
Learning noise-invariant representations for robust speech recognition

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Despite rapid advances in speech recognition, current models remain brittle to
2 superficial perturbations to their inputs. Small amounts of noise can destroy the
3 performance of an otherwise state-of-the-art model. To harden models against back-
4 ground noise, practitioners often perform data augmentation, adding artificially-
5 noised examples to the training set, carrying over the original label. In this paper,
6 we hypothesize that a clean example and its superficially perturbed counterparts
7 shouldn't merely map to the same *class* — they should map to the same *representa-*
8 *tion*. We propose invariant-representation-learning (IRL): At each training
9 iteration, for each training example, we sample a noisy counterpart. We then apply
10 a penalty term to coerce matched representations at each layer (above some chosen
11 layer). Our key results, demonstrated on the LibriSpeech dataset are the following:
12 (i) IRL significantly reduces character error rates (CER) on both 'clean' (3.3%
13 vs 6.5%) and 'other' (11.0% vs 18.1%) test sets; (ii) on several out-of-domain
14 noise settings (different from those seen during training), IRL's benefits are even
15 more pronounced. Careful ablations confirm that our results are not simply due to
16 shrinking activations at the chosen layers.

17 1 Introduction

18 Over the past several years, a series of papers have developed end-to-end deep learning systems
19 for automatic speech recognition (ASR), advancing the state of the art on a variety of benchmarks
20 [1, 2, 3, 4, 5, 6]. Typically, these models consist of either Recurrent Neural Networks (RNNs) with
21 Sequence-to-Sequence (Seq2Seq) architectures [7] and attention mechanisms [8, 9], RNN transducers
22 [10], transformer networks [11, 6], convolutional neural networks paired with transformer networks
23 [12, 13], or RNNs trained with CTC loss [14]. Often, these models act on spectral features, e.g.,
24 Mel-Frequency Cepstral Coefficients (MFCC) [15].

25 While these systems achieve impressive accuracy when trained and evaluated on clean data, they
26 suffer a well-documented sensitivity to changing noise levels and various noise types [16]. Perhaps
27 this vulnerability should not be surprising, given the significant impact that background noise can
28 have on MFCC features [16].

29 One simple strategy to combat the vulnerability of deep nets to background noise is a technique
30 known generally as *data augmentation*, and as *multi-condition training* in the speech recognition
31 community. Here, we augment the original data by applying transformations to which we want our
32 models to be invariant and assigning these perturbed data points the same label as the unperturbed
33 originals. While the computer vision literature has long focused on perturbations like random crops,
34 rotations, translations, and Gaussian noise [17, 18, 19, 20, 21], data augmentation papers in the ASR
35 literature commonly sample snippets of additive background noise from datasets such as MUSAN
36 [22], which contains environmental noise (dial tones, thunder, footsteps, animal noises, etc), music

37 (baroque, classical, romantic, jazz, bluegrass, hip-hop, etc.), and speech. ASR models trained with
38 such augmented data have demonstrated lower grapheme error rates on noisy data [23, 24, 25].

39 In this paper, we draw inspiration from the human ability to recognize not only that a clean clip and
40 its noisy counterpart belong to the same category but that they are produced from the same exact
41 recording. Thus, we propose models that map both clean inputs and their noisy counterparts onto the
42 same point in representation space, introducing this inductive bias via regularization terms, penalizing
43 differences between the hidden representations produced from real and noisy data. Throughout train-
44 ing, for each clean example, we synthesize one noisy counterpart, using a custom data augmentation
45 pipeline that first selects a random noise snippet and volume level, adding the two raw waveforms
46 and then generating the corresponding MFCC features on the fly. At each iteration, we apply the
47 original cross-entropy loss on the predictions for both clean and perturbed inputs and also penalize
48 the difference in hidden activations encouraging corresponding activations as quantified by both
49 cosine distance and L_2 distance.

50 Our experiments address the LibriSpeech dataset [26], building on a Seq2Seq baseline with cross-
51 entropy loss. To keep the empirical study clean, we do not use a language model. We run all
52 experiments both on the standard dev and test sets and also under a variety of out-of-domain noise
53 conditions. First, we show that while data augmentation improves generalization error on both
54 the original task and under out-of-domain noise, the models still suffer significant degradation in
55 performance. Next, we show that Invariant-Representation Learners (IRLs) improve significantly
56 over generic data augmentation models, both on the *clean* and *other* (the more challenging dataset
57 with higher word error rate) subsets of the LibriSpeech test set. Comparisons against an adversarial
58 approach proposed by [27] and the logit pairing approach due to [28] demonstrate the significant
59 advantage of IRL. We then demonstrate that on a variety of simulated out-of-domain noise conditions,
60 the IRL models are considerably more robust than all baselines. Finally, we perform ablation
61 experiments, showing that our models trained with the IRL algorithm outperform well-known
62 regularization tactics like weight decay applied on the same representations.

63 1.1 Related work

64 A number of proposed models address the goal of noise-robust speech recognition: [29] proposes a
65 method called *noise-aware training* that introduces information about the environment as additional
66 inputs to DNN-based acoustic models. [23] proposes augmenting training examples with additive
67 noise sampled from the DEMAND noise database training examples. [27] seeks noise-invariant
68 representations in DNN-HMM architectures through an adversarial learning setup. [30] shows the
69 training on multi-modal data leads to noise robust models. [31] demonstrates that modeling speech as
70 a linear combination of exemplars results in noise-robust ASR models. [32] proposes deep recurrent
71 autoencoders to denoise input features. [33] presents an overview of methods for noise-robust ASR,
72 including recursive cepstral mean and variance normalization [34], joint adaptive training [35], and
73 speaker adaptive training [36]. To our knowledge, no prior work in speech recognition employs our
74 simple approach of penalizing distance between the hidden representations corresponding to clean
75 and noisy signals.

76 In the most similar paper, [27] claimed that with adversarially trained DNN-HMM systems, the best
77 performance gain is achieved when a small number of noise types are available for training. When
78 using 6 noise classes (airport, babble, car, restaurant, street, and train), [27] found that there was no
79 significant difference between the adversarial and baseline models. In contrast, our models show a
80 CER improvement over baseline of 3.1% absolute on test-clean and 6.5% absolute on test-other using
81 hundreds of noise classes.

82 2 Noise-invariant representations

83 To begin, we formally describe our loss function for enforcing noise-invariant representations on
84 the outputs of a given layer. Because our first proposed model focuses noise-invariance in the
85 encoding layer, we dub models using such loss functions *IRL-E*. In other experiments, we apply a
86 cumulative penalty, additionally requiring noise-invariant representations at all subsequent layers,
87 naming this model *IRL-C*. We begin by describing IRL-E. Subsequently, extension to IRL-C will be
88 straightforward.

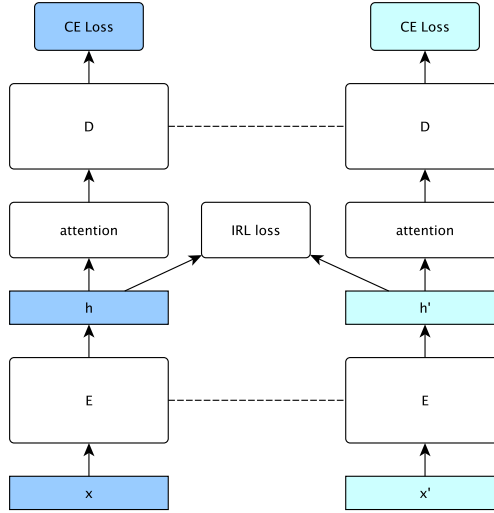


Figure 1: Diagram demonstrating the various terms in the IRL loss function as applied to a Seq2Seq attention model. Dotted lines represent shared weights.

89 2.1 IRL-E

90 The IRL algorithm is simple: First, during training, for each example \mathbf{x} , we produce a noisy version
 91 by sampling from $\mathbf{x}' \sim \nu(\mathbf{x})$, where ν is a stochastic function. In our experiments, this function takes
 92 a random snippet from a noise database, sets its amplitude by drawing from a normal distribution, and
 93 adds it to the original (in sample space), before converting to spectral features. We then incorporate a
 94 penalty term in our loss function to penalize the distance between the encodings of the original data
 95 point $\phi_e(\mathbf{x})$ and the noisy data point $\phi_e(\mathbf{x}')$, where ϕ_l is representation at layer l . In our experiments,
 96 we choose ϕ_e to be the output of the encoder in our Seq2Seq model. We illustrate the learning setup
 97 graphically in Figure 1. In short, our loss function consists of three terms, one to maximize the
 98 probability assigned to the clean example’s label, another to maximize the probability our model
 99 assigned to the noisy example’s (identical) label scaled by hyper-parameter α , and a penalty term to
 100 induce noise-invariant representations L_d . In the following equations, we express the loss calculated
 101 on a single example \mathbf{x} and its noisy counterpart \mathbf{x}' , omitting sums over the dataset for brevity.

$$L(\theta) = L_c(\mathbf{x}; \theta) + \alpha L_c(\mathbf{x}'; \theta) + L_d(\mathbf{x}, \mathbf{x}'; \theta),$$

102 where θ denotes our model parameters. Because our experiments address multiclass classification,
 103 our primary loss L_c is cross-entropy:

$$L_c(\mathbf{x}; \theta) = - \sum_{k=1}^C y_k \log \hat{y}_k(\mathbf{x}; \theta),$$

104 where C denotes the vocabulary size and \hat{y} is our model’s softmax output. To induce similar
 105 representations for clean and noised data, we apply a penalty consisting of two terms, the first
 106 penalizes the L_2 distance between $\phi_e(\mathbf{x})$ and $\phi_e(\mathbf{x}')$, the second penalizes their negative cosine
 107 distance.

$$L_d(\mathbf{x}, \mathbf{x}'; \theta) = \gamma (\phi_e(\mathbf{x}) - \phi_e(\mathbf{x}'))^2 - \lambda \frac{\phi_e(\mathbf{x}) \cdot \phi_e(\mathbf{x}')}{\|\phi_e(\mathbf{x})\| \cdot \|\phi_e(\mathbf{x}')\|}$$

108 We jointly penalize the L_2 and cosine distance for the following reason. It is possible to lower the
 109 L_2 distance between the two (clean and noisy) hidden representations simply by shrinking the scale
 110 of all encoded representations. Trivially, these could then be dilated again simply by setting large
 111 weights in the following layer. On the other hand, it is possible to assign high cosine similarity to
 112 the two vectors but for their magnitudes to vary significantly. By jointly penalizing L_2 and cosine
 113 distance, we require that both the clean and noisy representations point in the same direction and are
 114 close to each other in magnitude.

115 2.2 Applying IRL cumulatively across layers (IRL-C)

116 It is possible for representations to be close (but not identical) in the encoder layer, but to subsequently
117 be pushed apart in subsequent decoder layers. Thus, we introduce another model, *IRL-C* (C for
118 *cumulative*), that additionally applies the IRL penalty on all subsequent decoder layers. By requiring
119 noise-invariant representations in multiple layers, we ensure that each training example and its
120 randomly-sampled noisy counterpart have similar representations throughout the network. Note that
121 if the encodings of the clean and noisy examples are identical at the encoder layer, then all subsequent
122 layers will also be identical and thus those penalties will go to 0. We can express this loss as a sum
123 over successive representations ϕ_l of the clean $\phi_l(\mathbf{x})$ and noisy $\phi_l(\mathbf{x}')$ data:

$$L_d(\mathbf{x}, \mathbf{x}'; \theta) = \sum_{l=e}^L \left[\gamma(\phi_l(\mathbf{x}) - \phi_l(\mathbf{x}'))^2 - \lambda \frac{\phi_l(\mathbf{x}) \cdot \phi_l(\mathbf{x}')}{\|\phi_l(\mathbf{x})\| \cdot \|\phi_l(\mathbf{x}')\|} \right]$$

124 In our experiments, we find that IRL-C consistently gives a small improvement over results achieved
125 with IRL-E.

126 2.3 Application to recurrent speech models

127 As described to this point, our loss can be applied on any feedforward neural network with any
128 noise process ν . Applying our technique to recurrent neural networks requires just a few additional
129 considerations. Primarily, we must decide how to deal with the sequence structure. Two natural
130 choices are (i) to concatenate the representations for a given layer across time steps, and then to
131 apply our penalty on the concatenated representations and (ii) to apply the penalty separately at each
132 time step and then to sum (or equivalently, up to a scaling factor to average) over the time steps.
133 These approaches are identical for the L_2 penalty but not for the cosine distance penalty, owing to
134 the normalizing factor which may be different at each time step. In this work we take approach (i)
135 concatenating the representations across time steps and then calculating the penalty.

136 All of our models are based off of the sequence-to-sequence due to [9]. The input to the encoder is a
137 sequence of spectral features, here MFCC, which are encoded by several consecutive layers of LSTM
138 units. The encoder output states are then passed through an attention mechanism which computes the
139 similarity between the decoder hidden states and the encoder output states. The output is a softmax
140 over the vocabulary (here, characters) at each decoder time step.

141 In our experiments with IRL-E (penalty applied on a single layers), we use the output of the encoder
142 to calculate the penalty. Note that there is one output per step in the input sequence and thus we are
143 concatenating across the T_1 steps.

144 To calculate IRL-C, we also start with the encoder output concatenating across all T_1 sequence steps
145 to calculate the IRL penalty. However, for all subsequent layers, we are acting upon layers in the
146 decoder, and thus concatenating across the number of decoding sequence steps T_2 for calculating
147 these terms in the IRL-C penalty.

148 3 Datasets

149 **LibriSpeech** We evaluate all models on the LibriSpeech [26] dataset. This dataset consists of
150 roughly 1000 hours of audio split into training, dev and test partitions. The dataset was carefully
151 designed to ensure that no speaker (person) appears in multiple partitions. Within both the dev and the
152 test partitions, the data is further subdivided into “clean” and “other” subsets based on the speakers.
153 The “clean” portion contains those speakers for which a baseline model had the lowest CER, and the
154 “other” portion contains those speakers for whom the error rate was high. Following common practice
155 in the literature on these datasets, we evaluate all models on the dev-clean, dev-other, test-clean, and
156 test-other splits separately.

157 **The MUSAN noise dataset** For our additive noise, we draw upon samples from the MUSAN
158 noise dataset [22]. MUSAN was released under a flexible Creative Commons license and consists of
159 approximately 109 hours of noise sampled at 16kHz. The dataset contains music from several genres,
160 namely baroque, classical, romantic jazz, bluegrass, and hip-hop, among others, speech from twelve
161 languages, and a wide assortment of technical and non-technical noises. To generate noisy audio, we

162 first add MUSAN noise to the training data point at a signal-to-noise ratio drawn from a Gaussian
163 with a mean of 12dB and a variance of 8dB. This aligns roughly with the scale of noise employed in
164 other papers using multi-condition training [2].

165 4 Experiments

166 Before presenting our main results, we briefly describe the model architectures, training details, and
167 the various baselines that we compare against. We also present details on our pipeline for synthesizing
168 noisy speech and explain the experimental setup for evaluating on out-of-domain noise.

169 4.1 Model architecture

170 To facilitate reliable comparisons between our methods and various baseline training schemes, we
171 conduct all experiments using identical architectures and tuning schemes. Because we conduct a
172 large number of experiments and because of the computational expense of unrolling of long speech
173 sequences, we struck a balance between performance and speed when choosing the basic architecture.
174 The encoder for our base model consists of 4 layers: 2 encoder BLSTM layers with 320 hidden
175 units each, followed by 2 encoder LSTM layers with 320 hidden units each. Our decoder accesses
176 the encoded representations using dot product attention, and contains 4 decoder LSTM layers, with
177 320 hidden units each. Notably, our first encoder layer halves the sequence length by concatenating
178 adjacent inputs along the temporal axis. Each model across all of our comparisons has the exact same
179 number of trainable parameters. To keep things simple, we do not use an external language model.
180 Instead we decode predictions from all models via beam search with width 10.

181 To ensure fair comparisons, we perform hyper-parameter searches separately for each model and
182 account for variability due to initialization by training each model 5 times and keeping the best run
183 as determined on the dev-other partition. Specifically, we tune the weights on our losses by trying
184 each of the scale values (0.001, 0.01, 0.1, 1, 10, and 100). We found that an α of 1 (the weight on
185 the cross-entropy loss of the noised data), a γ of 0.01 (the weight on the L2 distance loss), and a λ of
186 0.01 (the weight on the cosine distance loss) worked well.

187 4.2 Training details

188 We train all models with the Adam optimizer with an initial learning rate of 0.001. We employ a
189 learning rate schedule similar to NewBob [37] that decreases by a factor of 2 if there is an increase
190 in validation perplexity epoch-over-epoch. We employ a stopping criterion that ends training if
191 validation perplexity does not decrease for three epochs in a row. We limit each models to a maximum
192 of 40 epochs, although our networks generally converge within 20 epochs.

193 The primary loss function for each model is cross-entropy loss and our primary evaluation metric to
194 evaluate all models is the character error rate. As described above, the additional loss terms for our
195 IRL models are L2 loss and cosine distance between representations of clean and noisy audio.

196 4.3 Baselines

197 Our baseline models include a model trained on the standard training data, a model trained with
198 noise-augmented data, a model trained with noise augmented data and weight decay, and a data
199 augmented model supervised with L2 loss to push activations of the encodings to 0. These ablation
200 tests provide evidence that our IRL algorithm isn't simply penalizing the norm of the encodings.

- 201 • **Baseline:** Our base model trains the baseline sequence-to-sequence model on the original
202 960 hours of LibriSpeech training data.
- 203 • **Data augmentation:** Our data augmentation model trains the sequence-to-sequence model
204 on both the examples from the 960 hour LibriSpeech training corpus and the randomly
205 generated noisy counterparts.
- 206 • **Adversarial:** The adversarial model consists of an adversarial noise discriminator trained
207 on top of the encoder outputs. The discriminator consists of 2 layers of 256 ReLU units and a
208 single unit sigmoid output. We train the discriminator to classify whether the representation
209 originates from clean or noised inputs. The encoder meanwhile is trained both to minimize

210 the classification loss and to fool the discriminator, in a scheme similar to the reverse gradient
 211 technique in the domain-adversarial approach due to [38] and applied to speech by [27].

- 212 • **Logit pairing:** Our final baseline consists of the logit pairing model due to [28] which
 213 applies L2 loss and cosine distance loss on the final decoder layer logits, enforcing noise-
 214 invariant representations but only on the output layer.

215 4.4 Synthesizing noise

216 We train all models on the LibriSpeech corpus, generating noisy data by adding randomly selected
 217 noise tracks from the MUSAN dataset with a signal to noise ratio drawn from a Gaussian distribution
 218 (12dB mean, 8dB standard deviation) and temporal shift drawn from a uniform distribution (with
 219 a range of 0 to 1000ms). For the data augmentation model, this result resembles the typical data
 220 augmentation (multi-condition training) procedure.

221 4.5 Out-of-domain noise

222 Next, we evaluate each of our models on a variety of noise conditions that were not seen at training
 223 time. In particular, we consider the following out-of-domain noise conditions: (i) augmenting the
 224 test-clean split with overlapping out-of-domain speech from the WSJ-0 dataset [39] to simulate
 225 multi-speaker environments, (ii) applying additive noise with various SNRdb to simulate varying
 226 noise levels, (iii) modulating the volume of the clean signal to simulate different levels of speaker
 227 loudness, (iv) convolving the original wave file with room impulse responses to simulate the effect
 228 of room reverberation on speech, and (v) re-sampling to 8kHz to simulate telephony data. For each
 229 setting, we measure CER on the out-of-domain noise-augmented test-clean data.

230 5 Results

Table 1: Evaluation and test set character error rate on the LibriSpeech corpus.

	Evaluation set		Test set	
	dev-clean	dev-other	test-clean	test-other
Baseline	6.7%	17.8%	6.5%	18.1%
Data aug.	6.4%	16.8%	6.4%	17.5%
Adversarial	6.7%	16.7%	6.5%	17.6%
Logit pairing	5.1%	14.5%	5.1%	14.8%
IRL-E	3.6%	11.0%	3.5%	11.2%
IRL-C	3.4%	10.7%	3.3%	11.0%

231 Our IRL-C model achieves the best CER on both test-clean and test-other 3.3% and 11%, respec-
 232 tively (Table 1). This compares baseline scores of 6.5% and 18.1%, respectively. We note that by
 233 comparison, conventional data augmentation is only marginally effective. Among the baselines
 234 that we consider, logit pairing performs best (5.1% and 14.8%) although the improvements are not
 235 comparable to either IRL model.

236 We found that weight decay slowed down network convergence and did not outperform pure data
 237 augmented training. However, [40] showed that weight decay is most effective with separate λ_{rec}
 238 and λ_{nonrec} hyper-parameters for determining the strength of regularization for the recurrent and
 239 non-recurrent weight matrices. We have not tried this in our experiments. Additionally, we discovered
 240 that applying multi-condition training while naively lowering the activations of hidden representations
 241 leads to nearly identical performance (on both the original and out-of-domain noise perturbed test
 242 data) and convergence trajectory as the base model trained on noise augmented data. These results
 243 support our hypothesis that models trained with the IRL algorithm do not trivially decrease the
 244 magnitude of intermediate representations.

245 Our final experiments test the effects of various out-of-domain noise on our models. The results are
 246 shown in Table 2. We found that our models trained with the IRL procedure had stronger results
 247 (and significantly less degradation) across all tasks compared to the baseline and the purely data

Table 2: Character error rate for test-clean augmented with noise

	CER on noisy data					
	Base	Data aug.	Adv.	Logit	IRL-E	IRL-C
Error on test-clean	6.5%	6.4%	6.5%	5.1%	3.5%	3.3%
In-domain (6SNRdB)	27.8%	10.8%	16.5%	8.7%	6.0%	5.7%
In-domain (12SNRdB)	13.5%	7.8%	12.1%	6.2%	4.2%	4.1%
Impulse convolve	24.1%	21.0%	28.3%	47.6%	18.0%	13.8%
Speech (6SNRdB)	91.5%	32.0%	67.7%	33.0%	16.4%	14.1%
Speech (12SNRdB)	77.8%	15.2%	34.7%	11.1%	7.6%	6.8%
Volume (+6 dB)	6.5%	6.4%	9.8%	5.1%	3.6%	3.5%
Volume (-6 dB)	6.5%	6.3%	9.6%	5.0%	3.6%	3.5%
Telephony	14.2%	12.2%	21.3%	10.3%	7.1%	6.4%

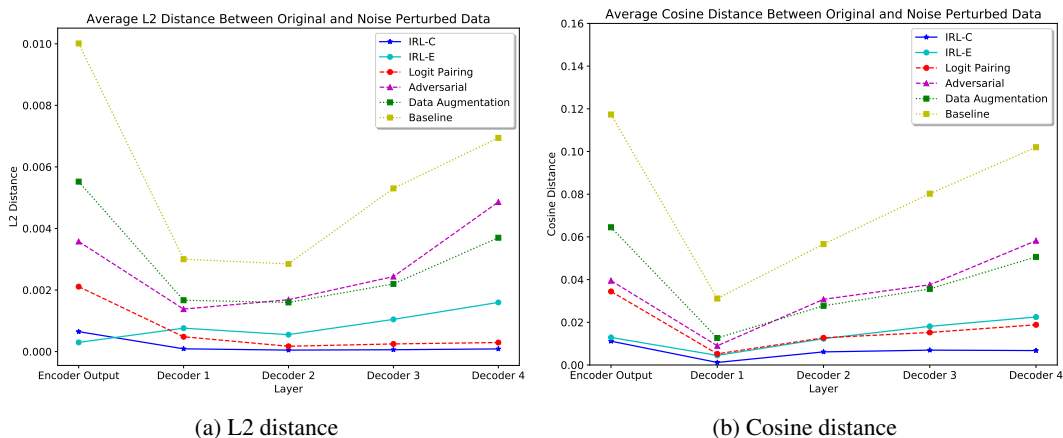


Figure 2: Average distance between original and noised data for various models (distinct lines) and various layers (x-axis). Subplot (a) depicts L2 distance and (b) depicts cosine distance.

248 augmented models. When applying various room reverberation on speech, we found that the IRL-C
 249 model had a character error rate of 13.8% compared to 21.0% on the data augmented model and
 250 24.1% on the baseline model. Our IRL-C model shows 14.1% character error rate on out-of-domain
 251 overlapping speech compared to 91.5% for the baseline and 32.0% on the data augmented model. We
 252 found that decreasing the signal-to-noise ratio also effected the baseline models more than the models
 253 trained on the IRL algorithm: our IRL-C model received a character error rate of 5.7% compared to
 254 27.8% for baseline and 10.8% for the purely data augmented model. We found that modifying the
 255 volume of the speaker did not effect the accuracy of any of the networks. Finally, we found that our
 256 models trained with the IRL algorithm performed better for re-sampled telephony data, achieving
 257 a character error rate of 6.4% for IRL-C compared to 14.2% for baseline and 12.2% for the purely
 258 data augmented model.

259 We also executed some empirical analysis to determine the effect of the various approaches on the
 260 distances between noisy examples and their clean counterparts in representation space. In general,
 261 our IRL models have the lowest L2 and cosine distances between noisy representations and the
 262 clean counterparts. In Figure 2, you can see that although the IRL-E and IRL-C model models have
 263 similarly close representations at the encoder layer, neither reaches 0 distance. Then for IRL-E over
 264 the subsequent layers, the clean and noisy representations diverge again, while for IRL-C they remain
 265 close throughout.

266 6 Conclusions

267 In this paper, we demonstrated that enforcing noise-invariant representations by penalizing differences
268 between pairs of clean and noisy data can increase model accuracy on the ASR task, produce models
269 that are robust to out-of-domain noise, and improve convergence speed. The performance gains
270 achieved by IRL come without any impact to inference throughput. We note that our core ideas
271 here can be applied broadly to deep networks for any supervised task. While the speech setting is
272 particularly interesting to us, our methods are equally applicable to other machine learning fields,
273 notably computer vision. One natural extension might be to experiment with various other loss
274 functions such as triplet losses, requiring that noisy data be both close to its clean counterpart and
275 further away from *different* clean data. Additionally, our approach might be well-suited to conferring
276 greater robustness to adversarial examples. The comparative improvements over requiring invariant
277 hidden representations vs. invariant logits here raises the possibility that we might be able to realize
278 similar gains over logit pairing in the adversarial setting.

279 References

- 280 [1] Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro, Greg Diamos, Erich Elsen, Ryan
281 Prenger, Sanjeev Sathesh, Shubho Sengupta, Adam Coates, et al., “Deep speech: Scaling up
282 end-to-end speech recognition,” p. arXiv:1412.5567, 2014.
- 283 [2] Dario Amodei, Rishita Anubhai, Eric Battenberg, Carl Case, Jared Casper, Bryan Catanzaro,
284 Jingdong Chen, Mike Chrzanowski, Adam Coates, Greg Diamos, and et al., “Deep speech 2:
285 End-to-end speech recognition in english and mandarin,” p. arXiv:1512.02595, 2015.
- 286 [3] Yajie Miao, Mohammad Gowayyed, and Florian Metze, “Eesen: End-to-end speech recognition
287 using deep rnn models and wfst-based decoding,” 2015.
- 288 [4] Dzmitry Bahdanau, Jan Chorowski, Dmitriy Serdyuk, Philemon Brakel, and Yoshua Bengio,
289 “End-to-end attention-based large vocabulary speech recognition,” in *IEEE International
290 Conference on Acoustics, Speech and Signal Processing*, 2016.
- 291 [5] Albert Zeyer, Kazuki Irie, Ralf Schlüter, and Hermann Ney, “Improved training of end-to-end
292 attention models for speech recognition,” *Interspeech*, 2018.
- 293 [6] Shiyu Zhou, Linhao Dong, Shuang Xu, and Bo Xu, “Syllable-based sequence-to-sequence
294 speech recognition with the transformer in mandarin chinese,” *arXiv preprint arXiv:1804.10752*,
295 2018.
- 296 [7] Ilya Sutskever, Oriol Vinyals, and Quoc V Le, “Sequence to sequence learning with neural
297 networks,” in *Advances in Neural Information Processing Systems*, 2014.
- 298 [8] Minh-Thang Luong, Hieu Pham, and Christopher D Manning, “Effective approaches to
299 attention-based neural machine translation,” p. arXiv:1508.04025, 2015.
- 300 [9] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align
301 and translate,” arXiv preprint, 2014, p. arXiv:1409.0473.
- 302 [10] Alex Graves, “Sequence transduction with recurrent neural networks,” *arXiv preprint
303 arXiv:1211.3711*, 2012.
- 304 [11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
305 Łukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” in *Advances in Neural
306 Information Processing Systems*, 2017.
- 307 [12] Shiyu Zhou, Linhao Dong, Shuang Xu, and Bo Xu, “A comparison of modeling units in
308 sequence-to-sequence speech recognition with the transformer on mandarin chinese,” *arXiv
309 preprint arXiv:1805.06239*, 2018.
- 310 [13] Ronan Collobert, Christian Puhresch, and Gabriel Synnaeve, “Wav2letter: An end-to-end
311 convnet-based speech recognition system,” *arXiv preprint arXiv:1609.03193*, 2016.
- 312 [14] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton, “Speech recognition with deep
313 recurrent neural networks,” in *International Conference on Acoustics, Speech and Signal
314 Processing*. IEEE, 2013.
- 315 [15] S. Davis and P. Mermelstein, “Comparison of parametric representations for monosyllabic word
316 recognition in continuously spoken sentences,” *IEEE Transactions on Acoustics Speech and
317 Signal Processing*, 1980.

- 318 [16] U. Bhattacharjee, Swapnanil Gogoi, and Rubi Sharma, "A statistical analysis on the impact
319 of noise on mfcc features for speech recognition," *2016 International Conference on Recent
320 Advances and Innovations in Engineering*, 2016.
- 321 [17] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, "Imagenet classification with deep
322 convolutional neural networks," in *Advances in Neural Information Processing Systems*, 2012.
- 323 [18] Guozhong An, "The effects of adding noise during backpropagation training on a generalization
324 performance," *Neural computation*, 1996.
- 325 [19] Yves Grandvalet and Stéphane Canu, "Comments on" noise injection into inputs in back
326 propagation learning",," *IEEE Transactions on Systems, Man, and Cybernetics*, 1995.
- 327 [20] Chris M Bishop, "Training with noise is equivalent to tikhonov regularization," *Neural
328 Computation*, 1995.
- 329 [21] Yves Grandvalet, Stéphane Canu, and Stéphane Boucheron, "Noise injection: Theoretical
330 prospects," *Neural Computation*, 1997.
- 331 [22] D. Snyder, G. Chen, and D. Povey, "Musan: A music, speech, and noise corpus," arXiv preprint,
332 2015, p. arXiv:1510.08484v1.
- 333 [23] S. Yin, C. Liu, Z. Zhang, Y. Lin, D. Wang, J. Tejedor, T. Zheng, and Y. Li, "Noisy training for
334 deep neural networks in speech recognition," *Journal on Audio, Speech, and Music Processing*,
335 2015.
- 336 [24] Hong Yu, Achintya Sarkar, Dennis Alexander Lehmann Thomsen, Zheng-Hua Tan, Zhanyu
337 Ma, and Jun Guo, "Effect of multi-condition training and speech enhancement methods on
338 spoofing detection," *IEEE First International Workshop on Sensing, Processing and Learning
339 for Intelligent Machines*, 2016.
- 340 [25] J Rajnoha, "Multi-condition training for unknown environment adaptation in robust asr under
341 real conditions," *Acta Polytechnica*, 2009.
- 342 [26] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on
343 public domain audio books," *IEEE International Conference on Acoustics, Speech and Signal
344 Processing*, 2015.
- 345 [27] D. Serdyuk, K. Audhkhasi, P. Brakel, B. Ramabhadran, S. Thomas, and Y. Bengio, "Invariant
346 representations for noisy speech recognition," *Computing Research Repository*, 2016, p.
347 abs/1612.01928.
- 348 [28] H. Kannan, A. Kurakin, and I. Goodfellow, "Adversarial logit pairing," arXiv preprint, 2018, p.
349 arXiv:1803.06373.
- 350 [29] Michael L Seltzer, Dong Yu, and Yongqiang Wang, "An investigation of deep neural networks
351 for noise robust speech recognition," in *Acoustics, Speech and Signal Processing (ICASSP),
352 2013 IEEE International Conference on*. IEEE, 2013.
- 353 [30] Jing Huang and Brian Kingsbury, "Audio-visual deep learning for noise robust speech recogni-
354 tion," *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2013.
- 355 [31] Jort F Gemmeke, Tuomas Virtanen, and Antti Hurmalainen, "Exemplar-based sparse represen-
356 tations for noise robust automatic speech recognition," *IEEE Transactions on Audio, Speech,
357 and Language Processing*, 2011.
- 358 [32] Andrew L Maas, Quoc V Le, Tyler M O'Neil, Oriol Vinyals, Patrick Nguyen, and Andrew Y Ng,
359 "Recurrent neural networks for noise reduction in robust asr," in *Conference of the International
360 Speech Communication Association*, 2012.
- 361 [33] Jinyu Li, Li Deng, Yifan Gong, and Reinhold Haeb-Umbach, "An overview of noise-robust
362 automatic speech recognition," *IEEE Transactions on Audio, Speech, and Language Processing*,
363 2014.
- 364 [34] Olli Viikki, David Bye, and Kari Laurila, "A recursive feature vector normalization approach
365 for robust speech recognition in noise," *IEEE International Conference on Acoustics, Speech
366 and Signal Processing*, 1998.
- 367 [35] Hank Liao and MJ F Gales, "Adaptive training with joint uncertainty decoding for robust
368 recognition of noisy data," *IEEE International Conference on Acoustics, Speech and Signal
369 Processing*, 2007.

- 370 [36] Tasos Anastasakos, John McDonough, Richard Schwartz, and John Makhoul, "A compact
371 model for speaker-adaptive training," IEEE International Conference on Spoken Language,
372 1996.
- 373 [37] ICSI Berkeley, "Quicknet." Available: <http://www1.icsi.berkeley.edu/Speech/qn.html>, 2000.
- 374 [38] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François
375 Laviolette, Mario Marchand, and Victor Lempitsky, "Domain-adversarial training of neural
376 networks," *The Journal of Machine Learning Research*, 2016.
- 377 [39] et al. Garofolo, John S., "Csr-i (wsj0) complete ldc93s6a.," *Web Download. Philadelphia:
378 Linguistic Data Consortium*, 1993.
- 379 [40] Markus Kliegl, Siddharth Goyal, Kexin Zhao, Kavya Srinet, and Mohammad Shoeybi, "Trace
380 norm regularization and faster inference for embedded speech recognition rnns," *arXiv preprint
381 arXiv:1710.09026*, 2017.