

Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images

Anh Nguyen¹, Jason Yosinski², Jeff Clune¹

¹Department of Computer Science, University of Wyoming. ²Department of Computer Science, Cornell University.

Deep neural networks (DNNs) have recently been achieving state-of-the-art performance on a variety of pattern-recognition tasks, most notably visual classification problems. Given that DNNs are now able to classify objects in images with near-human-level performance, questions naturally arise as to what differences remain between computer and human vision. A recent study [3] revealed that changing an image (e.g. of a lion) in a way imperceptible to humans can cause a DNN to label the image as something else entirely (e.g. mislabeling a lion a library). Here we show a related result: it is easy to produce images that are completely unrecognizable to humans, but that state-of-the-art DNNs believe to be recognizable objects with 99.99% confidence (e.g. labeling with certainty that white noise static is a peacock, Fig. 2). Specifically, we take convolutional neural networks (AlexNet [1] and LeNet [2]) trained to perform well on either the ImageNet or MNIST datasets and then find images with evolutionary algorithms or gradient ascent that DNNs label with high confidence as belonging to each dataset class (Fig. 4). It is possible to produce images (Figs. 1, 2, 3) totally unrecognizable to human eyes that DNNs believe with near certainty are familiar objects, which we call "fooling images" (more generally, fooling examples). Our results shed light on interesting differences between human vision and current DNNs, raise questions about the generality of DNN computer vision, and reveal potential security concerns for applications using DNNs.

- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In Advances in neural information processing systems, pages 1097–1105, 2012.
- [2] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings* of the IEEE, 86(11):2278–2324, 1998.
- [3] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199, 2013.



Figure 1: Evolved images that LeNet believes with 99.99% confidence are the digits 0 through 9. Each column is a type of digits and each row shows images produced by an evolutionary algorithm that is either directly (*top*) or indirectly (*bottom*) encoded.



Figure 2: Evolved images that are unrecognizable to humans, but that stateof-the-art DNNs trained on ImageNet believe with $\geq 99.6\%$ certainty to be a familiar object. This result highlights differences between how DNNs and humans recognize objects. Images are either directly (*top*) or indirectly (*bottom*) encoded.



Figure 3: Evolving images to match ImageNet classes produces a tremendous diversity of images. The diversity suggests that the images are nonrandom, but that instead evolution is producing discriminative features of each target class. The mean DNN confidence scores for these images is 99.12%.



Figure 4: Although state-of-the-art deep neural networks can increasingly recognize natural images (*left panel*), they also are easily fooled into declaring with near-certainty that unrecognizable images are familiar objects (*center*). Images that fool DNNs are produced by evolutionary algorithms (*right panel*) that optimize images to generate high-confidence DNN predictions for each class in the dataset the DNN is trained on (here, ImageNet).

This is an extended abstract. The full paper is available at the Computer Vision Foundation webpage.