

Jihyeon Lee^{1,3} Joseph Z. Xu¹ Kihyuk Sohn¹ Wenhan Lu¹ David Berthelot¹

Izzeddin Gur¹ Pranav Khaitan¹ Ke-Wei (Fiona) Huang² Kyriacos Koupparis²

Bernhard Kowatsch²

¹Google Research, ²United Nations World Food Programme, ³Stanford University

jihyeon@cs.stanford.edu, {jzxu, kihyuks, wenhan, dberth, izzeddin}@google.com

Abstract

To respond to disasters such as earthquakes, wildfires, and armed conflicts, humanitarian organizations require accurate and timely data in the form of damage assessments, which indicate what buildings and population centers have been most affected. Recent research combines machine learning with remote sensing to automatically extract such information from satellite imagery, reducing manual labor and turn-around time. A major impediment to using machine learning methods in real disaster response scenarios is the difficulty of obtaining a sufficient amount of labeled data to train a model for an unfolding disaster. This paper shows a novel application of semi-supervised learning (SSL) to train models for damage assessment with a minimal amount of labeled data and large amount of unlabeled data. We compare the performance of state-of-the-art SSL methods, including MixMatch [2] and FixMatch [28], to a supervised baseline for the 2010 Haiti earthquake, 2017 Santa Rosa wildfire [15], and 2016 armed conflict in Syria [32]. We show how models trained with SSL methods can reach fully supervised performance despite using only a fraction of labeled data and identify areas for further improvements.

1 Introduction

When a humanitarian crisis such as a natural disaster occurs, crisis responders need to know the locations of affected populations to facilitate relief efforts. The locations and density of damaged buildings serve as a useful proxy to estimate this information [6]. One approach to identify them is remote sensing: expert analysts compare pre- and post-disaster satellite imagery of the affected region and mark the locations of damaged buildings. However, this is time-consuming (fewer than 100 buildings assessed per hour per person), challenging to scale, and prone to error [29, 19].

Machine learning (ML) has been utilized as an efficient tool to automate the damage assessment process [4, 18, 9, 15, 32, 15, 16, 30]. Models are trained to distinguish between images of damaged and undamaged buildings using expert-labeled images of past disasters. The trained models can analyze entire cities in a matter of minutes when deployed in modern data centers.

It remains challenging to build accurate models for new disasters. ML has traditionally relied on datasets with sufficient variation and coverage so that models generalize to unseen examples at inference time. This is not possible in our domain because there is a limited number of past disasters, and the appearance of each disaster or geographical region varies widely in their layout of buildings,

construction material, vegetation, appearance of damage, etc. [24] (see Figure 3 in Appendix). Furthermore, even if the model has been pre-trained on the same disaster type and location, there is noise inherent to satellite imagery due to changing landscapes, seasonal variation, cloud cover, and other factors that causes the inference data to systematically differ from the training data. Therefore, models trained only on past disasters will likely under-perform in new disasters.

We can avoid the generalization problem by training models on data from a new disaster after it strikes. Models trained in this way must use only a small number of labeled training examples, since manual expert labeling is time-consuming. At the onset of a new disaster, labeled data is limited, but an abundant amount of unlabeled satellite imagery can be automatically extracted from the region. Recent advances in semi-supervised learning (SSL) techniques show that algorithms combining labeled and unlabeled training data can achieve performance comparable to fully supervised algorithms trained on orders of magnitude more labeled data [2, 28, 1]. By using SSL techniques, we can train accurate damage assessment models for new disasters without spending much time gathering manual labels.

In this paper, we apply two SSL techniques, MixMatch [2] and FixMatch [28], to train building damage detection models using a limited amount of labeled data. We evaluate on three different disasters, varying the amount of labeled examples. We show that the models can get close to the level of fully supervised results using only a fraction of the labeled data.

2 Related Work

Machine Learning in Damage Building Assessment Past studies from fully supervised settings have successfully applied machine learning approaches to building damage detection from satellite imagery. The public xBD dataset [13, 15] was released alongside the xView2 Challenge [14], providing large-scale satellite imagery, building polygons, and ordinal labels that denote damage level across 19 disasters with the task of per-pixel classification. The first-place approach [10] had two stages, initially training a localization model with only pre-disaster images and then using the weights to initialize a Siamese Neural Network for building classification that shares weights between the pre-disaster and post-disaster images. Gupta et al. [16] proposed an end-to-end approach, first extracting multi-scale image features, feeding them into a segmentation head to predict buildings independently on the pre- and post-disaster images, and finally classifying each pixel. Weber et al. [30] trained a single network, modeling the task as semantic segmentation.

The above methods generate training and validation data from the same distribution, so they do not address the problem of running inference for a newly unfolding disaster with minimal data. Xu et al. [32] developed models for the Haiti, Mexico, and Indonesia earthquakes and conducted cross-region generalization experiments, showing how models pretrained on past disasters did not perform well in a new region with no or minimal labeled data. We use this method as a baseline in our experiments. Valentjin et al. [29] ran experiments on 13 disasters differing in hazard type, geographical region, and satellite parameters and found performance varied significantly across test disasters, whether or not data from the test disaster was included in training.

Semi-Supervised Learning Approaches Semi-supervised learning (SSL) provides approaches to alleviate the need for large amounts of labeled data by leveraging unlabeled data. There is a class of SSL methods for deep networks that perform pseudo-labeling (or self-training) [23, 25, 31, 27], generating artificial labels from the unlabeled data and involving a minimal amount of human labor [34, 2, 1, 28]. These techniques rely on consistency regularization [26, 22], which encourages the model to output predictions of a similar distribution across perturbations of a given input. MixMatch [2] adds other types of regularization, using MixUp [34] to encourage convex behavior “between” examples by generating weighted combinations of labeled and unlabeled ones. FixMatch [28] presented a simpler approach that achieved state-of-the-art performance on common SSL benchmarks, such as CIFAR [20] and STL [3]. The methods differ in how they use the pseudo-labels to calculate loss. For a given unlabeled image, MixMatch creates a guess label based on its weakly augmented versions and calculates loss based on how well the model predicts that label. FixMatch also generates a pseudo-label from weak augmentations, but it calculates loss on whether the model is able to predict the label on strongly augmented versions. The driving assumptions of pseudo-labeling SSL methods should hold true for satellite imagery and provide a way to take advantage of unlabeled data.

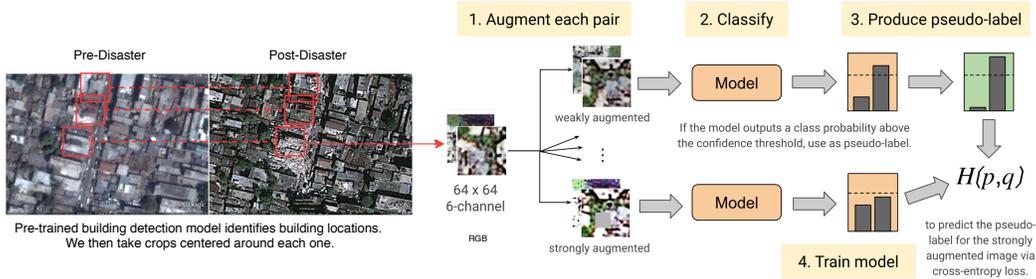


Figure 1: Pipeline that shows how the pre- and post-disaster images are stacked into a single 6-channel input and then augmented to produce pseudo-labels for training. This diagram includes strong augmentation to represent FixMatch, but MixMatch uses only weakly augmented images.

3 Data

We evaluated our approach on three disasters. In addition to the Santa Rosa wildfire [15], we generated our own datasets for the 2010 Haiti earthquake and 2016 snapshot of Aleppo as in [32]. First, we obtained imagery from before and after the disaster for each region, mainly from DigitalGlobe’s WorldView 2 and 3 satellites. For Haiti, candid flyover images were provided by the National Oceanic and Atmosphere Administration. We resampled all images to 0.3 meter resolution for consistency.

Next, we obtained positive ground truth labels from building damage assessments provided by UNOSAT, the operational satellite applications programme of the United Nations Institute for Training and Research (UNITAR), available on the Humanitarian Data Exchange website [17]. UNOSAT assessments use a 5-level scale to measure damage, but the labels were noisy and inconsistent across different datasets. Therefore, we grouped “Severe Damage” and “Destroyed” into a single “Damaged” class and formulated our problem as a binary classification problem to identify “Damaged” and “Undamaged” buildings. To acquire Undamaged examples, we used a pretrained building detection model [32] to identify buildings and filtered out the ones marked as damaged in UNOSAT assessments.

Finally, to create the training examples, we sampled crops centered around each building. We then aligned the pre- and post-disaster imagery and used Google Earth Engine [12] to spatially join the labels and cropped images. Each example in our dataset contains a 6-channel, 64 x 64 image and a classification label (0 for undamaged, 1 for damaged). There are 50,742 Haiti examples (44% positive), 12,897 Santa Rosa examples (27% positive), and 10,452 Aleppo examples (44% positive). In a newly unfolding disaster, this volume of labeled data would not be available. To simulate this, a random sample of the examples were considered labeled and the rest unlabeled (Section 5).

4 Approach

We formulate our task as binary classification, where the two classes are undamaged (0) or damaged (1). We define a batch of B labeled examples as $\mathcal{X} = \{(x_b, p_b) : b \in (1, \dots, B)\}$, where x_b contains the 6-channel example images and p_b are the one-hot labels. We also define a batch of μB unlabeled examples as $\mathcal{U} = \{u_b : 1, \dots, \mu B\}$, where μ is a hyperparameter that determines the relative sizes of \mathcal{X} and \mathcal{U} . We trained classification models using MixMatch and FixMatch, which use pseudo-labeling to create artificial labels for \mathcal{U} . Then, a convolutional neural network (CNN) is trained to predict on augmented examples produced from both \mathcal{X} and \mathcal{U} , using a loss function that has a labeled and unlabeled loss term. We compare the differences in pseudo-labeling and loss functions below.

Producing Pseudo-Labels We first review the label guessing step of **MixMatch** [2]. MixMatch initially applies weak augmentation, which consists of random flips, rotations, and shifts, to both \mathcal{X} and \mathcal{U} . Let $\alpha(\cdot)$ denote weak augmentation, such that, for a given unlabeled batch u_b , $\alpha(u_b)$ contains K augmentations for each example. We compute the average of the model’s predicted class distributions across all augmentations, which is $\bar{q}_b = \frac{1}{K} \sum_{k=1}^K p_{model}(y|\alpha(u_b, k))$. Then, the method applies sharpening, which reduces the entropy of the label distribution as introduced in [11]. $q_b = \text{Sharpen}(\bar{q}_b)$ acts as the target for the model’s prediction on an augmentation of u_b .

Finally, MixMatch takes the augmented versions of the labeled data $\alpha(\mathcal{X}) = ((\alpha(x_b), p_b); b \in (1, \dots, B))$ and unlabeled data $\alpha(\mathcal{U}) = ((\alpha(u_b), p_b); b \in (1, \dots, B))$ and applies MixUp [34], mixing both the labeled and unlabeled examples with label guesses. We refer to the final labeled dataset with augmentations and MixUp applied as \mathcal{X}' and the final unlabeled dataset as \mathcal{U}' .

Now, we review the pseudo-labeling step of **FixMatch** [28], which presented a relatively simpler method but performed better on benchmark datasets. Similarly to MixMatch, FixMatch creates the targets from the weakly augmented versions of a given unlabeled example. However, without sharpening the output distribution, it uses $\arg \max(q_b)$ as the pseudo-label. FixMatch uses not only weak augmentation but also strong augmentation, denoted by $\mathcal{A}(\cdot)$. In this work, we use FixMatch with RandAugment (RA) [1], which randomly applies a series of transformations. The types of transformations included changing the brightness, contrast, or saturation level, solarizing, posterizing, applying CutOut, and more [28]. FixMatch makes use of both weakly and strongly augmented images by encouraging the model to predict the pseudo-label for the strongly augmented versions of u_b , or $\mathcal{A}(u_b)$. We show the implications of this difference by comparing the loss functions below.

Loss Function For both methods, the loss function has a labeled or supervised loss term, \mathcal{L}_s , and an unlabeled loss term, \mathcal{L}_u . They are combined into a single training objective as $\mathcal{L} = \mathcal{L}_s + \lambda\mathcal{L}_u$, where λ is a hyperparameter that denotes the relative weight of the unlabeled loss. The labeled loss is similar for both methods. In MixMatch, for each batch $x'_b \in \mathcal{X}'$, the labeled loss is $\mathcal{L}_s = \frac{1}{B} \sum_{b=1}^B H(p_b, p_{\text{model}}(y|x'_b))$, where $H(p, q)$ is defined as the cross-entropy between two probability distributions. In FixMatch, where MixUp is not applied, $\mathcal{L}_s = \frac{1}{B} \sum_{b=1}^B H(p_b, p_{\text{model}}(y|\alpha(x_b)))$. The unlabeled loss term differs more significantly. For MixMatch, it is the squared L_2 loss on predictions and guessed labels as shown by eq. (1) for each batch $u'_b \in \mathcal{U}'$.

$$\mathcal{L}_u = \frac{1}{\mu B} \sum_{b=1}^B \|q_b - p_{\text{model}}(y|u'_b)\|_2^2 \quad (1)$$

Alternatively, FixMatch enforces cross-entropy loss against the strongly augmented unlabeled examples given the weakly augmented ones as shown in eq. (2).

$$\mathcal{L}_u = \frac{1}{\mu B} \sum_{b=1}^B \mathbb{1}(\max(q_b) \geq \tau) H(\arg \max(q_b), p_{\text{model}}(y|\mathcal{A}(u_b))) \quad (2)$$

where τ is a scalar hyperparameter that signifies the threshold above which we retain a pseudo-label. Although Fixmatch’s approach is simpler and empirically performed better, we ran experiments with both methods because of the shift in domain to satellite imagery, where augmentations may have different semantic implications. For example, applying a weak augmentation of a random shift may cause the model to predict a building as damaged because buildings can shift in an earthquake. Alternatively, due to the noisy nature of satellite imagery, strong augmentations could make it more challenging for the model to learn. Thus, we tried both MixMatch and FixMatch in our experiments.

5 Experiments & Results

We test the efficacy of SSL methods on classifying buildings in three different regions: Santa Rosa [15], Haiti, and Aleppo [32], using four different methods: the Twin Tower model that performed best in the past work by Xu et al. [32] as a baseline, fully-supervised learning with only labeled training data, and semi-supervised learning with MixMatch [2] and FixMatch [28].

Setting. We split datasets into 90% and 10% for train and test, respectively. For SSL experiments, we randomly sample a specific number of examples (e.g., 10, 50, 100, 500) evenly from each class as the labeled training set and consider the remainder as the unlabeled set. We conduct experiments 5 times with different splits of the labeled training set. We train a variant of Wide Residual Network (WRN) [33], containing four residual blocks with 32, 64, 128, and 256 filters, respectively. We closely follow the training strategies in [28], using momentum (0.9) SGD with cosine learning rate decay. We tune other hyperparameters, such as the learning rate or weight decay, for each method, but they are shared across datasets. We provide experimental details in the supplementary material.

Results. In Table 1, we report accuracy averaged over 5 runs, each of which is trained on different labeled data splits. We also train a fully supervised model on all labeled data (called "90% supervised") as a performance upper bound, although this would not be available in practice.

Dataset	# Labeled Data	Twin Tower	Fully Supervised	MixMatch	FixMatch
Haiti	10	0.53±0.03	0.58±0.02	0.60±0.11	0.56±0.05
	50	0.56±0.02	0.64±0.01	0.72±0.03	0.61±0.03
	100	0.56±0.05	0.69±0.02	0.75±0.04	0.75±0.10
	500	0.71±0.01	0.75±0.01	0.82±0.01	0.87±0.01
	90% supervised (45,667 data)	0.90	-	-	-
Santa Rosa	10	0.54±0.03	0.74±0.08	0.70±0.06	0.92±0.07
	50	0.54±0.03	0.91±0.03	0.85±0.08	0.96±0.02
	100	0.58±0.04	0.92±0.01	0.88±0.04	0.97±0.01
	500	0.69±0.03	0.96±0.01	0.96±0.01	0.98±0.00
	90% supervised (11,067 data)	0.99	-	-	-
Aleppo	10	0.52±0.04	0.55±0.05	0.65±0.08	0.72±0.15
	50	0.56±0.02	0.65±0.06	0.73±0.06	0.88±0.01
	100	0.57±0.04	0.71±0.06	0.78±0.04	0.89±0.01
	500	0.67±0.04	0.82±0.01	0.85±0.02	0.90±0.01
	90% supervised (9,406 data)	0.88	-	-	-

Table 1: Classification accuracy of the Twin Tower [32] baseline, fully-supervised, MixMatch, and FixMatch models with varying amounts of labeled training data, averaged over 5 runs (each with a different labeled data split). To provide an upper bound, one model is trained using all training data.

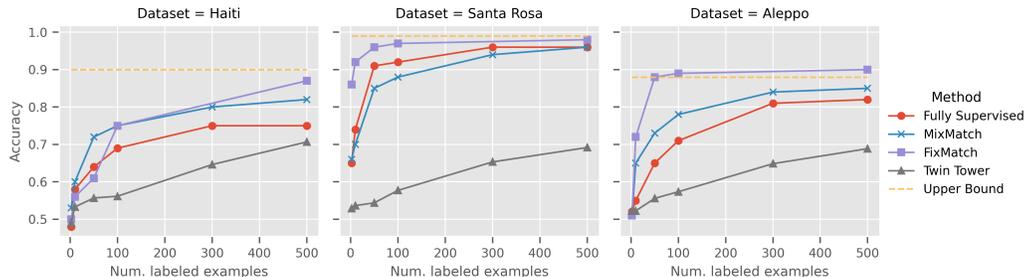


Figure 2: Performance of fully and semi-supervised models by number of labeled training examples.

SSL improves accuracy on all three disaster datasets. MixMatch shows consistent improvements over Twin Tower and fully-supervised models in Haiti and Aleppo, regardless of the number of labeled training data. FixMatch makes further significant improvements in all three regions, although it does require 100 or more labeled data to surpass MixMatch performance in Haiti. In Aleppo, FixMatch with 500 labeled data is able to outperform the fully-supervised model trained on all labeled data, possibly because strong augmentations capture a level of variance not in the dataset alone. We examine augmentation policies in Appendix C. Overall, our results demonstrate the generality of modern SSL methods based on data augmentation, such as MixUp [34], CTAugment [1], or RandAugment [5], in spite of the contrasting image statistics of satellite imagery compared to standard visual recognition benchmarks, such as CIFAR-10 [21] or ImageNet [7].

6 Conclusion

In this paper, we introduced a novel application of semi-supervised learning to automatically detect damaged buildings in satellite imagery with limited labeled data. We experimented with two recent techniques, MixMatch and FixMatch, and showed how they are able to achieve strong performance 100 labeled examples or fewer by leveraging unlabeled data. They consistently outperformed fully supervised models and even achieved performance close to that of a fully supervised setting with no data constraints. The results empirically showed how SSL approaches can be useful to train models when a new disaster is unfolding in an unseen region. For future work, we plan to investigate how to effectively incorporate data from past disasters; there may be region-independent transformations caused by a disaster that the models do not sufficiently capture or different types of augmentations and losses that are more robust to the noise inherent to satellite imagery.

Acknowledgments and Disclosure of Funding

This work is a collaboration between Google Research and the United Nations World Food Programme (WFP) Innovation Accelerator. The WFP Innovation Accelerator identifies, supports and scales high-potential solutions to hunger worldwide. We support WFP innovators and external start-ups and companies through financial support, access to a network of experts and a global field reach. We believe the way forward in the fight against hunger is not necessarily in building grand plans, but identifying and testing solutions in an agile way. The Innovation Accelerator is a space where the world can find out what works and what doesn't in addressing hunger - a place where we can be bold, and fail as well as succeed.

References

- [1] David Berthelot, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and Colin Raffel. Remixmatch: Semi-supervised learning with distribution matching and augmentation anchoring. In *International Conference on Learning Representations*, 2020.
- [2] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin Raffel. Mixmatch: A holistic approach to semi-supervised learning. In *NeurIPS*, 2019.
- [3] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In Geoffrey Gordon, David Dunson, and Miroslav Dudík, editors, *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 215–223, Fort Lauderdale, FL, USA, 11–13 Apr 2011. JMLR Workshop and Conference Proceedings.
- [4] Austin Cooner, Yang Shao, and James Campbell. Detection of urban damage using remote sensing and machine learning algorithms: Revisiting the 2010 haiti earthquake. *Remote Sensing*, 8(10):868, Oct 2016.
- [5] Ekin D. Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V. Le. Randaugment: Practical automated data augmentation with a reduced search space, 2019.
- [6] F. Dell'Acqua and P. Gamba. Remote sensing and earthquake damage assessment: Experiences, limits, and perspectives. *PIEEE*, 100(10):2876–2890, October 2012.
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [8] Terrance Devries and Graham W. Taylor. Improved regularization of convolutional neural networks with cutout. *CoRR*, abs/1708.04552, 2017.
- [9] D. Duarte, F.C. Nex, N. Kerle, and G. Vosselman. *Satellite Image Classification Of Building Damages Using Airborne And Satellite Image Samples In A Deep Learning Approach*, volume IV, pages 89–96. International Society for Photogrammetry and Remote Sensing (ISPRS), 2 edition, 5 2018.
- [10] Victor Durnov. xview2 challenge 1st place, 2020.
- [11] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [12] Noel Gorelick, Matt Hancher, Mike Dixon, Simon Ilyushchenko, David Thau, and Rebecca Moore. Google earth engine: Planetary-scale geospatial analysis for everyone. *Remote Sensing of Environment*, 2017.
- [13] Ritwik Gupta, Bryce Goodman, Nirav Patel, Ricky Hosfelt, Sandra Sajeew, Eric Heim, Jigar Doshi, Keane Lucas, Howie Choset, and Matthew Gaston. Creating xbd: A dataset for assessing building damage from satellite imagery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019.

- [14] Ritwik Gupta, Bryce Goodman, Nirav Patel, Ricky Hosfelt, Sandra Sajeev, Eric Heim, Jigar Doshi, Keane Lucas, Howie Choset, and Matthew Gaston. xview2 dataset, 2019.
- [15] Ritwik Gupta, Richard Hosfelt, Sandra Sajeev, Nirav Patel, Bryce Goodman, Jigar Doshi, Eric Heim, Howie Choset, and Matthew Gaston. xbd: A dataset for assessing building damage from satellite imagery, 2019.
- [16] Rohit Gupta and Mubarak Shah. Rescuenet: Joint building segmentation and damage assessment from satellite imagery, 2020.
- [17] Humanitarian data exchange. <https://data.humdata.org>. Accessed: 2019-09-01.
- [18] Min Ji, Lanfa Liu, and Manfred Buchroithner. Identifying collapsed buildings using post-earthquake satellite imagery and convolutional neural networks: A case study of the 2010 haiti earthquake. *Remote Sensing*, 10(11):1689, Oct 2018.
- [19] Norman Kerle. Satellite-based damage mapping following the 2006 indonesia earthquake—how accurate was it? *International Journal of Applied Earth Observation and Geoinformation*, 12(6):466 – 476, 2010. Geospatial Technologies for Disaster Management.
- [20] Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 05 2012.
- [21] Alex Krizhevsky. Learning multiple layers of features from tiny images. *University of Toronto*, 05 2012.
- [22] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. *CoRR*, abs/1610.02242, 2016.
- [23] G. J. McLachlan. Iterative reclassification procedure for constructing an asymptotically optimal rule of allocation in discriminant analysis. *Journal of the American Statistical Association*, 70(350):365–369, 1975.
- [24] Francesco Nex, Diogo Duarte, Fabio Giulio Tonolo, and Norman Kerle. Structural building damage detection with deep learning: Assessment of a state-of-the-art cnn in operational conditions. *Remote sensing*, 11(23):2765, 2019.
- [25] C. Rosenberg, M. Hebert, and H. Schneiderman. Semi-supervised self-training of object detection models. In *2005 Seventh IEEE Workshops on Applications of Computer Vision (WACV/MOTION'05) - Volume 1*, volume 1, pages 29–36, 2005.
- [26] Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. In *Proceedings of the 30th International Conference on Neural Information Processing Systems, NIPS'16*, page 1171–1179, Red Hook, NY, USA, 2016. Curran Associates Inc.
- [27] H. Scudder. Probability of error of some adaptive pattern-recognition machines. *IEEE Transactions on Information Theory*, 11(3):363–371, 1965.
- [28] Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Ekin D Cubuk, Alex Kurakin, Han Zhang, and Colin Raffel. Fixmatch: Simplifying semi-supervised learning with consistency and confidence. *arXiv preprint arXiv:2001.07685*, 2020.
- [29] Tinka Valentijn, Jacopo Margutti, Marc van den Homberg, and Jorma Laaksonen. Multi-hazard and spatial transferability of a cnn for automated building damage assessment. *Remote Sensing*, 12(17), 2020.
- [30] Ethan Weber and Hassan Kané. Building disaster damage assessment in satellite imagery with multi-temporal fusion, 2020.
- [31] Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V. Le. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

- [32] Joseph Z. Xu, Wenhan Lu, Zebo Li, Pranav Khaitan, and Valeriya Zaytseva. Building damage detection in satellite imagery using convolutional neural networks. *CoRR*, abs/1910.06444, 2019.
- [33] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- [34] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018.

Appendix A Visualization of Satellite Imagery

In this work, we consider evaluating on three datasets: Santa Rosa [15], Haiti and Aleppo [32]. In this section, we provide visualization of full and cropped satellite imagery for each dataset.

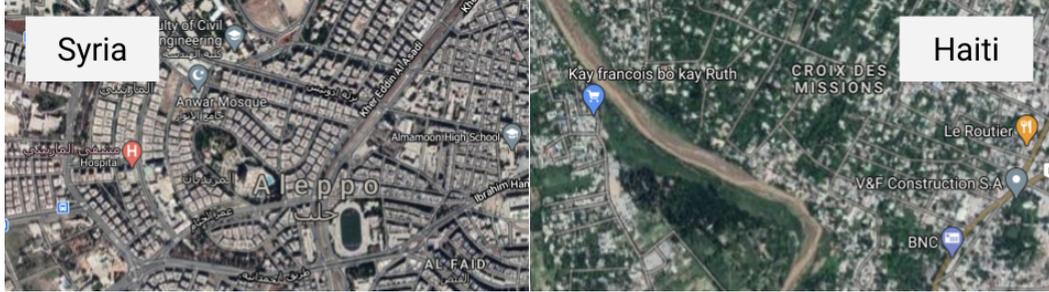


Figure 3: Comparison of imagery in two different regions, Syria and Haiti. Building layout, construction material, level of vegetation, appearance of damage (based on disaster type and buildings), and other factors make it difficult for models to generalize to unseen regions.



Figure 4: Example of images cropped around buildings from Haiti.

Appendix B Details of Experiments

We provide additional details of experiments for reproducible research.

B.1 Hyperparameters

For the fully supervised and MixMatch experiments, we used a batch size of 64, learning rate of 0.002, and weight decay of 0.002. We used a Beta distribution of $Beta(0.5)$ for mixup. For FixMatch, we

RandAugment Transformation	Cutout	Accuracy	Haiti	SoCal	Syria
none	yes	mean	0.60	0.97	0.83
		stdev	0.06	0.01	0.03
color	no	mean	0.47	0.98	0.61
		stdev	0.07	0.01	0.18
color	yes	mean	0.66	0.92	0.87
		stdev	0.06	0.02	0.02
geo	no	mean	0.56	0.97	0.73
		stdev	0.01	0.01	0.15
geo	yes	mean	0.68	0.97	0.87
		stdev	0.05	0.01	0.02
color + geo	no	mean	0.63	0.98	0.82
		stdev	0.02	0.01	0.07
color + geo	yes	mean	0.75	0.97	0.89
		stdev	0.10	0.01	0.01

Table 2: Ablation study of applying different transformations of RandAugment and Cutout. We report the average and standard deviation over 5 runs per setting.

used a batch size of 64, learning rate of 0.03, and weight decay of 0.0005. Additionally, the ratio for unlabeled data $\mu = 3$, pseudo-label loss weight $\lambda = 1$, and the confidence threshold $\tau = 0.95$.

For MixMatch and FixMatch experiments, rather than decaying the learning rate, we evaluate models using an exponential moving average of their parameters with a decay rate of 0.999.

Appendix C Ablation Study on Augmentation Policy

We used the RandAugment [5] augmentation policy for FixMatch, which randomly selects transformations for each sample in a minibatch from a collection of transformations (e.g. color inversion and geometric changes, such as rotation and translation). We also apply Cutout [8], which sets a random square patch of pixels to gray. We conduct an ablation study to examine which operations play a key role in FixMatch, testing 7 settings:

1. Cutout only
2. RandAugment only with color transformations
3. RandAugment and Cutout with color transformations
4. RandAugment only with geometric transformations
5. RandAugment and Cutout with geometric transformations
6. RandAugment only with color and geometric transformations
7. RandAugment and Cutout with color and geometric transformations

To clarify, the setting used in the Results section is RandAugment and Cutout with color and geometric transformations. We use 50 labeled and 50 unlabeled examples in training for all experiments. We run each setting 5 times and report the average and standard deviation of classification performance in Table 2. Generally, Cutout seemed to play a more significant role than RandAugment, similarly to the finding in the original FixMatch paper where Cutout enabled strong results on CIFAR-10 [28]. Between geometric and color transformations, the former seemed more influential, though the difference becomes marginal when Cutout is applied.