
Targeted Adversarial Examples for Black Box Audio Systems

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The application of deep recurrent networks to audio transcription has led to im-
2 pressive gains in automatic speech recognition (ASR) systems. Many have demon-
3 strated that small adversarial perturbations can fool deep neural networks into
4 incorrectly predicting a specified target with high confidence. Current work on
5 fooling ASR systems have focused on white-box attacks, in which the model archi-
6 tecture and parameters are known. In this paper, we adopt a black-box approach
7 to adversarial generation, combining the approaches of both genetic algorithms
8 and gradient estimation to solve the task. We achieve a 89.25% targeted attack
9 similarity after 3000 generations while maintaining 94.6% audio file similarity.

1 Introduction

11 Although neural networks have incredible expressive capacity, which allow them to be well suited
12 for a variety of machine learning tasks, they have been shown to be vulnerable to adversarial attacks
13 over multiple network architectures and datasets [7]. These attacks can be done by adding small
14 perturbations to the original input so that the network misclassifies the input but a human does not
15 notice the difference.

16 So far, there has been much more work done in generating adversarial examples for image inputs
17 than for other domains, such as audio. Voice control systems are widely used in many products
18 from personal assistants, like Amazon Alexa and Apple Siri, to voice command technologies in cars.
19 One main challenge for such systems is determining exactly what the user is saying and correctly
20 interpreting the statement. As deep learning helps these systems better understand the user, one
21 potential issue is targeted adversarial attacks on the system, which perturb the waveform of what
22 the user says to the system to cause the system to behave in a predetermined inappropriate way.
23 For example, a seemingly benign TV advertisement could be adversely perturbed to cause Alexa to
24 interpret the audio as “Alexa, buy 100 headphones.” If the original user went back to listen to the
25 audio clip that prompted the order, the noise would be almost undetectable to the human ear.

26 There are multiple different methods of performing adversarial attacks depending on what information
27 the attacker has about the network. If given access to the parameters of a network, white box attacks
28 are most successful, such as the Fast Gradient Sign Method [7] or DeepFool [11]. However, assuming
29 an attacker has access to all the parameters of a network is unrealistic in practice. In a black box
30 setting, when an attacker only has access to the logits or outputs of a network, it is much harder to
31 consistently create successful adversarial attacks.

32 In certain special black box settings, white box attack methods can be reused if an attacker creates
33 a model that approximates the original targeted model. However, even though attacks can transfer
34 across networks for some domains, this requires more knowledge of how to solve the task than
35 the original model is solving than an attacker may have [10, 13]. Instead, we propose a novel

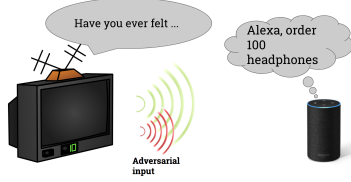


Figure 1: Example of targeted adversarial attack on speech to text systems in practice

36 combination of genetic algorithms and gradient estimation to solve this task. The first phase of the
 37 attack is carried out by genetic algorithms, which are a gradient-free method of optimization that
 38 iterate over populations of candidates until a suitable sample is produced. In order to limit excess
 39 mutations and thus excess noise, we improve the standard genetic algorithm with a new momentum
 40 mutation update. The second phase of the attack utilizes gradient estimation, where the gradients
 41 of individual audio points are estimated, thus allowing for more careful noise placement when the
 42 adversarial example is nearing its target. The combination of these two approaches provides a 89.25%
 43 average targeted attack similarity with a 94.6% audio file similarity after 3000 generations.

44 1.1 Problem statement

45 Adversarial attacks can be created given a variety of information about the neural network, such as
 46 the loss function or the output probabilities. However in a natural setting, usually the neural network
 47 behind such a voice control system will not be publicly released so an adversary will only have access
 48 to an API which provides the text the system interprets given a continuous waveform. Given this
 49 constraint, we use the open sourced Mozilla DeepSpeech implementation as a black box system,
 50 without using any information on how the transcription is done.

51 We perform our black box targeted attack on a model M given a benign input x and a target t by
 52 perturbing x to form the adversarial input $x' = x + \delta$, such that $M(x') = t$. To minimize the audible
 53 noise added to the input, so a human cannot notice the target, we maximize the cross correlation
 54 between x and x' . A sufficient value of δ is determined using our novel black box approach, so we
 55 do not need access to the gradients of M to perform the attack.

56 1.2 Prior work

57 Compared to images, audio presents a much more significant challenge for models to deal with. While
 58 convolutional networks can operate directly on the pixel values of images, ASR systems typically
 59 require heavy pre-processing of the input audio. Most commonly, the Mel-Frequency Cepstrum
 60 (MFC) transform, essentially a fourier transform of the sampled audio file, is used to convert the
 61 input audio into a spectrogram which shows frequencies over time. Models such as DeepSpeech (Fig.
 62 2) use this spectrogram as the initial input.

63 In a foundational study for adversarial attacks, Cisse et al. [5] developed a general attack framework
 64 to work across a wide variety of models including images and audio. When applying their method to
 65 audio samples, they ran into the roadblock of backpropagating through the MFC conversion layer.
 66 Thus, they were able to produce adversarial spectrograms but not adversarial .wav files.

67 Carlini and Wagner [3] overcame this challenge by developing a method of passing gradients through
 68 the MFC layer, a task which was previously proved to be difficult [5]. They applied their method
 69 to the Mozilla DeepSpeech model, which is a complex, recurrent, character-level network that can
 70 decode translations at up to 50 characters per second. With a gradient connection all the way to the
 71 raw input, they were able to achieve impressive results, including generating samples over 99.9%
 72 similar with a targeted attack accuracy of 100%. While the success of this attack opens new doors
 73 for white box attacks, adversaries in a real-life setting commonly do not have knowledge of model
 74 architectures or parameters.

75 Alzantot et al. [1] have demonstrated that black-box approaches for targeted attacks on ASR systems
 76 are possible. Using a genetic algorithm approach, they were able to iteratively apply noise to audio
 77 samples, pruning away poor performers at each generation, and ultimately end up with a perturbed
 78 version of the input that successfully fooled a classification system. This attack was conducted

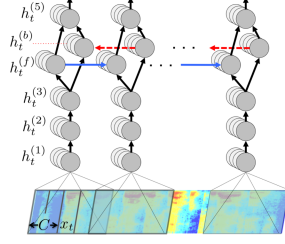


Figure 2: Diagram of Baidu’s DeepSpeech model [8]

on the Speech Commands classification model [1], which is a lightweight convolutional model for classifying up to 50 different single-word phrases.

Extending the research done by [1], we propose a genetic algorithm and gradient estimation approach to create targeted adversarial audio, but on the more complex DeepSpeech system. The difficulty of this task comes in attempting to apply black-box optimization to a deeply-layered, highly nonlinear decoder model that has the ability to decode phrases of arbitrary length. Nevertheless, the combination of two differing approaches as well as the momentum mutation update bring new success to this task.

1.3 Background

Dataset For the attack, we follow [3] and take the first 100 audio samples from the CommonVoice test set. For each, we randomly generate a 2-word target phrase and apply our black-box approach to construct an adversarial example. More details on evaluation can be found in section 3. Each sample in the dataset is a .wav file, which can easily be deserialized into a numpy array. Our algorithm operates directly on the numpy arrays, thus bypassing the difficulty of dealing with the MFC conversion.

Victim model The model we attack is Baidu’s DeepSpeech model [8], implemented in Tensorflow and open-sourced by Mozilla.¹ Though we have access to the full model, we treat it as if in a black box setting and only access the output logits of the model. In line with other speech to text systems [4, 5], DeepSpeech accepts a spectrogram of the audio file. After performing the MFC conversion, the model consists 3 layers of convolutions, followed by a bi-directional LSTM, followed by a fully connected layer. This layer is then fed into the decoder RNN, which outputs logits over the distribution of output characters, up to 50 characters per second of audio. The model is illustrated in figure 2.

Connectionist temporal classification While the DeepSpeech model is designed to allow arbitrary length translations, there is no given labeled alignment of the output and input sequences during training time. Thus the connectionist temporal classification loss (CTC Loss) is introduced, as it allows computing a loss even when the position of a decoded word in the original audio is unknown [3].

DeepSpeech outputs a probability distribution over all characters at every frame, for 50 frames per second of audio. In addition to outputting the normal alphabet a-z and space, it can output special character ϵ . Then CTC decoder $C(\cdot)$ decodes the logits as such: for every frame, take the character with the max logit. Then first, remove all adjacent duplicate characters, and then second, remove any special ϵ characters. Thus *aabeeb* will decode to *abb* [3].

As we can see, multiple outputs can decode to the same phrase. Following the notation in [3], for any target phrase p , we call π an alignment of p if $C(\pi) = p$. Let us also call the output distribution of our model y . Now, in order to find the likelihood of alignment π under y :

$$Pr(p|y) = \sum_{\pi|C(\pi)=p} Pr(\pi|y) = \sum_{\pi|C(\pi)=p} \prod_i y_{\pi}^i$$

¹ <https://github.com/mozilla/DeepSpeech>

113 as noted by [3]. This is the objective we use when scoring samples from the populations in each
 114 generation of our genetic algorithm as well as the score used in estimating gradients.

115 **Greedy decoding** As in traditional recurrent decoder systems, DeepSpeech typically uses a beam
 116 search of beam width 500. At each frame of decoding, 500 of the most likely π will be evaluated,
 117 each producing another 500 candidates for a total of 2500, which are pruned back down to 500 for
 118 the next timestep. Evaluating multiple assignments this way increases the robustness of the model
 119 decoding. However, following work in [3], we set the model to use greedy decoding. At each timestep
 120 only 1 π is evaluated, leading to a greedy assignment:

$$decode(x) = C(\underset{\pi}{\operatorname{argmax}} Pr(y(x)|\pi))$$

121 Thus, our genetic algorithm will focus on creating perturbations to the most likely sequence (if
 122 greedily approximated).

123 2 Black box attack algorithm

Algorithm 1 Black box algorithm for generating adversarial audio sample

Input: Original benign input x Target phrase t

Output: Adversarial Audio Sample x'

population $\leftarrow [x] * populationSize$

while iter $< maxIters$ and $Decode(best) \neq t$ **do**

 scores $\leftarrow -CTCLoss(population, t)$

 best $\leftarrow population[Argmax(scores)]$

if $EditDistance(t, Decode(best)) > 2$ **then**

 // phase 1 - do genetic algorithm

while populationSize children have not been made **do**

 Select *parent1* from *topk*(population) according to *softmax*(their score)

 Select *parent2* from *topk*(population) according to *softmax*(their score)

 child $\leftarrow Mutate(Crossover(parent1, parent2), p)$

end while

 newScores $\leftarrow -CTCLoss(newPopulation, t)$

 p $\leftarrow MomentumUpdate(p, newScores, scores)$

else

 // phase 2 - do gradient estimation

 top-element $\leftarrow top(population)$

 grad-pop $\leftarrow n$ copies of top-element, each mutated slightly at one index

 grad $\leftarrow (-CTCLoss(grad-pop) - scores)/mutation\text{-}delta$

 pop $\leftarrow top\text{-}element + grad$

end if

end while

return best

124 2.1 Genetic algorithm

125 As mentioned previously, Alzantot et al. [1] demonstrated the success of a black-box adversarial
 126 attack on speech-to-text systems using a standard genetic algorithm. The basic premise of our
 127 algorithm is that it takes in the benign audio sample and, through trial and error, adds noise to the
 128 sample such that the perturbed adversarial audio is similar to the benign input yet is decoded as the
 129 target, as shown in Figure 3. A genetic algorithm works well for a problem of this nature because it
 130 is completely independent of the gradients of the model. Alzantot et al. [1] used a limited dataset
 131 consisting of audio samples with just one word and a classification with a predefined number of
 132 classes. In order to extend this algorithm to work with phrases and sentences, as well as with CTC
 133 Loss, we make modifications to the genetic algorithm and introduce our novel momentum mutation.

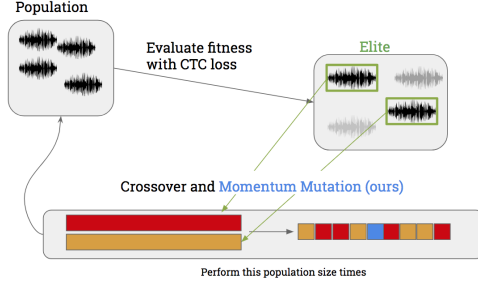


Figure 3: Diagram of our genetic algorithm

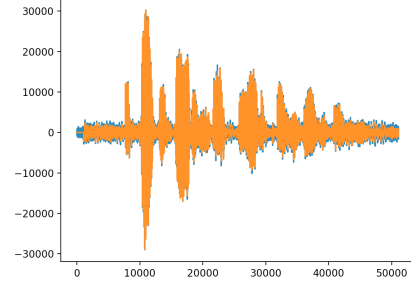


Figure 4: Overlapping of adversarial (blue) and original (orange) audio sample waveforms. The perturbation is barely noticeable

134 The genetic algorithm works by improving on each iteration, or generation, through evolutionary
 135 methods such as Crossover and Mutation [9]. For each iteration, we compute the score for each
 136 sample in the population to determine which samples are the best. Our scoring function was the
 137 CTC-Loss, which as mentioned previously, is used to determine the similarity between an input audio
 138 sequence and a given phrase. We then form our elite population by selecting the best scoring samples
 139 from our population. The elite population contains samples with desirable traits that we want to carry
 140 over into future generations. We then select parents from the elite population and perform Crossover,
 141 which creates a child by taking around half of the elements from *parent1* and the other half from
 142 *parent2*. The probability that we select a sample as a parent is a function of the sample’s score.
 143 With some probability, we then add a mutation to our new child. Finally, we update our mutation
 144 probabilities according to our momentum update, and move to the next iteration. The population will
 145 continue to improve over time as only the best traits of the previous generations as well as the best
 146 mutations will remain. Eventually, either the algorithm will reach the max number of iterations, or
 147 one of the samples is exactly decoded as the target, and the best sample is returned.

148 2.2 Momentum mutation

Algorithm 2 Mutation

Input: Audio Sample x

Mutation Probability p

Output: Mutated Audio Sample x'

```

for all  $e$  in  $x$  do
  noise  $\leftarrow \text{Sample}(\mathcal{N}(\mu, \sigma^2))$ 
  if  $\text{Sample}(\text{Unif}(0, 1)) < p$  then
     $e' \leftarrow e + \text{filter}_{\text{highpass}}(\text{noise})$ 
  end if
end for
return  $x'$ 

```

149 The mutation step is arguably the most crucial component of the genetic algorithm and is our only
 150 source of noise in the algorithm. In the mutation step, with some probability, we randomly add noise
 151 to our sample. Random mutations are critical because it may cause a trait to appear that is beneficial
 152 for the population, which can then be proliferated through crossover. Without mutation, very similar
 153 samples will start to appear across generations; thus, the way out of this local maximum is to nudge it
 154 in a different direction in order to reach higher scores.

155 Furthermore, since this noise is perceived as background noise, we apply a filter to the noise before
 156 adding it onto the audio sample. The audio is sampled at a rate of $f_s = 16kHz$, which means that
 157 the maximum frequency response $f_{max} = 8kHz$. As seen by Reichenbach and Hudspeth [14], given
 158 that the human ear is more sensitive to lower frequencies than higher ones, we apply a highpass filter

159 at a cutoff frequency of $f_{cutoff} = 7kHz$. This limits the noise to only being in the high-frequency
 160 range, which is less audible and thus less detectable by the human ear.

161 While mutation helps the algorithm overcome local maxima, the effect of mutation is limited by the
 162 *mutation probability*. Much like the step size in SGD, a low mutation probability may not provide
 163 enough randomness to get past a local maximum. If mutations are rare, they are very unlikely to
 164 occur in sequence and *add on* to each other. Therefore, while a mutation might be beneficial when
 165 accumulated with other mutations, due to the low mutation probability, it is deemed as not beneficial
 166 by the algorithm in the short term, and will disappear within a few iterations. This parallels the step
 167 size in SGD, because a small step size will eventually converge back at the local minimum/maximum.
 168 However, too large of a mutation probability, or step size, will add an excess of variability and prevent
 169 the algorithm from finding the global maximum/minimum. To combat these issues, we propose
 170 **Momentum Mutation**, which is inspired by the Momentum Update for Gradient Descent. With this
 171 update, our mutation probability changes in each iteration according to the following exponentially
 172 weighted moving average update:

$$p_{new} = \alpha \times p_{old} + \frac{\beta}{|currScore - prevScore|}$$

173 With this update equation, the probability of a mutation increases as our population fails to adapt
 174 meaning the current score is close to the previous score. The momentum update adds acceleration
 175 to the mutation probability, allowing mutations to accumulate and add onto each other by keeping
 176 the mutation probability high when the algorithm is stuck at a local maximum. By using a moving
 177 average, the mutation probability becomes a smooth function and is less susceptible to outliers in the
 178 population. While the momentum update may overshoot the target phrase by adding random noise,
 179 overall it converges faster than a constant mutation probability by allowing for more acceleration in
 180 the right directions.

Algorithm 3 Momentum Mutation Update

Input: Mutation Probability p

Scores for the new population $newScores$

Scores for the previous population $scores$

Output: Updated mutation probability p_{new}

$currScore = \max(newScores)$

$prevScore = \max(scores)$

$p_{new} = \alpha \times p_{old} + \frac{\beta}{|currScore - prevScore|}$

return p_{new}

181 2.3 Gradient estimation

182 Genetic algorithms work well when the target space is large and a relatively large number of mutation
 183 directions are potentially beneficial; the strength of these algorithms lies in being able to search
 184 large amounts of space efficiently [6]. When an adversarial sample nears its target perturbation,
 185 this strength of genetic algorithms turn into a weakness, however. Close to the end, adversarial
 186 audio samples only need a few perturbations in a few key areas to get the correct decoding. In this
 187 case, gradient estimation techniques tend to be more effective. Specifically, when edit distance of
 188 the current decoding and the target decoding drops below some threshold, we switch to phase 2.
 189 When approximating the gradient of a black box system, we can use the technique proposed by Nitin
 190 Bhagoji et al. [12]:

$$FD_x(x, \delta) = \begin{bmatrix} (g(x + \delta_1) - g(x))/\delta \\ \vdots \\ (g(x + \delta_n) - g(x))/\delta \end{bmatrix}$$

191 Here, x refers to the vector of inputs representing the audio file. δ_i refers to a vector of all zeros,
 192 except at the i^{th} position in which the value is a small δ . $g(\cdot)$ represents the evaluation function,

193 which in our case is CTCLoss. Essentially, we are performing a small perturbation at each index and
 194 individually seeing what the difference in CTCLoss would be, allowing us to compute a gradient
 195 estimate with respect to the input x .

196 However, performing this calculation in full would be prohibitively expensive, as the audio is sampled
 197 at $16kHz$ and so a simple 5-second clip would require 80,000 queries to the model for just one
 198 gradient evaluation! Thus, we only randomly sample 100 indices to perturb each generation when
 199 using this method. When the adversarial example is already near the goal, gradient estimation makes
 200 the tradeoff for more informed perturbations in exchange for higher compute.

201 3 Evaluation

202 3.1 Metrics

203 We tested our algorithm by running it on a 100 sample subset of the Common Voice dataset. For each
 204 audio sample, we generated a single random target phrase by selecting two words uniformly without
 205 replacement from the set of 1000 most common words in the English language. The algorithm
 206 was then run for each audio sample and target phrase pair for 3000 generations to produce a single
 207 adversarial audio sample.

208 We evaluated the performance of our algorithm in two primary ways. The first method is determining
 209 the accuracy with which the adversarial audio sample gets decoded to the desired target phrase. For
 210 this, we use the Levenshtein distance, or the minimum character edit distance, between the desired
 211 target phrase and the decoded phrase as the metric of choice. We then calculated the percent similarity
 212 between the desired target and the decoded phrase by calculating the ratio of the Levenshtein distance
 213 and the character length of the original input, ie. $1 - \frac{\text{Levenshtein}(M(x'),t)}{\text{len}(M(x))}$. The second method is
 214 determining the similarity between the original audio sample and the adversarial audio sample. For
 215 this, we use the accepted metric of the cross correlation coefficient between the two audio samples.

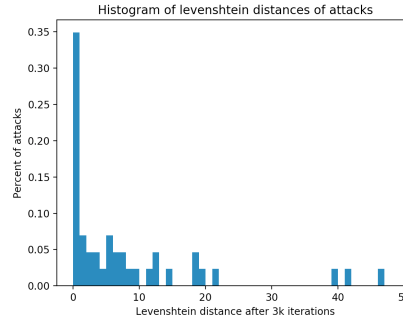


Figure 5: Histogram of levenshtein distances of attacks.

216 3.2 Results

217 Of the audio samples for which we ran our algorithm on, we achieved a 89.25% similarity between the
 218 final decoded phrase and the target using Levenshtein distance, with an average of 94.6% correlation
 219 similarity between the final adversarial sample and the original sample. The average final Levenshtein
 220 distance after 3000 iterations is 2.3, with 35% of the adversarial samples achieving an exact decoding
 221 in less than 3000 generations, and 22% of the adversarial samples achieving an exact decoding in less
 222 than 1000 generations.

223 One thing to note is that our algorithm was 35% successful in getting the decoded phrase to match
 224 the target exactly; however, noting from figure 5, the vast majority of failure cases are only a few edit
 225 distances away from the target. This suggests that running the algorithm for a few more iterations
 226 could produce a higher success rate, although at the cost of correlation similarity. Indeed, it becomes
 227 apparent that there is a tradeoff between success rate and audio similarity such that this threshold
 228 could be altered for the attacker’s needs.

A comparison of white box targeted attacks, black box targeted attacks on single words (classification), and our method:

Metric	White Box Attacks	Our Method	Single Word Black Box
Targeted attack success rate	100%	35%	87%
Average similarity score	99.9%	94.6%	89%
Similarity score method	cross-correlation	cross-correlation	human study
Loss used for attack	CTC	CTC	Softmax
Dataset tested on	Common Voice	Common Voice	Speech Commands
Target phrase generation	Single sentence	Two word phrases	Single word

One helpful visualization of the similarity between the original audio sample and the adversarial audio sample through the overlapping of both waveforms, as shown in figure 4. As the visualization shows, the audio is largely unchanged, and the majority of the changes to the audio is in the relatively low volume noise applied uniformly around the audio sample. This results in an audio sample that still appears to transcribe to the original intended phrase when heard by humans, but is decoded as the target adversarial phrase by the DeepSpeech model.

That 35% of random attacks were successful in this respect highlights the fact that black box adversarial attacks are definitely possible and highly effective at the same time.

4 Conclusion

In combining genetic algorithms and gradient estimation we are able to achieve a black box adversarial example for audio that produces better samples than each algorithm would produce individually. By initially using a genetic algorithm as a means of exploring more space through encouragement of random mutations and ending with a more guided search with gradient estimation, we are not only able to achieve perfect or near-perfect target transcriptions on most of the audio samples, we were able to do so while retaining a high degree of similarity. While this remains largely as a proof-of-concept demonstration, this paper shows that targeted adversarial attacks are achievable on black box models using straightforward methods.

Furthermore, the inclusion of momentum mutation and adding noise exclusively to high frequencies improved the effectiveness of our approach. Momentum mutation exaggerated the exploration at the beginning of the algorithm and annealed it at the end, emphasizing the benefits intended by combining genetic algorithms and gradient estimation. Restricting noise to the high frequency domain improved upon our similarity both subjectively by keeping it from interfering with human voice as well as objectively in our audio sample correlations. By combining all of these methods, we are able to achieve our top results.

In conclusion, we introduce a new domain for black box attacks, specifically on deep, nonlinear ASR systems that can output arbitrary length translations. Using a combination of existing and novel methods, we are able to exhibit the feasibility of our approach and open new doors for future research.

References

- [1] M. Alzantot, B. Balaji, and M. Srivastava. Did you hear that? Adversarial Examples Against Automatic Speech Recognition. *ArXiv e-prints*, January 2018.
- [2] N. Carlini and D. Wagner. Towards Evaluating the Robustness of Neural Networks. *ArXiv e-prints*, August 2016.
- [3] N. Carlini and D. Wagner. Audio Adversarial Examples: Targeted Attacks on Speech-to-Text. *ArXiv e-prints*, January 2018.
- [4] C.-C. Chiu, T. N. Sainath, Y. Wu, R. Prabhavalkar, P. Nguyen, Z. Chen, A. Kannan, R. J. Weiss, K. Rao, E. Gonina, N. Jaitly, B. Li, J. Chorowski, and M. Bacchiani. State-of-the-art Speech Recognition With Sequence-to-Sequence Models. *ArXiv e-prints*, December 2017.
- [5] M. Cisse, Y. Adi, N. Neverova, and J. Keshet. Houdini: Fooling Deep Structured Prediction Models. *ArXiv e-prints*, July 2017.

- 272 [6] Patrice Godefroi and Sarfraz Khurshid. Exploring Very Large State Spaces Using Genetic
273 Algorithm. *MIT*, unknown.
- 274 [7] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and Harnessing Adversarial Examples.
275 *ArXiv e-prints*, December 2014.
- 276 [8] A. Hannun, C. Case, J. Casper, B. Catanzaro, G. Diamos, E. Elsen, R. Prenger, S. Satheesh,
277 S. Sengupta, A. Coates, and A. Y. Ng. Deep Speech: Scaling up end-to-end speech recognition.
278 *ArXiv e-prints*, December 2014.
- 279 [9] J. H. Holland. Genetic algorithms. *Scientific american*, 1992.
- 280 [10] Y. Liu, X. Chen, C. Liu, and D. Song. Delving into Transferable Adversarial Examples and
281 Black-box Attacks. *ArXiv e-prints*, November 2016.
- 282 [11] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard. DeepFool: a simple and accurate method
283 to fool deep neural networks. *ArXiv e-prints*, November 2015.
- 284 [12] A. Nitin Bhagoji, W. He, B. Li, and D. Song. Exploring the Space of Black-box Attacks on
285 Deep Neural Networks. *ArXiv e-prints*, December 2017.
- 286 [13] Nicolas Papernot, Patrick D. McDaniel, Ian J. Goodfellow, Somesh Jha, Z. Berkay Celik, and
287 Ananthram Swami. Practical black-box attacks against deep learning systems using adversarial
288 examples. *CoRR*, abs/1602.02697, 2016. URL <http://arxiv.org/abs/1602.02697>.
- 289 [14] T. Reichenbach and A. J. Hudspeth. Discrimination of Low-Frequency Tones Employs Temporal
290 Fine Structure. *PLoS ONE*, 7:e45579, September 2012. doi: 10.1371/journal.pone.0045579.