

1                   **Predicting tree failure likelihood for utility risk mitigation via artificial**  
2                   **intelligence**

3                   Jimi Oke<sup>1</sup>, Nasko Apostolov<sup>1</sup>, Ryan Suttle<sup>2</sup>, Sanjay Arwade<sup>1</sup>, and Brian Kane<sup>2</sup>

4                   <sup>1</sup>Department of Civil and Environmental Engineering, University of Massachusetts Amherst, MA  
5                   01003, USA

6                   <sup>2</sup>Department of Environmental Conservation, University of Massachusetts Amherst, MA 01003,  
7                   USA

8                   **ABSTRACT**

9                   Critical to the resilience of utility power lines, tree failure assessments have historically been  
10                  performed via manual and visual inspections. In this paper, we develop a convolutional neural  
11                  network (CNN) to predict tree failure likelihood categories under four classification scenarios.  
12                  Starting with an original set of 525 expert-labeled images, we perform preprocessing and aug-  
13                  mentation tasks to increase the number of samples. We optimize several hyperparameters in our  
14                  CNN and test its performance for three different image resolutions. The trained classifier has an  
15                  average validation accuracy of at least 0.94 across all scenario-resolution combinations. Thus,  
16                  via this novel framework, we demonstrate the potential of artificial intelligence to automate and  
17                  consequently reduce the costs of tree failure likelihood assessments, thereby promoting sustainable  
18                  infrastructure.

19                   **INTRODUCTION**

20                  Despite extensive efforts by utilities to prevent them, contacts between tree parts and power  
21                  lines cause outages that annually result in tens of billions of dollars in economic costs throughout  
22                  the United States. Presently, the identification of potential contact between trees and power lines  
23                  is labor intensive and time-consuming. This paper describes an artificial intelligence and machine  
24                  learning approach that automatically classifies trees, using only a single photograph and with a high  
25                  degree of accuracy, into categories used by utility arborists to describe the likelihood of tree failure:  
26                  probable, possible, and improbable. This preliminary study demonstrates the possible efficacy of  
27                  AI approaches to tree risk assessment and, following further development of the approach, has the  
28                  potential to reduce power outages and utility costs by allowing utilities to more effectively target  
29                  their pruning and mitigation efforts.

30                  Contact between tree parts and power lines can take three forms: tree branches can grow into  
31                  lines; branches can fail and fall onto lines; whole-tree failure can occur due to uprooting or trunk  
32                  failure. A study in Connecticut, USA provides some context for the amount of economic disruption,  
33                  documenting annual disruptions of \$8.3 billion between 2005 and 2015 ([Graziano et al. 2020](#)). That  
34                  extremely high cost occurred despite extensive efforts on the part of utilities to mitigate conflicts  
35                  between trees and power lines through active and aggressive pruning programs that, on their own  
36                  cost billions of dollars annually ([Guggenmoos 2003](#)).

37                  Despite its high cost, pruning trees to maintain clearance from power lines is an effective way  
38                  to reduce outages due to so-called “preventable” contacts between trees and power lines. For

example, in Massachusetts, USA, where tree failure was responsible for 40% of preventable tree-caused outages, pruning was able to improve reliability by 20% to 30% (Simpson and Van Bossuyt 1996); similar results were found in a study conducted in Connecticut (Parent et al. 2019). The efficacy of pruning has also been shown in a study of two states in the Gulf Coast region of the USA that showed wind-induced power outage prediction models becoming less uncertain when pruning was included in the model (Nateghi et al. 2014).

Even effective pruning cannot, however, completely eliminate tree-caused outages. Failure of trees outside the right-of-way can still impact the lines and cause outages (Guggenmoos 2003). The proportion of outages caused by failure of trees outside the right-of-way has not been rigorously quantified. Guggenmoos (2011) estimated that 95% of tree-caused outages in the Pacific Northwest region of the USA, were due to tree failure, and Wismer (2018) reported approximately 25% of interruptions in Illinois, USA, were caused by trees that uprooted or broke in the stem.

Predicting the likelihood of failure is an inexact science, but tree risk assessment best management practices have been developed (Smiley et al. 2017; Goodfellow 2020). Estimating tree risk includes assessing the likelihood of tree failure, the likelihood of impact of the failed tree (or tree part) on a target, and the severity of consequences of the impact. The likelihood of failure depends on the anticipated loads on the tree and its load-bearing capacity. The likelihood of impact depends on proximity to the target (the lines, poles, and other hardware—“infrastructure”—in the case of utility tree risk assessment), the target’s occupancy rate (which is constant for utility lines) and whether the target is sheltered, for example by neighboring trees. Severity of consequences depends on the damage done to the infrastructure—which, in turn, is partially related to the size of the tree or tree part that fails, and how much momentum it has when it impacts the infrastructure—and, more importantly in some cases, the economic costs and disruption associated with electrical outages.

Individual tree risk assessment can be costly because of the time it requires. In some situations, a less time-consuming assessment may be justified to reduce costs, i.e. a “Level 1” assessment (Smiley et al. 2017). Studies in Rhode Island, USA (Rooney et al. 2005) and Florida, USA (Koeser et al. 2016) have shown that, compared to more time-consuming risk assessments, Level 1 risk assessments successfully identified trees with a higher degree of risk—precisely the trees that arborists prioritize for risk mitigation. The utility of Level 1 assessments demonstrated in these studies suggests that artificial intelligence (AI) tools may be an effective way to reduce the cost of tree risk assessment while still identifying high risk trees.

The method described in the paper uses convolutional neural networks (CNN) to classify images of trees among three categories of failure likelihood: probable, possible, and improbable. The data used for training, testing and illustration of the method consists of 505 tree images that have been classified by the authors according to best management practices used by utility arborists (Goodfellow 2020).

The remainder of the paper provides a brief history and background of AI and its use in infrastructure risk assessment and tree identification (section 2); describes the methods used to train and validate a novel CNN to categorize likelihood of tree failure (section 3); and presents and discusses the output of the novel CNN (sections 4 and 5). The goal is to further demonstrate an innovative automated approach to tree risk assessment using an AI tool that can be readily deployed for use in various locations and also continually improved through subsequent training on new datasets.

82 **BACKGROUND**

83 AI-based image analysis is relatively widely used, even in engineering applications, such as  
84 earthquake risk assessment ([Jiao and Alavi 2020](#); [Salehi and Burgueño 2018](#)) and structural health  
85 monitoring ([Spencer et al. 2019](#); [Wang et al. 2019](#)). Neural networks, which comprise a major  
86 category of AI frameworks, have been widely applied in the field of earthquake risk assessment  
87 (an excellent review is provided by [Xie et al. \(2020\)](#)), but the authors are not aware of attempts  
88 to operate directly on, for example, building images in the absence of technical structural data to  
89 predict seismic risk. Neural networks have also been used to interrogate remote sensing data of  
90 the landscape to assess landslide risk ([Su et al. 2020](#)). A few recent efforts have demonstrated the  
91 potential for AI-based tree recognition from drone imagery ([dos Santos et al. 2019](#); ?). Furthermore,  
92 an application of a convolutional neural network (CNN) to tree species identification was  
93 recently demonstrated by [Fricke et al. \(2019\)](#). Yet, AI has yet to be applied to the problem of  
94 tree-utility line risk assessment—one that is complicated by the very large number of tree species  
95 to be considered, seasonal variation in tree appearance and associated risk and local meteorological  
96 conditions.

97 The groundbreaking study of [Hubel and Wiesel \(1959\)](#) showed that visual perception in cats  
98 was a result of the activation or inhibition of groups of cells in the visual cortex known as “receptive  
99 fields.” Further, they attempted to map the cortical architecture in cats and monkeys ([Hubel and  
100 Wiesel 1962](#); [Hubel and Wiesel 1965](#); [Hubel and Wiesel 1968](#)). Subsequent attempts were then  
101 made to model neural networks that could be trained to automatically recognize visual patterns with  
102 modest performance ([Rosenblatt 1962](#); [Kabrisky 1966](#); [Giebel 1971](#); [Fukushima 1975](#)). However,  
103 the breakthrough came with the “neocognitron” ([Fukushima 1980](#)), which was a self-learning  
104 neural network for pattern recognition that was robust to changes in position and shape distortion, a  
105 problem that plagued earlier efforts, including the “cognitron” also proposed by [Fukushima \(1975\)](#).

106 A few notable efforts demonstrated the neural networks for handwritten digit recognition  
107 ([Fukushima 1988](#); [Denker et al. 1988](#)), but these required significant preprocessing and feature  
108 extraction. [LeCun et al. \(1989\)](#) soon afterward introduced a multilayer neural network that mapped  
109 a feature in each neuron (representing a “local receptive field”) via convolution. This network could  
110 also be trained by backpropagation like other existing neural networks and featured pooling operations  
111 for better distortion and translation invariance. Further developments from this milestone  
112 yielded the LeNet-5 convolutional neural network which attained accuracy levels that rendered it  
113 commercially viable.

114 The big data revolution coupled with technological advancements that have made it possible to  
115 capture and store high resolution images have raised challenges that continue to be surmounted with  
116 successively high-performing architectures. Over the past decade, some of these efforts resulted in  
117 significant breakthroughs in performance. AlexNet ([Krizhevsky et al. 2012](#)), with 5 convolutional  
118 layers and 3 dense layers—one of the largest CNNs of its time, won the ILSVRC-2012<sup>1</sup> competition  
119 with a top-5 error rate of 15.3% and served as a landmark in the Deep Learning subdomain. [Zeiler  
120 and Fergus \(2014\)](#) then introduced ZFNet, besting the performance of AlexNet, and pioneered  
121 visualization techniques that were foundational for model inference and interpretability. In the  
122 same year, GoogLeNet, a 22-layer network, was proposed ([Szegedy et al. 2014](#)), featuring the  
123 novel “Inception module,” which allowed for efficiency and accuracy in a very deep network.  
124 Subsequent improvements have been proposed to the original inception framework ([Szegedy et al.](#)

---

<sup>1</sup>ImageNet Large Scale Visual Recognition Challenge; held annually from 2010 through 2017.

125        2015; Szegedy et al. 2016). VGGNet (Simonyan and Zisserman 2015) also pushed the boundaries  
126        of depth with up 19 layers, achieving state-of-the-art performance at ILSVRC-2014. Finally, ResNet  
127        (He et al. 2015) addressed the accuracy degradation problem that arises with increasing depth in a  
128        network by successively fitting smaller sets of layers to the residual and employing skip connections.  
129        With these innovations, an unprecedented level of depth was achieved. Implementations with with  
130        34, 50, 101 and 152 layers were demonstrated. ResNet-152 won first place in ILSVRC-2015.

131        Along with these developments in their architectures, CNNs have demonstrated viability for  
132        applications ranging from image classification, object and text detection to document tracking,  
133        labeling, speech, among several other related fields (Gu et al. 2018). In this study, we show that  
134        a relatively simple CNN architecture coupled with state-of-the-art approaches for model training  
135        and regularization is capable of efficiently and effectively predicting tree failure classes.

## 136 DATA AND METHODS

### 137 Image data description

138        The training dataset consisted of 505 images, each having an original size of  $4032 \times 3024$  pixels.  
139        Images were captured over a single field season in Massachusetts, USA, between May and  
140        September 2020 to limit any potential influence of changes in tree appearance due to seasonal leaf  
141        senescence on image processing. ESRI ArcMaps was used to randomly distribute sampling sites  
142        across the state. Field assessments of trees to classify likelihood of failure followed the “Level  
143        1” methods outlined in the second edition of the International Society of Arboriculture’s (ISA)  
144        Tree Risk Assessment Best Management Practices (Smiley et al. 2017) and ISA’s Utility Tree Risk  
145        Assessment Best Management Practices (Goodfellow 2020). This method is commonly used to  
146        assess trees in the United States. A Level 1 assessment was selected for this study because: (1)  
147        individual risk assessments may be prohibitively expensive at higher orders, i.e. Level 2 or Level 3  
148        (Smiley et al. 2017), given the hundreds of thousands of trees utilities must manage across territory  
149        areas; (2) utility right-of-way (ROW) easements may not allow utility inspectors full access to trees  
150        in practical application of higher order risk assessment procedure if the trees are beyond the edge  
151        of the ROW (Goodfellow 2020); and (3) studies have shown reasonable efficacy of limited basic  
152        visual assessment techniques in identifying more severe tree defects (Rooney et al. 2005; Koeser  
153        et al. 2016) leading to greater likelihood of failure ratings. The four categories of likelihood of tree  
154        failure, which are always considered in a stated time frame, are defined as follows (Smiley et al.  
155        2017):

- 156        • **Improbable:** failure unlikely either during normal or extreme weather conditions;
- 157        • **Possible:** failure expected under extreme weather conditions; but unlikely during normal  
158        weather conditions;
- 159        • **Probable:** failure expected under normal weather conditions within a given time frame;
- 160        • **Imminent:** failure has started or is most likely to occur in the near future, even if there is no  
161        significant wind or increased load. This is a rare occurrence for a risk assessor to encounter,  
162        and may require immediate action to protect targets from impact.

163        In this study, only images of trees assigned to the 3 likelihood of failure categories of *Improbable*,  
164        *Possible*, and *Probable* were included in modeling, while images of *Imminent* trees were excluded  
165        due to their rarity. Typical examples are shown for each category in Figure 1. In the original set of  
166        training images, the class distribution is given in Table 1.



(a) Probable

(b) Possible

(c) Improbable

**Fig. 1.** Examples of training images in each of the three tree risk categories considered in this study: Probable, Possible and Improbable. Two are shown in each case.

Category	Number of images
Improbable	322
Possible	80
Probable	56
Total	505

**TABLE 1.** Category distribution of images in the original set of input images

### Classification scenarios

In order to investigate the efficacy of an AI classifier to distinguish the failure-liability categories, we defined four classification scenarios in Table 2 for our experiments. Each scenario represents a unique grouping of each of the three categories, with a minimum of two derived classes in each case. Thus,

### Image pre-processing and augmentation

Data augmentation refers to the variety of methods that are employed for synthetically generating more samples in a training dataset in order to improve model performance (Wong et al. 2016).

Scenario	Description	No. classes
Pr_Po_Im	{Probable, Possible, Improbable}	3
Pr_Im	{Probable, Improbable}	2
PrPo_Im	{Probable + Possible, Improbable}	2
Pr_PoIm	{Probable, Possible + Improbable}	2

TABLE 2. Classification scenarios

Augmentation is desired, particularly in situations where the number of original observations is small, and the effectiveness of various relevant techniques in this domain has been amply demonstrated (Shorten and Khoshgoftaar 2019).

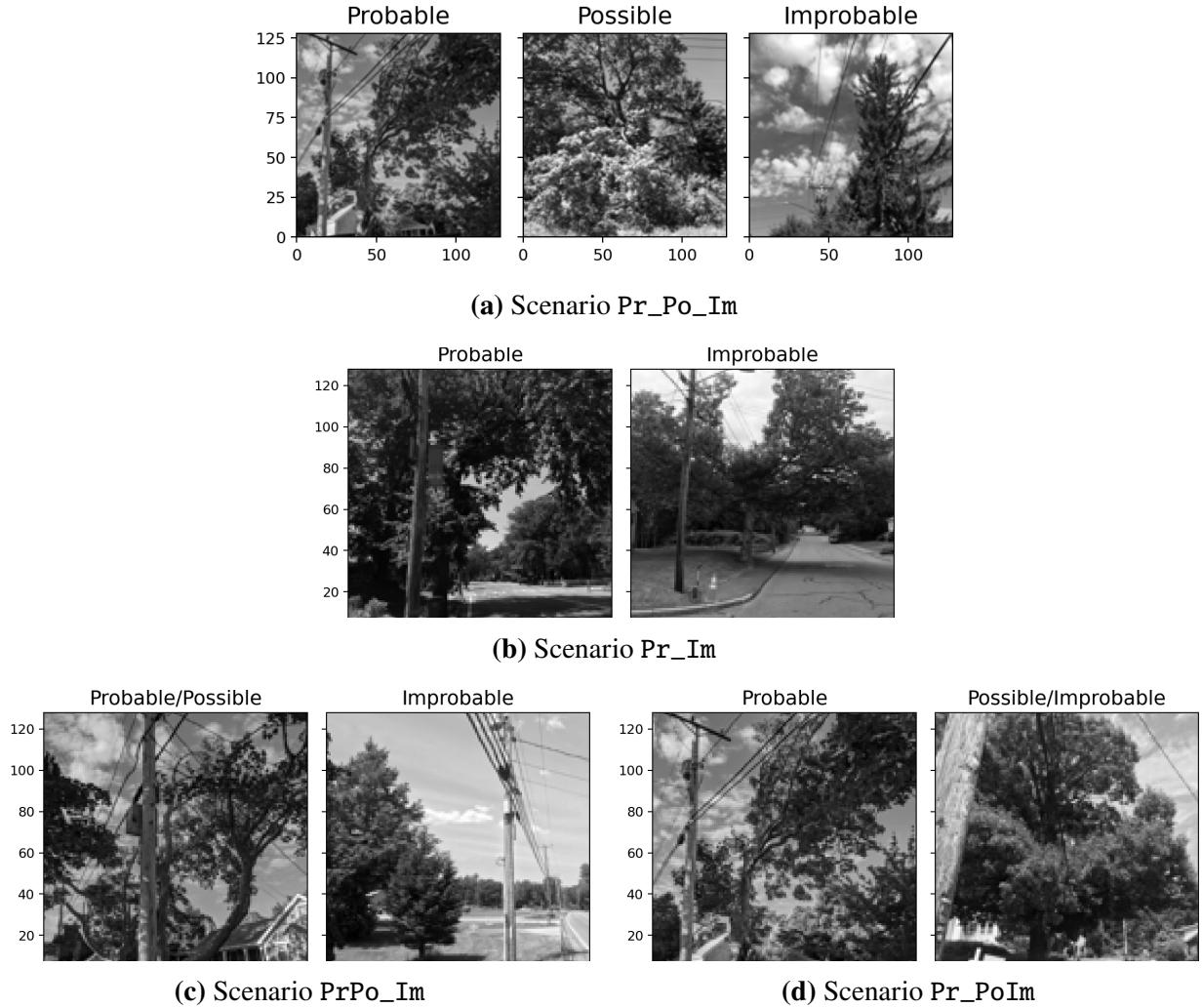
In order to achieve robustness in our model, and given the relatively small number of training images, we randomly cropped each image on either axis to  $3024 \times 3024$  pixels, generating five instances for each one. Thus, we increased the size of our training set from 505 to 2525 images. Further, we performed horizontal flipping with a 50% probability on each of the generated images. For efficiency, we converted the images to grayscale and scaled the pixel values from 0 to 1. Finally, we downsampled the images to the following resolutions (pixels):  $64 \times 64$ ,  $128 \times 128$  and  $224 \times 224$ , creating a training set for each case. Random sets of images from each class across each of the four classification scenarios are shown in Figure 2.

### Convolutional neural network

We employ a convolutional neural network (CNN) as the AI framework for tree risk failure likelihood prediction. Like other neural networks, the CNN is an arrangement of neurons within layers, each neuron performing an operation that maps from a pixel in an input image to the final output. The input into each neuron is a weighted sum from the previous layer, while the output from each neuron is modulated by an activation function. The activation function in the final layer of an CNN is typically the softmax function, which gives the class probabilities of a given input image. A class is assigned to the image based on the one with the corresponding maximum probability.

Unlike other neural networks, however, the CNN performs fundamental pixel-mapping operations in its convolutional layers. Each convolutional layer is defined by a stack of feature maps, which result from a dot product of a filter and correspondingly-sized local receptive fields from the input image or preceding layer. The numeric values of the filters correspond to weights whose optimal values are learned during the training of the CNN. The size of the filter in each convolutional layer is given by the *kernel size*.

In this paper, we employ a relatively simple CNN architecture, historically inspired by AlexNet (Krizhevsky et al. 2012). The structure of the CNN is shown in Figure 3 (generated using an automated framework (Bäuerle et al. 2021)). After the input layer (a matrix of pixels from the input image), we use a 64-filter convolutional layer. The output is downsampled using a pooling layer that returns the maximum output from a  $2 \times 2$  subsample from the previous layer. Next, we stack two successive convolutional layers each with a depth of 128 filters. We follow these with a maximum-pooling layer and then two further 256-filter convolutional layers. A final maximum-pooling layer is used before we "flatten" all outputs into a one-dimensional (fully-connected) array of neurons. After flattening the outputs, we use a dense layer to reduce the number of outputs a specified number of units. Batch normalization (Ioffe and Szegedy 2015) is employed after the first



**Fig. 2.** Random examples of the processed training images under each classification scenario; 128 pixels.

dense layer to improve training efficiency. A second dense layer is used prior to the output layer, with the number of outputs corresponding to the number of classes in the dataset. A regularization technique known as "dropout" is used after each dense layer. In each dropout layer, a proportion of the neurons are randomly zeroed during training in order to improve the robustness of the model.

The various hyperparameters in the model are summarized in Table 3.

#### Hyperparameter optimization

We optimized eight of the CNN hyperparameters using Hyperband (Li et al. 2018), an efficient guided grid-search algorithm. Twelve searches were performed for each classification scenario and image resolution combination. Each search was conducted using 90 trials of unique hyperparameter combinations. The specified range of each parameter along with the search results are shown in Table 4. For the kernel size in the first convolutional layer, we allowed for a choice between a  $5 \times 5$  and a  $7 \times 7$  kernel. The activation function in both dense layers was specified as a choice



**Fig. 3.** Diagram of convolutional neural network structure. Hyperparameters that are optimized include the number of units in the penultimate dense layers.

Layer	Layer #	No. Filters	Kernel Size	Strides	Activation	Rate	No. Units
Convolutional	1	64	$k^*$	2	ReLU		
Max. Pooling	1		2				
Convolutional	2	128	3	1	ReLU		
Convolutional	3	128	3	1	ReLU		
Max. Pooling	2		2				
Convolutional	4	256	3	1	ReLU		
Convolutional	5	256	3	1	ReLU		
Max. Pooling	3		2				
Flatten							
Dense	1				$a_1^*$		$u_1^*$
Dropout	1					$r_1^*$	
Dense	2				$a_2^*$		$u_2^*$
Dropout	2					$r_2^*$	

**TABLE 3.** Summary of the convolutional neural network hyperparameters. Those indicated by an asterisked symbol are optimized using a guided search.

between the rectified linear unit (ReLU) function and the hyperbolic tangent (tanh). The ReLU was introduced to address the so-called "vanishing gradient" problem and has been shown to improve performance in CNNs (Glorot et al. 2011). Nevertheless, the tanh function remains a viable option, as well. The dropout rates were allowed to range from 0 to 5 in steps of 0.05, while the number of neurons or units in each dense layer varied from 32 to 512 in steps of 32. Finally, we uniformly sampled learning rates  $\lambda$  for the optimizer in the  $\log_{10}$  space of  $[10^{-4}, 10^{-2}]$ .

### Model training and assessment

The softmax activation function  $f(\cdot)$  in the output layer returns the class prediction probabilities for a given observation  $i$ . It is defined in terms of the class-specific score  $s_c$ :

$$f(s_c) = \frac{e^{s_c}}{\sum_{c' \in C} e^{s_{c'}}} \quad (1)$$

Scenario	Hyperparameter	Range	Resolution		
			64	128	224
Pr_Im	1st conv. kernel size, $k$	{5, 7}	7	7	7
	1st dense activation, $a_1$	{ReLU, tanh}	ReLU	ReLU	ReLU
	2nd dense activation, $a_2$	{ReLU, tanh}	ReLU	ReLU	ReLU
	1st dropout rate, $r_1$	{0, .05, ..., 5}	0.1	0.1	0.1
	2nd dropout rate, $r_2$	{0, .05, ..., 5}	0.3	0.3	0.3
	1st dense layer units, $u_1$	{32, 64, ..., 512}	480	480	480
	2nd dense layer units, $u_2$	{32, 64, ..., 512}	448	448	448
	learning rate, $\lambda$	[ $10^{-4}, 10^{-2}$ ]	$1.03 \cdot 10^{-4}$	$1.03 \cdot 10^{-4}$	$1.03 \cdot 10^{-4}$
PrPo_Im	1st conv. kernel size, $k$	{5, 7}	7	5	7
	1st dense activation, $a_1$	{ReLU, tanh}	ReLU	tanh	ReLU
	2nd dense activation, $a_2$	{ReLU, tanh}	ReLU	ReLU	ReLU
	1st dropout rate, $r_1$	{0, .05, ..., 5}	0.1	0.25	0.1
	2nd dropout rate, $r_2$	{0, .05, ..., 5}	0.3	0.35	0.3
	1st dense layer units, $u_1$	{32, 64, ..., 512}	480	384	480
	2nd dense layer units, $u_2$	{32, 64, ..., 512}	448	256	448
	learning rate, $\lambda$	[ $10^{-4}, 10^{-2}$ ]	$1.03 \cdot 10^{-4}$	$1.09 \cdot 10^{-4}$	$1.03 \cdot 10^{-4}$
Pr_PoIm	1st conv. kernel size, $k$	{5, 7}	7	7	7
	1st dense activation, $a_1$	{ReLU, tanh}	ReLU	ReLU	ReLU
	2nd dense activation, $a_2$	{ReLU, tanh}	tanh	ReLU	tanh
	1st dropout rate, $r_1$	{0, .05, ..., 5}	0.25	0.1	0.2
	2nd dropout rate, $r_2$	{0, .05, ..., 5}	0.3	0.3	0.15
	1st dense layer units, $u_1$	{32, 64, ..., 512}	128	480	416
	2nd dense layer units, $u_2$	{32, 64, ..., 512}	320	448	416
	learning rate, $\lambda$	[ $10^{-4}, 10^{-2}$ ]	$1.76 \cdot 10^{-4}$	$1.03 \cdot 10^{-4}$	$1.15 \cdot 10^{-4}$
Pr_Po_Im	1st conv. kernel size, $k$	{5, 7}	7	5	7
	1st dense activation, $a_1$	{ReLU, tanh}	ReLU	tanh	ReLU
	2nd dense activation, $a_2$	{ReLU, tanh}	tanh	ReLU	ReLU
	1st dropout rate, $r_1$	{0, .05, ..., 5}	0.25	0.25	0.1
	2nd dropout rate, $r_2$	{0, .05, ..., 5}	0.3	0.35	0.3
	1st dense layer units, $u_1$	{32, 64, ..., 512}	128	384	480
	2nd dense layer units, $u_2$	{32, 64, ..., 512}	320	256	448
	learning rate, $\lambda$	[ $10^{-4}, 10^{-2}$ ]	$1.76 \cdot 10^{-4}$	$1.09 \cdot 10^{-4}$	$1.03 \cdot 10^{-4}$

**TABLE 4.** Optimal hyperparameters found using the Hyperband search algorithm for 12 classification scenario and input image resolution combinations.

Thus for the  $i$ th observation, the softmax activation returns the predicted probability  $p_{i,c}$  that the  $i$ th observation belongs to class  $c$ . The CNN is trained using a variant of the stochastic gradient algorithm, Adam (Kingma and Ba 2017). The goal of the training procedure is to learn the optimal weights and bias terms for the CNN by minimizing a loss function. In this case, we use the

232  
233  
234  
235

236 categorical cross-entropy loss function, which for a single observation can be simply defined as:

$$237 \quad L_i^{CE} = -\log(p_{i,c}) = -\log(f(s_c^i)) \quad (2)$$

238 Training is iteratively performed, with gradient of the loss function computed and averaged over a  
 239 batch of input images. Here, we use a batch size of 32. The learning rate  $\lambda$  of the optimization  
 240 algorithm is an important hyperparameter that affects training performance. We optimized for this  
 241 in the hyperparameter search as discussed. Furthermore, the CNN is trained over multiple passes  
 242 through the entire training set. Each such pass is referred to as an epoch.

In real terms, we measured the performance of the trained CNN by how accurately it predicts the classes in a validation set excluded from the training set. For this paper, we used a randomly sampled validation that was 20% of the size of the input dataset of 2525 images in each training instance. Thus, we define the accuracy as the overall proportion of correct predictions across all classes. This metric was computed both for the training and validation sets in each epoch. In addition to overall accuracy of making correct classifications, we assessed the models trained under these scenarios based on the macro-averages of the precision, recall and  $F_1$  score metrics computed over the validation set in each epoch. The precision score  $Pr$  captures the proportion of correct predictions for a certain class relative to all the predictions for that class, and is an important measure of how good a classifier is. The recall  $Re$  captures the ability of a classifier to correctly predict observations for a certain class relative to all the true observations in that class. The  $F_1$  metric is given as the harmonic mean of the precision and recall, and is thus more sensitive than the overall accuracy score. These three metrics are macro-averaged. Thus, each category is given equal weight, ensuring that misclassifications within the smaller classes ("Probable" and "Possible") are adequately represented in the aggregation. These metrics are formally defined as follows:

$$\text{Macro-average precision: } Pr^m = \frac{1}{|C|} \sum_{c \in C} \left( \frac{TP_c}{TP_c + FP_c} \right) \quad (3)$$

$$\text{Macro-average recall: } Re^m = \frac{1}{|C|} \sum_{c \in C} \left( \frac{TP_c}{TP_c + FN_c} \right) \quad (4)$$

$$\text{Macro-average } F_1 \text{ score: } F_1^m = \frac{1}{|C|} \sum_{c \in C} \left( \frac{2Pr_c Re_c}{Pr_c + Re_c} \right) \quad (5)$$

where  $c$  is the index of a class in the set  $C$  and  $|C|$  the number of classes in the dataset. The class-specific prediction metrics are given by:

$$\text{True positives for class } c: TP_c = \sum_{i \in c} I(\hat{y}_i = y_i) \quad (6)$$

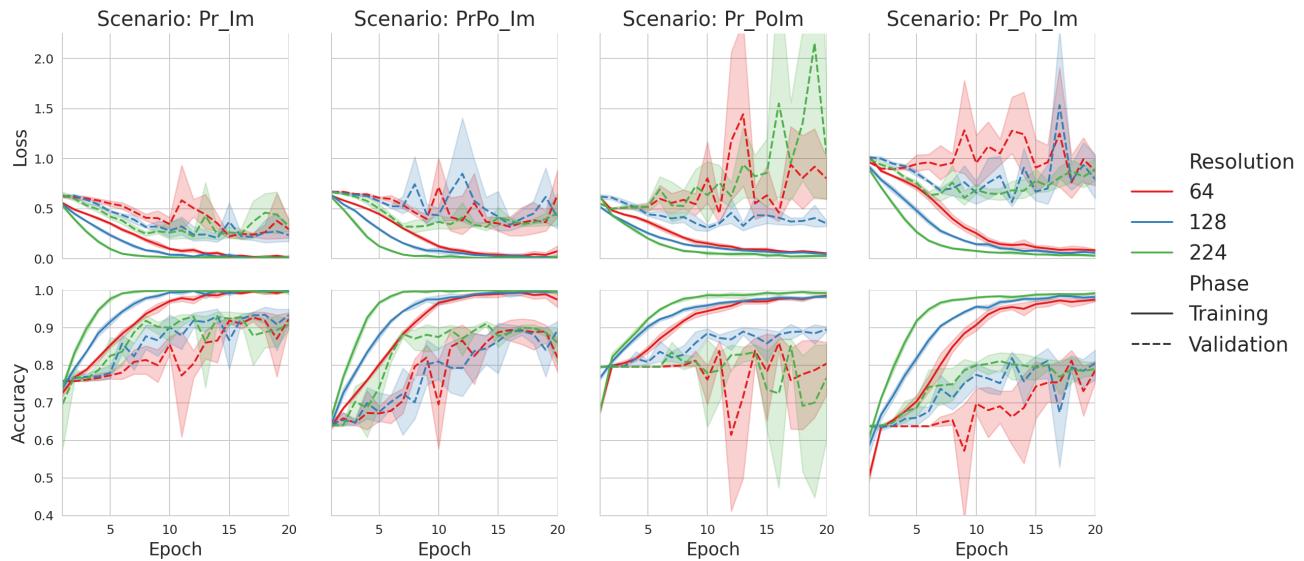
$$\text{False positives for class } c: FP_c = \sum_{i \in c} I(\hat{y}_i \in c | y_i \notin c) \quad (7)$$

$$\text{False negatives for class } c: FN_c = \sum_{i \in c} I(\hat{y}_i \notin c | y_i \in c) \quad (8)$$

243 where  $\hat{y}_i$  is the predicted class and  $y_i$  the observed (true) class for a given image  $i$  in class  $c$ . The  
 244 indicator function  $I(\cdot)$  returns 1 if the corresponding condition is true, and 0 otherwise.

245 **RESULTS**

246 We trained the CNN for each of the 4 classification scenarios. In each scenario, we also trained  
 247 on three image resolution sets (64, 108 and 224). The goal was to determine the best performing  
 248 resolution, given the tradeoff between performance and computational expenditure. Thus, we  
 249 generated twelve learning cases. In each case, we performed a 5-fold cross-validation (resulting in  
 250 a training-validation ratio of 80:20 in each fold). We trained the CNN over 20 epochs in all the  
 251 cases. The trajectories of the loss and accuracy are shown in Figure 4 for both the training and  
 252 validation sets.



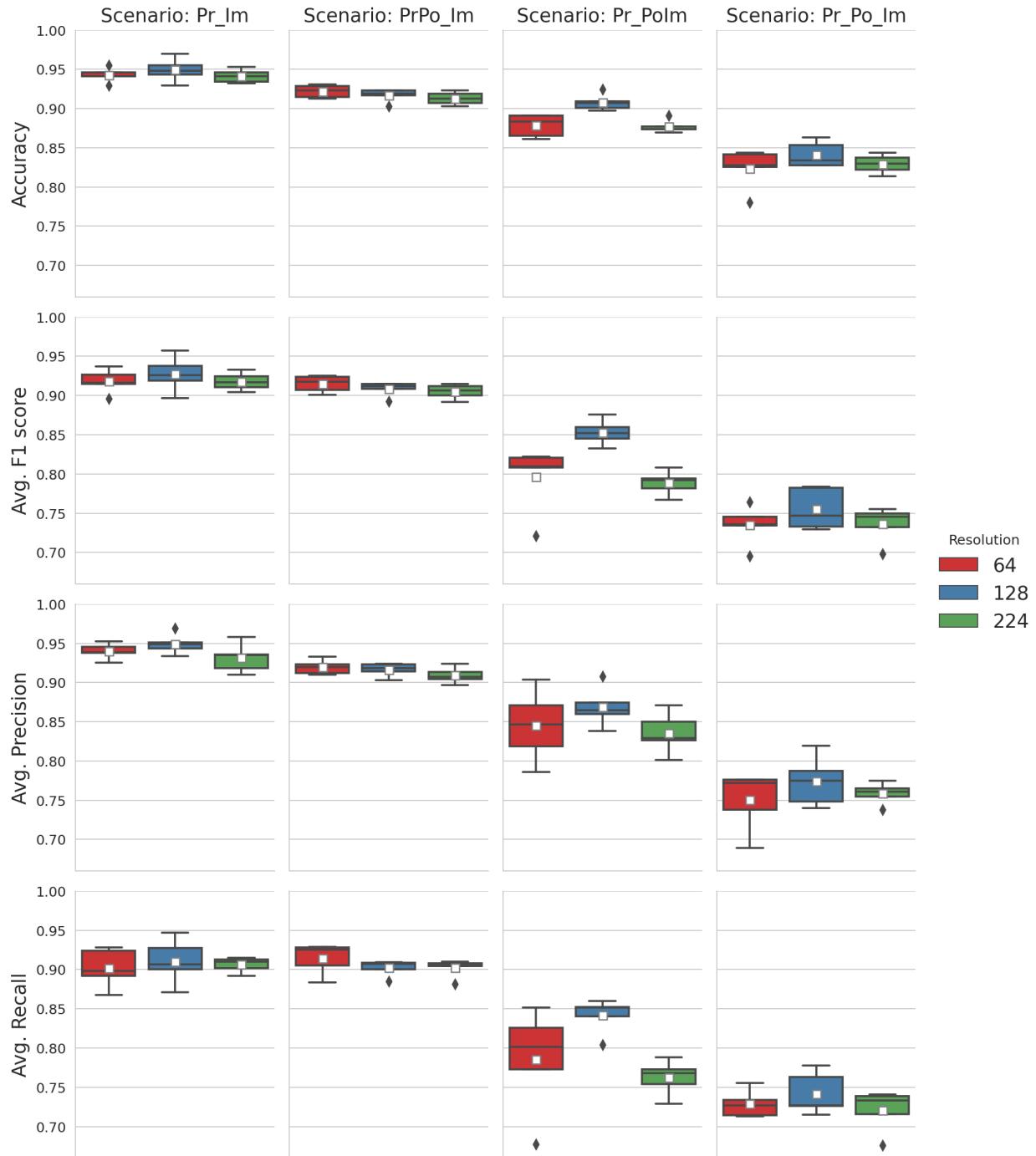
**Fig. 4.** Performance metrics for each scenario and input resolution instance. Average metrics and 95% confidence intervals (trials,  $n = 5$ ) are shown for 20 epochs in each case.

253 **Sensitivity to training resolution**

254 Figure 5 shows the boxplots of the validation performance metrics for all twelve cases con-  
 255 sidered. The metrics are: accuracy and the macro-averages of precision, recall and the  $F_1$  score.  
 256 Taking all metrics into consideration, Welch pairwise tests showed that there are no significant  
 257 differences in performance among the 3 training resolutions, with one exception. In the scenario  
 258 Pr\_PoIm, the accuracy,  $F_1^m$ ,  $Re^m$  for the 128-pixel case are greater than those for the 224-pixel case  
 259 ( $p < .001$ ). This difference in performance is not as stark between the 128-pixel and 64-pixel cases  
 260 in the same scenario. This outcome implies that we can achieve efficiency by training at lower  
 261 resolutions without significant losses in performance.

262 **Scenario sensitivity**

263 We conducted pairwise tests between each scenario combination for all resolutions. The results  
 264 indicated that considering all metrics, the scenarios are statistically significantly different from each  
 265 other ( $p < .005$ ) except Pr\_Im compared to PrPo\_Im and Pr\_PoIm compared to Pr\_Po\_Im, with  
 266 respect to certain metrics. If the accuracy is not taken into account, then there is no strong evidence  
 267 that Pr\_Im differs in performance when compared with PrPo\_Im. In the case of Pr\_PoIm versus  
 268 Pr\_Po\_Im, there is greater evidence to support their similarity in performance with respect to  $F_1^m$ .



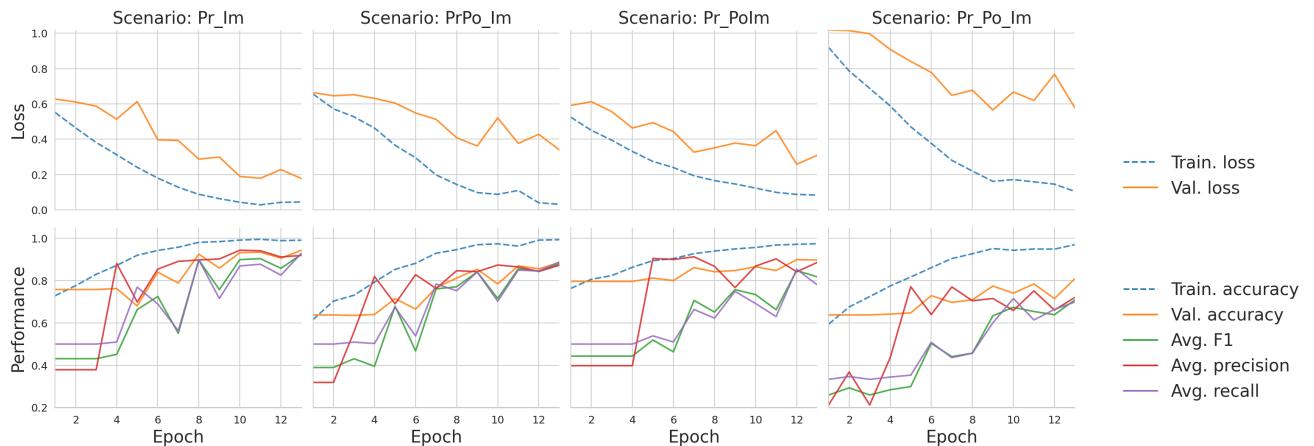
**Fig. 5.** Boxplots of validation performance metrics for each scenario and resolution combination. The upper and lower bounds of each box are the first and third quartiles; the whiskers are three standard deviations apart; and the diamonds are outliers. Mean values are depicted as white squares in each box.

269 and  $Re^m$ . Otherwise, the performance between these two scenarios is significantly different. From  
 270 Figure 5, we observe the best performances in Pr\_Im and PrPo\_Im, with all metrics greater than or

271 equal to 0.90 The performance in Pr\_PoIm is lower, and Pr\_Po\_Im has the lowest performance of  
 272 all four.

### 273 Analysis of classification strategies

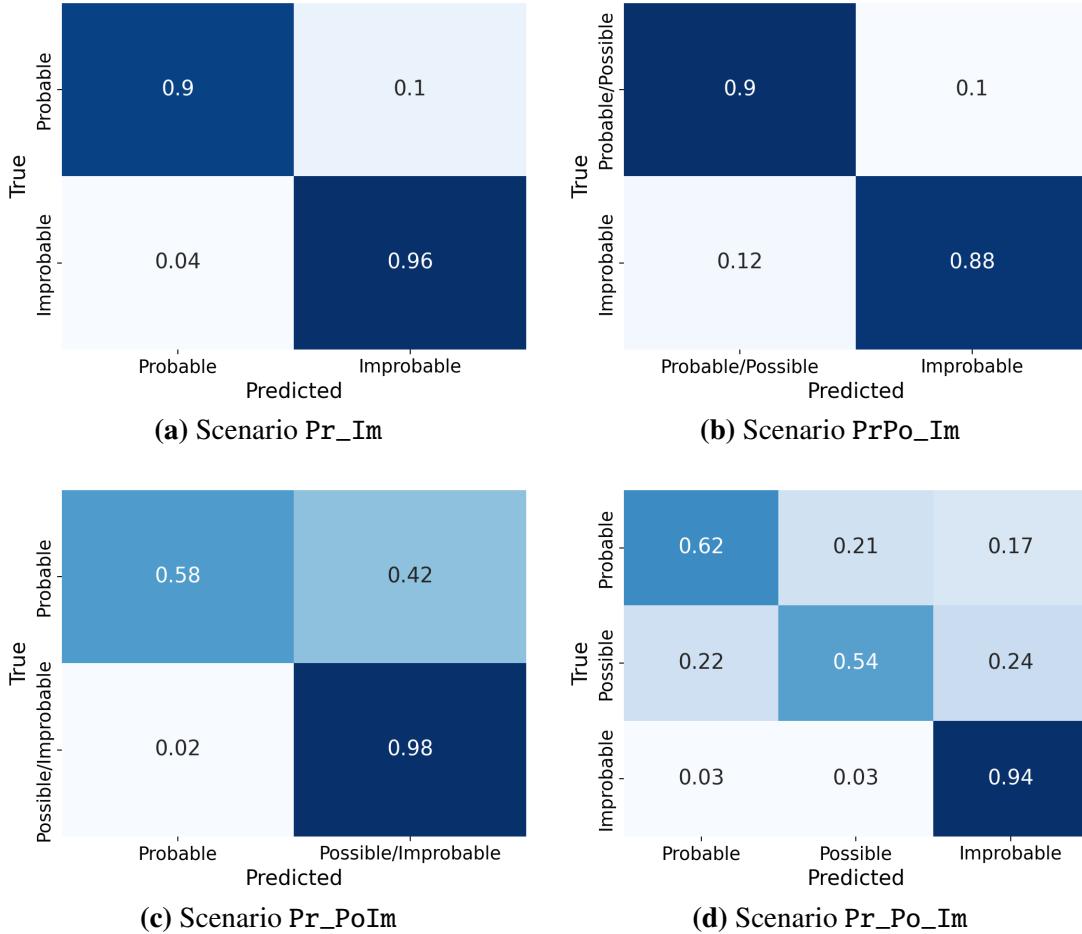
274 Given that there was no significant difference in performance among the three resolutions tested,  
 275 we focus our attention on the 128-px case, as a tradeoff between efficiency and performance. We  
 276 re-trained the model on all the four scenarios and then evaluated performance on the following  
 277 validation metrics: average precision, average recall and accuracy. The goal of this analysis was to  
 278 explore the efficiencies of various classification strategies for tree failure risk. Using a randomly  
 279 sampled validation set that was 20% of the augmented dataset as before, we trained the CNN with  
 280 the respective optimal hyperparameters for a single instance in order to compare the performance  
 281 across the scenarios. The metrics are shown in Figure 6. In each scenario, we trained the model  
 282 over 20 epochs, but we employed early stopping using validation loss as the criterion with a patience  
 283 of 3. That is, the training routine would set the weights of the final model at the epoch from which  
 284 there was no further improvement in validation loss after three subsequent iterations. From this  
 285 figure, we see that Pr\_Po\_Im performed significantly worse than the other three. This indicates  
 286 the uncertainty surrounding the subjective risk classification of the trees to begin with. Given the  
 287 results from the sensitivity tests, however, it is possible that a different training instance may have  
 288 provided better results, or perhaps, further iterations were required to achieve convergence.



**Fig. 6.** Model performance across four classification scenarios for a single training instance in each case. Training resolution was 128px. All weighted average metrics were computed on the validation set.

289 We also plotted the confusion matrix for each scenario in Figure 7. Each row of the confusion  
 290 matrix indicates the proportion of observations in a given class that are predicted to be in the  
 291 classes across the columns. Thus, the diagonal entry in each matrix represents the class-specific  
 292 recall score,  $Re_c$ . The matrices show that generally, the “Possible” category is difficult to predict  
 293 accurately as an individual class. The classifier performed best when it only had to distinguish  
 294 between “Probable” and “Improbable” (Figure 7a). Performance only suffered slightly when  
 295 “Possible” was grouped with “Probable” in the PrPo\_Im scenario (Figure 7b).

296 The detailed performance metrics for the 128-pixel case are summarized in Table 5.



**Fig. 7.** Confusion matrices for CNN performance using the 128px images across all scenarios. Each row of a confusion matrix indicates the proportional distribution of class predictions for the true members of each class. Thus, the diagonals indicate the recall for each class  $Re_c$ . The average of the diagonal values gives the macro-average  $Re^m$  for each scenario.

## CONCLUSION

We have demonstrated the efficacy of an artificial intelligence framework for predicting tree failure likelihood (Probable, Possible and Improbable) with respect to utility infrastructure. Specifically, we developed a convolutional neural network with state-of-the-art configurations. We applied data augmentation and pre-processing strategies to increase the size of our dataset by a factor of 5, thus generating 2525 images. From an initial resolution of  $4032 \times 3024$  pixels, we created three sets of image resolutions:  $64 \times 64$ ,  $128 \times 128$  and  $224 \times 224$ . We then optimized eight of our CNN hyperparameters for 12 classification scenario and image resolution combinations. We conducted ten model training trials in each of the 12 cases, using a 20% validation set to measure model performance. In these trials, the average accuracy of classification was  $\geq 0.94$ , with a maximum standard deviation of 0.7. There was no statistically significant difference between the performance across image resolutions.

We further conducted more detailed classification performance analyses over single training

<b>Scenario</b>	<b>Averaging method</b>	<b>Precision</b>	<b>Recall</b>	<b>F1-score</b>	<b>Support</b>
Pr_Im	Probable	0.9	0.75	0.81	103
	Improbable	0.92	0.97	0.95	322
	accuracy			0.92	425
	weighted avg	0.92	0.92	0.92	425
Pr_Po_Im	Probable	0.82	0.56	0.67	103
	Possible	0.43	0.74	0.54	80
	Improbable	0.93	0.86	0.89	322
	accuracy			0.78	505
Pr_PoIm	weighted avg	0.83	0.78	0.79	505
	Probable	0.77	0.75	0.76	103
	Possible/Improbable	0.94	0.94	0.94	402
	accuracy			0.9	505
PrPo_Im	weighted avg	0.9	0.9	0.9	505
	Probable/Possible	0.85	0.96	0.9	183
	Improbable	0.97	0.9	0.94	322
	accuracy			0.92	505
	weighted avg	0.93	0.92	0.92	505

**TABLE 5.** Detailed performance metrics across the four scenarios in the 128-pixel case (single trial)

instances for the 128-pixel case. Our results indicate that the CNN performed best at recalling Probable vs. Improbable cases. The Possible category appeared to be a confounding case. This was not unexpected, given the uncertainty and subjectivity in assessing these trees to begin with. The CNN, however, performed better at distinguishing between Possible/Improbable and Probable, than between Probable/Possible and Improbable. This might be an indicator that the Possible cases are more likely identified as Improbable. When all three failure-liability categories were predicted individually, the Possible case had the lowest recall score.

Nevertheless, given the relatively small input dataset of original images, these preliminary results are extremely promising for future improvements. First, we can train better models with more data. We also plan to rigorously quantify the uncertainty in ground truth category assignments by incorporating predictions from multiple experts on the same images. In order to better understand how the CNN is classifying each image, we will conduct extensive visual inference in further work, for instance, using the gradient-weighted class activation mapping approach ([Zeiler and Fergus 2014](#)).

#### DATA AVAILABILITY STATEMENT

All data, models, or code generated or used during the study are available in a GitHub repository online at <https://github.com/narslab/tree-risk-ai>.

#### ACKNOWLEDGMENTS

The authors acknowledge the support of Eversource.

- 329 **REFERENCES**
- 330 Bäuerle, A., van Onzenoodt, C., and Ropinski, T. (2021). “Net2Vis – A Visual Grammar for  
331 Automatically Generating Publication-Tailored CNN Architecture Visualizations.” *IEEE Trans-*  
332 *actions on Visualization and Computer Graphics*, 1–1.
- 333 Denker, J., Gardner, W., Graf, H., Henderson, D., Howard, R., Hubbard, W., Jackel, L. D., Baird,  
334 H., and Guyon, I. (1988). “Neural Network Recognizer for Hand-Written Zip Code Digits.”  
335 *Advances in Neural Information Processing Systems*, 1, 323–331.
- 336 dos Santos, A. A., Marcato Junior, J., Araújo, M. S., Di Martini, D. R., Tetila, E. C., Siqueira, H. L.,  
337 Aoki, C., Eltner, A., Matsubara, E. T., Pistori, H., Feitosa, R. Q., Liesenberg, V., and Gonçalves,  
338 W. N. (2019). “Assessment of CNN-Based Methods for Individual Tree Detection on Images  
339 Captured by RGB Cameras Attached to UAVs.” *Sensors*, 19(16), 3595.
- 340 Fricker, G. A., Ventura, J. D., Wolf, J. A., North, M. P., Davis, F. W., and Franklin, J. (2019). “A  
341 Convolutional Neural Network Classifier Identifies Tree Species in Mixed-Conifer Forest from  
342 Hyperspectral Imagery.” *Remote Sensing*, 11(19), 2326.
- 343 Fukushima, K. (1975). “Cognitron: A self-organizing multilayered neural network.” *Biological  
344 Cybernetics*, 20(3), 121–136.
- 345 Fukushima, K. (1980). “Neocognitron: A self-organizing neural network model for a mechanism  
346 of pattern recognition unaffected by shift in position.” *Biological Cybernetics*, 36(4), 193–202.
- 347 Fukushima, K. (1988). “Neocognitron: A hierarchical neural network capable of visual pattern  
348 recognition.” *Neural Networks*, 1(2), 119–130.
- 349 Giebel, H. (1971). “Feature Extraction and Recognition of Handwritten Characters by Homoge-  
350 neous Layers.” *Zeichenerkennung Durch Biologische Und Technische Systeme / Pattern Recog-  
351 nition in Biological and Technical Systems*, O.-J. Grüsser and R. Klinke, eds., Berlin, Heidelberg,  
352 Springer, 162–169.
- 353 Glorot, X., Bordes, A., and Bengio, Y. (2011). “Deep Sparse Rectifier Neural Networks.” *Proceed-  
354 ings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, JMLR  
355 Workshop and Conference Proceedings, 315–323 (June).
- 356 Goodfellow, J. W. (2020). “Best Management Practices - Utility Tree Risk Assessment.” *Report  
357 No. P1321*, International Society of Arboriculture.
- 358 Graziano, M., Gunther, P., Gallaher, A., Carstensen, F. V., and Becker, B. (2020). “The wider  
359 regional benefits of power grids improved resilience through tree-trimming operations evidences  
360 from Connecticut, USA.” *Energy Policy*, 138, 111293.
- 361 Gu, J., Wang, Z., Kuen, J., Ma, L., Shahroud, A., Shuai, B., Liu, T., Wang, X., Wang, G., Cai, J.,  
362 and Chen, T. (2018). “Recent advances in convolutional neural networks.” *Pattern Recognition*,  
363 77, 354–377.
- 364 Guggenmoos, S. (2003). “EFFECTS OF TREE MORTALITY ON POWER LINE SECURITY.”  
365 *Journal of Arboriculture*, 29(4), 181–196.
- 366 Guggenmoos, S. (2011). “Tree-related Electric Outages Due To Wind Loading.” *Arboriculture and  
367 Urban Forestry*, 37(4), 147–151.
- 368 He, K., Zhang, X., Ren, S., and Sun, J. (2015). “Deep Residual Learning for Image Recognition.”  
369 *arXiv:1512.03385 [cs]*.
- 370 Hubel, D. H. and Wiesel, T. N. (1959). “Receptive fields of single neurones in the cat’s striate  
371 cortex.” *The Journal of Physiology*, 148(3), 574–591.
- 372 Hubel, D. H. and Wiesel, T. N. (1962). “Receptive fields, binocular interaction and functional

- 373 architecture in the cat's visual cortex." *The Journal of Physiology*, 160(1), 106–154.2.
- 374 Hubel, D. H. and Wiesel, T. N. (1965). "Receptive fields and functional architecture in two nonstriate  
375 visual areas (18 and 19) of the cat." *Journal of Neurophysiology*, 28(2), 229–289.
- 376 Hubel, D. H. and Wiesel, T. N. (1968). "Receptive fields and functional architecture of monkey  
377 striate cortex." *The Journal of Physiology*, 195(1), 215–243.
- 378 Ioffe, S. and Szegedy, C. (2015). "Batch Normalization: Accelerating Deep Network Training by  
379 Reducing Internal Covariate Shift." *arXiv:1502.03167 [cs]*.
- 380 Jiao, P. and Alavi, A. H. (2020). "Artificial intelligence in seismology: Advent, performance and  
381 future trends." *Geoscience Frontiers*, 11(3), 739–744.
- 382 Kabrisky, M. (1966). *A Proposed Model for Visual Information Processing in the Human Brain*.  
383 University of Illinois Press.
- 384 Kingma, D. P. and Ba, J. (2017). "Adam: A Method for Stochastic Optimization." *arXiv:1412.6980  
385 [cs]*.
- 386 Koeser, A. K., McLean, D. C., Hasing, G., and Allison, R. B. (2016). "Frequency, severity,  
387 and detectability of internal trunk decay of street tree Quercus spp. in Tampa, Florida, U.S.." *Arboriculture & Urban Forestry*, 42(4), 217–226.
- 388 Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). "ImageNet Classification with Deep  
389 Convolutional Neural Networks." *Advances in Neural Information Processing Systems*, 25, 1097–  
390 1105.
- 391 LeCun, Y., Boser, B. E., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W. E., and Jackel,  
392 L. D. (1989). "Handwritten Digit Recognition with a Back-Propagation Network." *Advances in  
393 Neural Information Processing Systems* 2, 9.
- 394 Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., and Talwalkar, A. (2018). "Hyperband: A  
395 Novel Bandit-Based Approach to Hyperparameter Optimization." *Journal of Machine Learning  
396 Research*, 18(185), 1–52.
- 397 Nateghi, R., Guikema, S., and Quiring, S. M. (2014). "Power Outage Estimation for Tropical  
398 Cyclones: Improved Accuracy with Simpler Models." *Risk Analysis*, 34(6), 1069–1078.
- 399 Parent, J. R., Meyer, T. H., Volin, J. C., Fahey, R. T., and Witharana, C. (2019). "An analysis of  
400 enhanced tree trimming effectiveness on reducing power outages." *Journal of Environmental  
401 Management*, 241, 397–406.
- 402 Rooney, C. J., Ryan, H. D., Bloniarz, D. V., and Kane, B. (2005). "THE RELIABILITY OF  
403 A WINDSHIELD SURVEY TO LOCATE HAZARDS IN ROADSIDE TREES." *Journal of  
404 Arboriculture*, 31(2).
- 405 Rosenblatt, F. (1962). *Principles of Neurodynamics; Perceptrons and the Theory of Brain Mechanisms*. Spartan Books, Washington.
- 406 Salehi, H. and Burgueño, R. (2018). "Emerging artificial intelligence methods in structural engi-  
407 neering." *Engineering Structures*, 171, 170–189.
- 408 Shorten, C. and Khoshgoftaar, T. M. (2019). "A survey on Image Data Augmentation for Deep  
409 Learning." *Journal of Big Data*, 6(1), 60.
- 410 Simonyan, K. and Zisserman, A. (2015). "Very Deep Convolutional Networks for Large-Scale  
411 Image Recognition." *arXiv:1409.1556 [cs]*.
- 412 Simpson, P. and Van Bossuyt, R. (1996). "TREE-CAUSED ELECTRIC OUTAGES." *Journal of  
413 Arboriculture*, 22, 117–121.
- 414 Smiley, E. T., Matheny, N., and Lilly, S. (2017). "Best Management Practices - Tree Risk Assess-  
415 ment, Second Edition." *Report No. P1542*, International Society of Arboriculture.

- 418 Spencer, B. F., Hoskere, V., and Narazaki, Y. (2019). “Advances in Computer Vision-Based Civil  
419 Infrastructure Inspection and Monitoring.” *Engineering*, 5(2), 199–222.
- 420 Su, Z., Chow, J. K., Tan, P. S., Wu, J., Ho, Y. K., and Wang, Y.-H. (2020). “Deep convolutional  
421 neural network-based pixel-wise landslide inventory mapping.” *Landslides*.
- 422 Szegedy, C., Ioffe, S., Vanhoucke, V., and Alemi, A. (2016). “Inception-v4, Inception-ResNet and  
423 the Impact of Residual Connections on Learning.” *arXiv:1602.07261 [cs]*.
- 424 Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and  
425 Rabinovich, A. (2014). “Going Deeper with Convolutions.” *arXiv:1409.4842 [cs]*.
- 426 Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. (2015). “Rethinking the Inception  
427 Architecture for Computer Vision.” *arXiv:1512.00567 [cs]*.
- 428 Wang, N., Zhao, X., Wang, L., and Zou, Z. (2019). “Novel System for Rapid Investigation and  
429 Damage Detection in Cultural Heritage Conservation Based on Deep Learning.” *Journal of  
430 Infrastructure Systems*, 25(3), 04019020.
- 431 Wismer, S. (2018). “Targeted Tree Trimming Offers Reliability Benefits.” *T&D World* (May).
- 432 Wong, S. C., Gatt, A., Stamatescu, V., and McDonnell, M. D. (2016). “Understanding data aug-  
433 mentation for classification: When to warp?” *arXiv:1609.08764 [cs]*.
- 434 Xie, Y., Ebad Sichani, M., Padgett, J. E., and DesRoches, R. (2020). “The promise of implementing  
435 machine learning in earthquake engineering: A state-of-the-art review.” *Earthquake Spectra*,  
436 36(4), 1769–1801.
- 437 Zeiler, M. D. and Fergus, R. (2014). “Visualizing and Understanding Convolutional Networks.”  
438 *Computer Vision – ECCV 2014*, D. Fleet, T. Pajdla, B. Schiele, and T. Tuytelaars, eds., Lecture  
439 Notes in Computer Science, Cham, Springer International Publishing, 818–833.