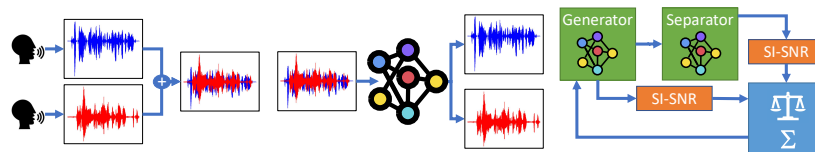


Using generative adversarial networks to make robust speech separation systems

Bc. Ján Pavlus



Abstract

Speech separation is a task of separating single signals from the given mixture of multiple speakers. Speech separation systems are trained on artificial mixtures generated from single speaker's signals. These signals are then used as targets for the training. Neural networks trained this way work well on artificial data but they often fail on real-world examples. To improve their behavior on real-world mixtures it is possible to use training data augmentations for example noise addition. Nevertheless, the power of these augmentations is limited as they have to be manually designed. Using generative adversarial networks (GAN) could improve this process by generating augmentations for data depending on the success of confusing the separation system using these data. Speech separation could be then made more and more robust with each generator and separator training step.

This paper describes experiments that are used to find the right parameters and their combination for the GAN model training. Although the experiments do not yet lead to a more robust speech separation, they provide an analysis of the pitfalls of training the GAN, which is the necessary first step towards a successful system. These experiments show that training the GAN model to the stable state is difficult by adjusting the exact number of batches, after which the separator and generator training is switched. On the other hand, adjusting the to-be-achieved scores of the generator or separator training move could work much better and train the GAN model properly. Other experiments have to be done to prove the correctness of these parameters and their settings.

Keywords: speech-separation — GAN — robust — adversarial augmentations

Supplementary Material: [Downloadable Code](#)

*xpavlu10@stud.fit.vutbr.cz, Faculty of Information Technology, Brno University of Technology

1. Introduction

Speech separation systems are useful as a pre-processing for speech recognition systems which often fail on more overlapped speech. In these cases, the speech separation system could improve the result of the recognition system by separating individual signals from the mixed speech. However, as speech separation systems today are based on neural networks, they need to

be trained on the mixed signals for which single signals are well known. For real-world mixtures, single-speaker signals are usually unavailable, and thus it is necessary to use artificial mixtures. This leads to a problem with bad performance of speech separation systems on the real-world mixtures. This creates the need to make the speech separation systems more robust towards the real-world mixtures.

9
10
11
12
13
14
15
16

17 The more robust speech separation system could
 18 be achieved by using different data augmentations
 19 on the training data. Classic data augmentations are
 20 performed by well-known methods, that are more de-
 21 scribed in Section 3. The disadvantage of these meth-
 22 ods is that they do not cover all possible augmentations
 23 and each new augmentation needs to be manually de-
 24 signed. Generative adversarial networks could be used
 25 to perform data augmentations for speech separation
 26 systems automatically. Their advantage is that they can
 27 generate augmentations depending on the response of
 28 the speech separation system.

29 In this paper, a modified version of the generative
 30 adversarial network is used. It consists of the generator
 31 network generating the augmentations, the separator
 32 network that should be trained to be more robust, and
 33 the similarity loss function that constraints the genera-
 34 tor network. The Separator network and the similarity
 35 loss function represent the discriminator role. For both
 36 networks, the ConvTasNet [1] architecture (with dif-
 37 ferent parameters) has been used.

38 2. Speech separation using neural net- 39 works

39 This section introduces the speech separation task and
 40 the way neural networks are used to tackle it.

41 2.1 Speech separation task

42 The speech separation task could be explained as a
 43 Cocktail party problem. Imagine a cocktail party where
 44 a lot of people talk over each other. The listener present
 45 at the party is trying to focus on one specific speech.
 46 The human ear and brain are well adapted to solve
 47 this task, but for computer systems, it is very difficult.
 48 More formally, there is a mixture defined as:

$$y_t = \sum_{n=1}^N s_{t,n} \quad (1)$$

49 where y_t is the mixture to be separated, $s_{t,n}$ is the
 50 speech signal of a single speaker or noise, t is the
 51 time index, n is the source index, and N is the number
 52 of sources. The main task in speech separation is to
 53 reconstruct signals $s_{t,n}$ from the mixture y_t with no
 54 information about the signals $s_{t,n}$.

55 In the past, there were attempts to solve this task
 56 with classical methods such as principal component
 57 analysis [2] or independent component analysis [3].
 58 These classical methods usually work well when the
 59 task is greatly simplified, but they fail when silent
 60 blocks, echoes, and delays are present. Nowadays
 61 neural networks are used for speech separation tasks

62 and they work well. The most used neural network
 63 architectures are convolutional neural networks and
 64 recurrent neural networks.

65 2.2 Training neural networks for speech sepa- 66 ration

67 Neural network is used to estimate the single signals
 68 from the given mixture. During training, estimated
 69 signals are compared with the known original ones us-
 70 ing the scale-invariant signal-to-noise-ratio (SI-SNR)
 71 function [4], which is defined as:

$$\vec{s}_{\text{target}} := \frac{\langle \hat{s}, \vec{s} \rangle \vec{s}}{\|\vec{s}\|^2} \quad (2)$$

$$\vec{e}_{\text{noise}} := \hat{s} - \vec{s}_{\text{target}} \quad (3)$$

$$\text{SI-SNR}(\vec{s}, \hat{s}) := 10 \log_{10} \frac{\|\vec{s}_{\text{target}}\|^2}{\|\vec{e}_{\text{noise}}\|^2} \quad (4)$$

72 where $\hat{s} \in \mathbb{R}^{1 \times T}$ is the estimated source. $\vec{s} \in \mathbb{R}^{1 \times T}$
 73 is the original source signal used as the target. The
 74 $\|\vec{s}\|^2 = \langle \hat{s}, \vec{s} \rangle$ denotes the signal power, where $\langle \hat{s}, \vec{s} \rangle$
 75 denotes the dot product between estimated and original
 76 source. The function is scale-invariant because the
 77 scale of the estimated signal does not influence the
 78 result. The neural network is trained to maximize this
 79 function.

80 The neural network estimates the N separated sig-
 81 nals from the given mixture to the N outputs. Never-
 82 theless the neural network can output the signals of
 83 individual speakers in any order. This gives rise to a
 84 permutation problem.

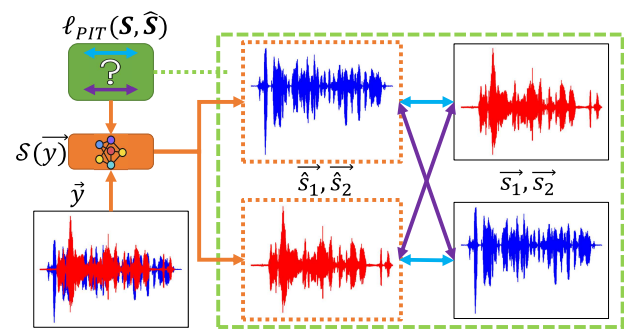


Figure 1. Example of permutation invariant training (PIT) with two estimated outputs by neural network in the orange dotted boxes and two ground truth targets on the right from them. PIT compares estimations and ground truth in all permutations and chooses the permutation with the best result.

85 The solution is the permutation invariant training
 86 (PIT) method [5] shown in Figure 1. This method

87 computes the loss function between all permutations
 88 of original and estimated signals. The best-computed
 89 value corresponds to the right permutation of the es-
 90 timated outputs and this value is also used for the
 91 training of the neural network. This method is defined
 92 as:

$$\hat{s}_1, \hat{s}_2, \dots, \hat{s}_N = \mathcal{S}(\vec{y}) \quad (5)$$

$$l_{PIT}(\mathbf{S}, \hat{\mathbf{S}}) = \min_i \sum_{j=1}^N -\text{SI-SNR}(\vec{s}_{\sigma_{i,j}}, \hat{s}_j) \quad (6)$$

93 where $\mathcal{S}(\vec{y})$ is separator function, that estimates sepa-
 94 rated signal as outputs from the given mixture \vec{y} . These
 95 signals are represented by vectors $\hat{s}_1, \hat{s}_2, \dots, \hat{s}_N$. The
 96 $l_{PIT}(\mathbf{S}, \hat{\mathbf{S}})$ function takes two parameters: \mathbf{S} , which
 97 is matrix of vectors of target signals and matrix $\hat{\mathbf{S}}$,
 98 which consists of vectors of estimated separated sig-
 99 nals. Variable N represents the number of single speak-
 100 ers present in the mixture. Permutation $\sigma_{i,j}$ is the index
 101 of j -th target signal in the i -th permutation of target
 102 vectors given by the matrix \mathbf{S} .

103 3. Robust speech separation

104 To train the speech separation system on real-world
 105 mixtures is really difficult as it is hard to obtain single
 106 speech signals from them. Therefore, the artificial
 107 mixtures are used instead with known single speech
 108 signals.

109 Using solely artificial mixtures during training usu-
 110 ally leads to worse system results on real-world mix-
 111 tures. Real-world mixtures contain echoes, noises, and
 112 reverberations obtained from real-world spaces like
 113 concert halls, public places, congress rooms, airports,
 114 etc. The neural network does not experience all the
 115 variety during training, so it often leads to improper
 116 separation of the speech signals.

117 To prevent the neural network from confusion there
 118 are methods to make the speech separation system
 119 more robust. The common way is to use data augmen-
 120 tations for the training data. There are many classical
 121 data augmentations [6] such as polarity inversion aug-
 122 mentation, frequency filters, decreasing or increasing
 123 the volume, adding noises, adding reverberation etc.

124 3.1 Generative adversarial networks

125 The above mentioned data augmentation methods are
 126 well known and they help the speech separation sys-
 127 tem to work better on real-world data. These methods,
 128 however, can not cover all the features of real-world
 129 mixtures, so there is a potential to use generative ad-
 130 versarial networks (GAN).

The principle of a GAN models is to use two neural 131
 networks, where the first one is used as the generator 132
 and the second one is used as the discriminator. These 133
 two networks are then trained by playing a min-max 134
 game against each other. This principle is slightly mod- 135
 ified in this paper. There are also two neural networks: 136
 the first one is the generator network, which takes an 137
 artificial mixture signal as input and provides an aug- 138
 mented signal, while the second network is the speech 139
 separation neural network which should be made ro- 140
 bust. The basic idea is to make speech separation 141
 more robust epoch by epoch by generating better and 142
 better-generated data augmentations. These generated 143
 augmentations are constrained by maximizing the simi- 144
 larity between the original mixture and the augmented 145
 mixture. 146

The training of GAN model consists of two steps: 147

1. Generator training, which is shown in Figure 2.
 The weights of the separator network are locked.
 Mixtures are given to the generator network,
 which generates the augmented mixtures on the
 output. Augmented mixtures are then submitted
 to a separator network, which provides separated
 signals. The l_{PIT} loss function which is defined
 in Equation 6 is then computed between the
 separated signals and the target ones while the
 generator is trained to maximize it. In parallel,
 the SI-SNR is computed between the augmented
 mixtures and the original ones, which gives a
 similarity value the generator tries to maximize.
 This is necessary to prevent the generator from
 completely destroying the information in the
 mixture. The loss value used for training of the
 generator network is computed as a weighted
 sum of the similarity value and the separator loss
 value, which is defined as:

$$\vec{g} = \mathcal{G}(\vec{y}) \quad (7)$$

$$x_{\text{sim}} = -\text{SI-SNR}(\vec{y}, \vec{g}) \quad (8)$$

$$x_{\text{sep}} = l_{\text{pit}}(\mathbf{S}, \mathcal{S}(\vec{g})) \quad (9)$$

$$l_{\text{gen}}(\vec{g}, \vec{y}, \mathbf{S}) = -w_{\text{sep}} * x_{\text{sep}} + w_{\text{sim}} * x_{\text{sim}} \quad (10)$$

where \vec{g} is the generated augmented mixture by 148
 the generator function $\mathcal{G}(\vec{y})$, which takes mix- 149
 ture \vec{y} as an input. The x_{sim} is the similarity value 150
 computed by the SI-SNR function between the 151
 generated augmented mixture \vec{g} and the original 152
 mixture \vec{y} . The x_{sep} is the value computed by 153
 the PIT loss function between the target signals 154
 in the matrix \mathbf{S} and the separated signals by the 155
 separator function $\mathcal{S}(\vec{g})$, that takes the gener- 156
 ated augmented mixture \vec{g} as an input. Then 157

158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175

the generator loss function $(\vec{g}, \vec{y}, \mathbf{S})$, which takes the generated augmented mixture \vec{g} , the original mixture \vec{y} and the matrix of the target signals \mathbf{S} as an input is computed by the weighted sum of the similarity x_{sim} and the x_{sep} values. The weights are set by the parameters w_{sep} , which sets the importance of the separator loss, and w_{sim} , which sets the importance of the similarity value in the generator loss function.

2. Separator training, which is shown in Figure 3. The weights of the generator network are locked in this step. The separator neural network is trained in the classical speech separation way, but it receives a certain amount of augmented mixtures, for example, 50%. The ratio of the augmented mixtures will be denoted as r_{aug} . In this step, the separator network is trained to handle the augmented mixtures to be more robust.

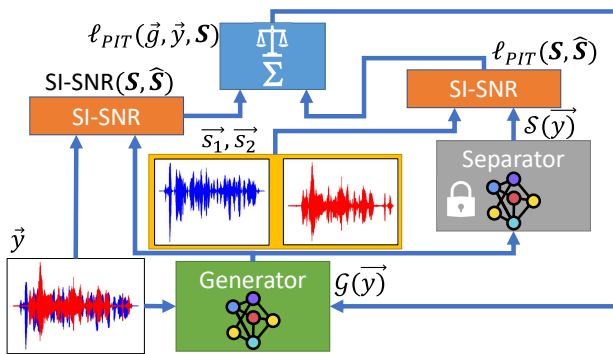


Figure 2. Architecture of generative adversarial network training used in this paper. This figure shows the step where the generator is trained.

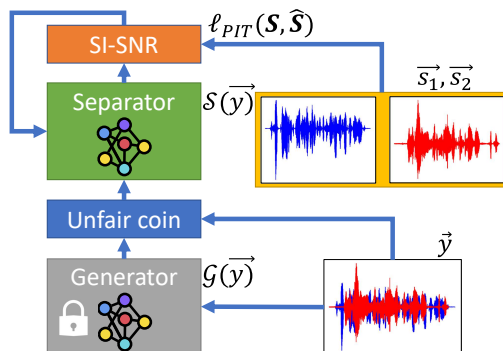


Figure 3. Architecture of generative adversarial network training used in this paper. This figure shows the step where the separator is trained.

3.2 Problems with generative adversarial networks 176 177

There are several problems with training GAN models [7]. The first of them is non-convergence, which means that the GAN weights oscillate and never converge to the one best state. The second problem is mode collapse, where GAN is collapsing to the few generating modes. For example, the model which is trained to generate numbers only generates numbers two and five. Another problem is called diminished gradient. In this problem, the discriminator gets too successful so the generator is unable to learn anything. In this paper, we experienced a problem with the imbalance between the separator and generator networks. The imbalance problem and the problem of finding the right parameters are described in Section 4 in a more detailed manner. 178
179
180
181
182
183
184
185
186
187
188
189
190
191
192

Another problem that appears when training the separator together with the generator is the right choice of the best model. In the classical neural networks training methods, the best model is chosen depending on the best cross-validation loss value during the training. This is not applicable in this paper. It is necessary to find the right separator model, which works well on both the original and the augmented data. The generated augmentations are changing during the training due to the changes in the generator. Therefore, it is not possible to compare the validation loss values through separator models from different epochs. This causes a problem with the best separator network selection. 193
194
195
196
197
198
199
200
201
202
203
204
205

The solution is to save generator and separator models from each epoch during training. They can be used to generate augmented mixtures by randomly choosing a generator from a uniform distribution for each mixture of the cross-validation set. Thus augmented mixtures now represent all possible augmentations learned during training. Now it is necessary to choose the right separator, which will be the most robust one. This task could be achieved by evaluating each separator on mentioned augmented mixtures. The separator with the best evaluation result should be the most robust one. 206
207
208
209
210
211
212
213
214
215
216
217

4. Experiments 218

Experiments with a generative adversarial network used to make the speech separation more robust are performed with the setup described below. The first experiments try to find the applicable combinations of parameters, which do not lead to one of the GAN training problems. 219
220
221
222
223
224

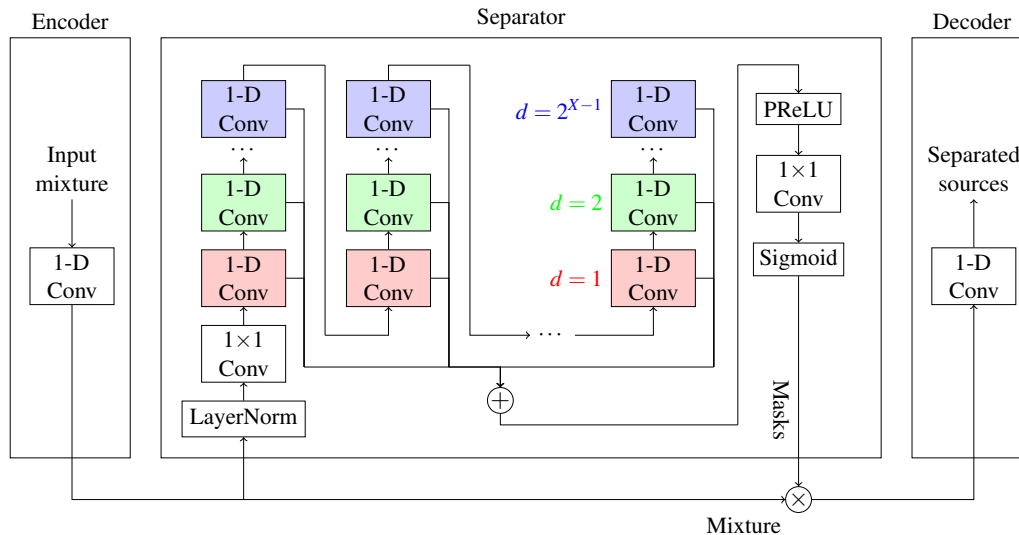


Figure 4. ConvTasNet neural network architecture. Image adapted from [1]

225 4.1 Dataset

226 Experiments are provided on the Wall Street Journal
 227 dataset (WSJ) [8]. It consists of three parts, which
 228 contain training, cross-validation, and testing data.
 229 The dataset contains both mixtures and parallel single-
 230 speaker recordings. Speakers are randomly mixed with
 231 various signal-to-noise ratios (SNR) between 0 dB and
 232 5 dB. For training there are 20000 mixtures correspond-
 233 ing to 30 hours, for cross-validation, there are 5000
 234 mixtures corresponding to 10 hours, and there are 3000
 235 mixtures corresponding to 5 hours for testing.

236 The second dataset used in the experiment is WSJ0
 237 Hipster Ambient Mixtures (WHAM) [9]. This dataset
 238 pairs each two-speaker mixture in the WSJ dataset
 239 with the unique noise background scene. The sizes
 240 of training, cross-validation, and testing parts are the
 241 same as in the raw WSJ dataset.

242 4.2 ConvTasNet

243 The neural network architecture used for both networks
 244 (separator and generator) is ConvTasNet [10]. This
 245 architecture consists of three parts as it is shown in
 246 Figure 4.

247 The first part is the encoder which consists of one
 248 convolutional layer. This layer takes the mixture signal
 249 as an input and provides pseudo short-time Fourier
 250 transform (STFT) transformation of the given signal.
 251 An encoded signal is then given to the second part of
 252 the architecture.

253 This part is called separator and provides separa-
 254 tion as the name reveals. The separation part consists
 255 of the blocks of the convolutional layers that are ap-
 256 plied to the larger and larger context. The output of

257 this part is a series of masks for each speaker. These
 258 masks are applied to the encoded signal and given to
 259 the last part, decoder, as an input.

260 The decoder part consists of one convolutional
 261 layer like the encoder part. The task of this part is to
 262 reassemble the signal from the encoded pseudo STFT
 263 format. Parameters of the ConvTasNet are shown in
 264 Table 1. The table also shows the parameter values for
 265 the generator and separator neural networks used in
 266 the experiments.

267 4.3 Initial separator and generator networks 268 setups

269 The separator network used for training has been pre-
 270 trained on one of the above mentioned datasets or their
 271 combination. So there are three separator networks
 272 used in experiments, each pretrained on one of them.
 273 The baseline pretrained scores are shown in Table 2.
 274 From the given results it is obvious that the WHAM
 275 dataset is much more difficult than the WSJ. This is
 276 caused by the noises added to the mixtures. If the sys-
 277 tem is trained on the WSJ and tested on the WHAM,
 278 then the results are quite poor.

279 The generator network is much smaller than the
 280 separator one and it has been pretrained for the self-
 281 identity task. Since the aim of the GAN training is
 282 not to train the encoder and decoder parts of the Con-
 283 vTasNet, pretraining the self-identity task removes the
 284 encoder and decoder training problem.

285 4.4 Adjustable parameters

286 The first task of the experiments is to find the right
 287 combination of parameters that will train GAN prop-
 288 erly. The adjustable parameters are:

Table 1. Hyperparameters of the ConvTasNet network [1]

Symbol	Description	Generator	Separator
F	Number of filters in autoencoder	128	128
L	Length of the filters (in samples)	40	40
B and the residual paths' 1×1 -conv blocks	Number of channels in bottleneck	128	128
H	Number of channels in convolutional blocks	192	192
P	Kernel size in convolutional blocks	3	3
X	Number of convolutional blocks in each repeat	3	7
R	Number of repeats	1	3
O	Number of outputs	1	2

Table 2. Baseline results of pretrained separator neural networks. Results are computed by SI-SNR loss function using PIT method. Datasets in rows are the training ones. Testing datasets are in columns.

	WSJ	WHAM
WSJ	12.46	-2.99
WHAM	9.04	6.09
WSJ + WHAM	12.34	6.45

r_{aug} and the similarity loss SI-SNR cap c_{sim} , which sets the value that is used to clip the similarity loss to a maximum value. This serves to prevent the similarity function to be too strong in comparison with the separator loss function results.

4.5 Initial experiment

The initial experiment is set with following parameters:

- w_{sep} and $w_{\text{sim}} = 1.0$,
- c_{sep} and $c_{\text{gen}} = 10$,
- $r_{\text{aug}} = 0.5$ and
- $c_{\text{sim}} = 40$.

The purpose of the initial experiment was to inspect the basic behavior of the loss functions during the training. The results are shown in Figure 5. The separator loss function curve shows that the generator network managed to completely confuse the separator network. Although this is the task of the generator, in this case, the generator completely dominated the training to the point that the separator was unable to adapt to the augmented mixtures. The strength of the generator network is possibly caused by:

1. Too much emphasis on the separator network confusion, which could be adjusted by the parameters w_{sep} and w_{sim} . These adjustments will be examined in Section 4.6, or
2. Too much training space for the generator network, which could be adjusted by the parameters c_{gen} and c_{sep} . These adjustments will be examined in Section 4.7.

4.6 Generator loss weights

The base experiments lead to the imbalance between the generator and separator network. There is a chance to solve this imbalance problem by finding the correct weights w_{sep} and w_{sim} .

Therefore, other experiments are set with different combinations of values of the weight parameters. Thus the separator weight w_{sep} value is reduced by tenths to 0.1 with similarity weight w_{sim} locked at 1.0. This

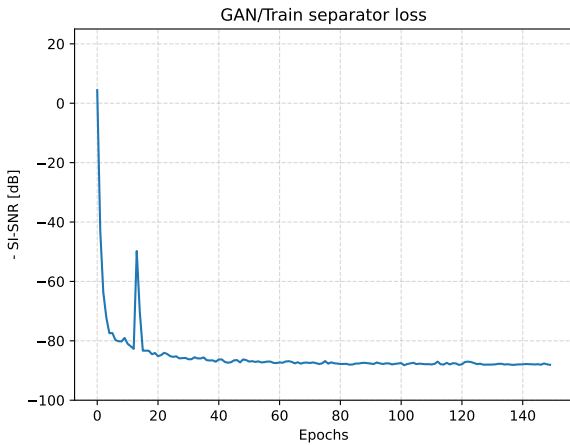
- Separator loss weight w_{sep} , which sets the importance of the separator loss in generator training. The generator training loss function is defined by Equation 10. The separator loss value is computed during generator training on the generated augmented mixtures. The generator is trained to maximize the separator loss in order to confuse the separator as much as possible.
- Similarity loss weight w_{sim} is used to indicate the importance of the similarity between the generated augmented mixture and the original one. This loss function is also computed during training and its role is to constrict the generator so that it would not generate complete nonsense.

GAN model is switching between the separator and generator training during each epoch. Two parameters control this:

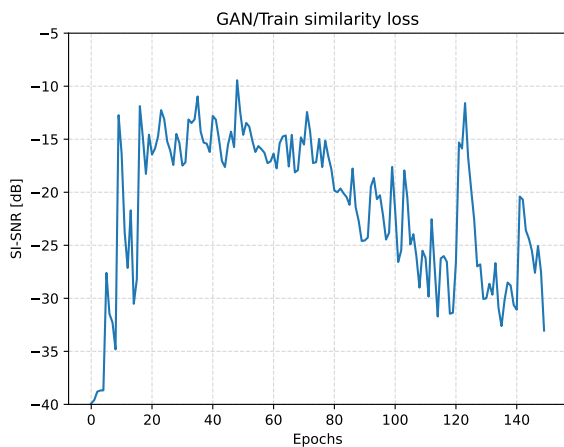
- Separator batch cap c_{sep} , which sets how many batches will be used in separator training turn.
- Generator batch cap c_{gen} , which sets how many batches will be used in generator training turn.

For example, when c_{sep} and c_{gen} are set to 10, the generator will be trained on the first ten batches. After this, the training is switched to the separator training, which uses other ten batches and then switches back. The number of batches for each model can significantly influence the training and these parameters are difficult to set properly.

The last two adjustable parameters are the ratio of the augmented mixtures during the separator training

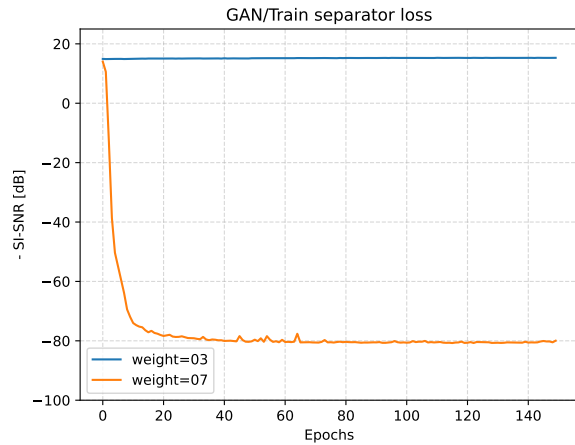


(a) Training curve of the separator loss function

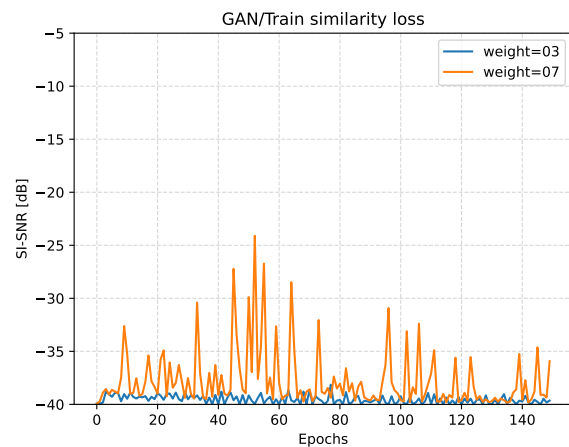


(b) Training curve of the similarity loss function

Figure 5. Training curves of the separator and similarity loss computed during the generator training move each epoch. These curves shows collapse of the GAN model with initial parameters setting to the imbalance state, where the generator is too strong for the separator.



(a) Training curve of the separator loss function



(b) Training curve of the similarity loss function

Figure 6. Training curves of the separator and similarity loss computed during the generator training move each epoch. These curves show the collapse of the GAN model with $w_{sep} = 0.7$ to the imbalance state, where the generator is too strong for the separator.

357 could reduce the generator strength and help to a better
358 system balance.

359 Experiments collapse to two modes, where the
360 generator network:

- 361 1. Is too strong and overwhelms the separator net-
362 work, or
- 363 2. Generates very similar mixtures to the original
364 ones and does not make any changes.

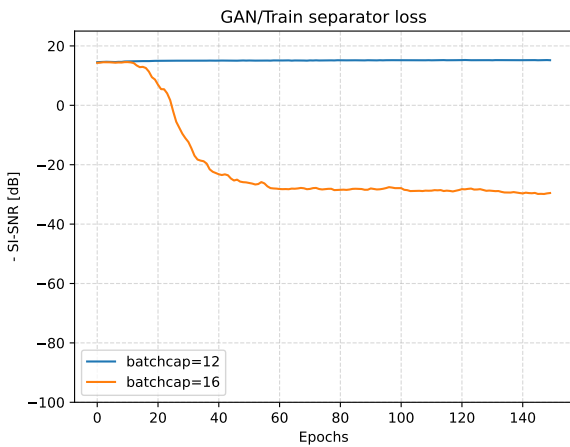
365 The first mode is achieved when the c_{sep} is above
366 the value 0.5 as shown by the orange curves in Fig-
367 ure 6. Lower values collapse to the second mode,
368 where the similarity loss function has a big influence
369 on the generator as it is shown by the blue curves in
370 Figure 6.

4.7 Separator and generator batch caps

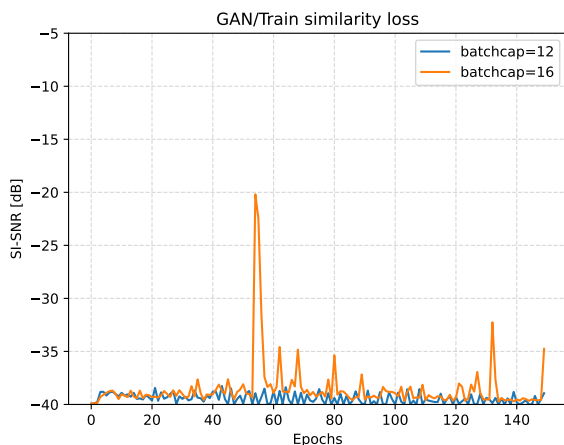
371

Another idea is to solve the collapsing to the first mode
by the right combination of the separator and genera-
tor batch cap parameter values. The main issue is that
the generator confuses the separator too much. There-
fore, increasing the separator batch cap c_{sep} could help
the separator to better adapt to the augmented mix-
tures. The fixed value $w_{sep} = 0.7$ was used in batch
cap experiments. This value has been chosen because
although the system with these settings collapses to
the first mode, it reduces the impact of the separator
loss function on the generator. The c_{sep} value is incre-
mented by the unit. Nevertheless, the system does not
stabilize again, it collapses to the second mode when
the $c_{sep} \geq 13$ as it is shown by the blue curves in
372
373
374
375
376
377
378
379
380
381
382
383
384
385

386 Figure 7. With lower values of c_{sep} , the system stands
 387 in the first collapse mode as it is shown by the orange
 388 curves in Figure 7. Therefore, it means that the batch
 389 cap parameters are not distinguished finely enough.



(a) Training curve of the separator loss function



(b) Training curve of the similarity loss function

Figure 7. Training curves of the separator and similarity loss computed during the generator training move each epoch. The blue curves show the collapse of the GAN model with $c_{sep} = 12$ to the imbalance state, where the generator does not generate any augmentations. The orange curves show the collapse of the GAN model with $c_{sep} = 13$ to the imbalance state, where the generator is too strong for the separator.

390 4.8 Automatic separator and generator batch 391 caps

392 After these experiments, it turned out that adjusting the
 393 batch caps c_{sep} and c_{sim} is not enough to stabilize the
 394 training. The constant batch cap values still lead to one
 395 of the two collapse modes described in Section 4.6,
 396 i.e. one of the models is too strong while the other
 397 does not learn anything. Here, we explore another way

to balance the training, where the number of batches 398
 for each model is adjusted dynamically, based on the 399
 SI-SNR value of the separator. The generator is thus 400
 trained until the separator gains loss values higher than 401
 parameter $c_{snr_{gen}}$ and the separator is trained until the 402
 separator does not achieve a loss value higher than 403
 parameter $c_{snr_{sep}}$ on the augmented mixtures. 404

Experiments using this method are initially set 405
 with parameters: 406

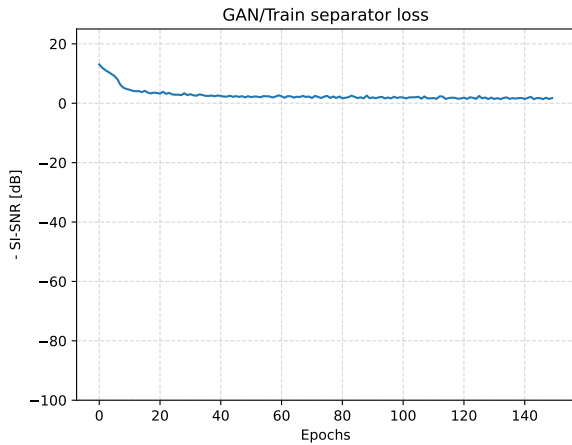
- $w_{sep} = 0.6, 0.7, 0.9$ 407
- $w_{sim} = 1.0,$ 408
- $c_{snr_{gen}} = 0$ 409
- $c_{snr_{sep}} = 5.0,$ 410
- $r_{aug} = 0.5$ and 411
- $c_{sim} = 20.$ 412

The c_{sim} value follows the knowledge from the pre- 413
 vious experiment, where the similarity values around 414
 the 40dB overweight values of the separator loss func- 415
 tion during the generator training. These experiments 416
 use the pretrained separator model on the WSJ dataset. 417
 From training curves shown in Figure 8 it is evident, 418
 that systems trained by using this method do no longer 419
 collapse to the modes mentioned in Section 4.6. The 420
 training curve that stands for the level of the separator 421
 confusion converges to the value set by the parame- 422
 ter $c_{snr_{gen}}$. The second chart shows that the training 423
 curve of the separator gained values on the augmented 424
 data, which converges to the value set by the parameter 425
 $c_{snr_{sep}}$. 426

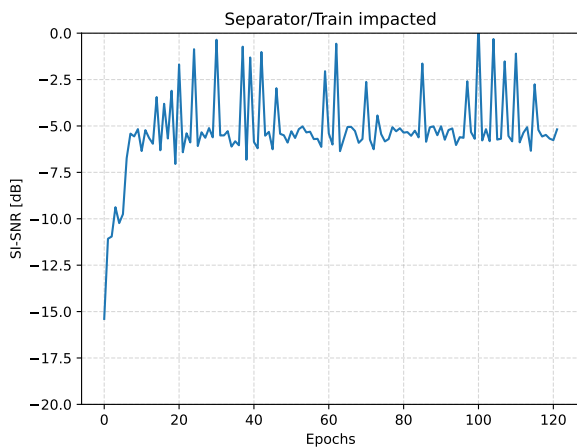
4.9 Final results 427

The best-trained separator model has to be found. This 428
 is provided by generating augmented mixtures by ran- 429
 domly chosen trained generators. Then each tenth sep- 430
 arator model is evaluated on all generated augmented 431
 mixtures. The model with the best result is then cho- 432
 sen for the evaluation using the test data. The results 433
 of the separator selection are shown in Figure 9. The 434
 first chart shows the overall results of each evaluated 435
 separator model. The bar charts show how the separa- 436
 tor was successful on the generated augmentations by 437
 each generator. Bars represent groups of generators, 438
 which are grouped by their epoch number. 439

The separator with the best score is chosen as the 440
 best one. This separator is then evaluated on the WSJ 441
 and WHAM datasets. The results are shown in Table 3. 442
 They show that the trained GAN model with the above 443
 mentioned parameters setting does not train the separa- 444
 tor to be more robust. This is tested by evaluating the 445
 best separator model on the testing part of the WHAM 446
 dataset. The SI-SNR result achieved by the separator 447
 model pretrained on the WSJ dataset is similar to the 448



(a) Training curve of the separator loss function



(b) Training curve of the similarity loss function

Figure 8. The training curve of the separator loss computed during the generator training moves each epoch and the training curve of the computed separator loss function during the separator training on the augmented mixtures. These curves show, that the curves converge to the set parameters c_{snrgen} and c_{snrsep} . The c_{snrsep} parameter value is inverted during the training.

Table 3. Final results of the experiments with automated c_{gen} and c_{sep} parameter settings. In first column there are different w_{sep} parameter settings. Other columns represent results from the evaluation on the tested part of the WSJ or WHAM dataset. The columns with original annotation contain evaluation results of the raw pretrained separator model on the WSJ dataset. The columns with the augmented annotation contain the evaluation results of the best separator model chosen from the GAN training.

	WSJ original	WSJ robust	WHAM original	WHAM robust
$w_{\text{sep}} = 0.6$	12.46	12.18	-2.99	-2.89
$w_{\text{sep}} = 0.7$	12.46	11.95	-2.99	-2.78
$w_{\text{sep}} = 0.9$	12.46	11.36	-2.99	-3.05

The main obstacle to training such a system is finding the correct parameters. These parameters should be chosen experimentally. The presented experiments show, that the system collapses to the two modes. Other experiments are set to solve this collapsing show that adjusting exactly the amount of the training space, which is given for both separator and generator networks during the training epoch is ineffective. Therefore, it is better to use the automated adjusting method of the amount of the training space. The automated method sets the goals of the generator and separator networks that should be achieved by each of them during their training turn.

The model using this automated method no longer collapses to one of the instability modes. However, as the final results reveal it does not train the separator network to be more robust properly. Further experiments have to be done to find the right parameters combination, that will train the separator network to be more robust. Adjusting the weight of the similarity function or different values of the generator and separator goals parameters could improve these results.

If such a parameters combination will be discovered, it will fit only for the current model with current settings and the current dataset. If there will be any change in of those three things, it is necessary to find the right parameters combination again. Therefore, the GAN parameters are very sensitive, it is very difficult to find such a combination. Nevertheless it is possible to use some algorithms to find the right parameters such as evolutionary algorithms [11] or bayesian hyperparameter optimization [12]. If these methods will work how they should then the using a GAN to make speech separation system robust could be used widely.

449 model trained by the GAN, when evaluating on both
450 the WSJ and WHAM datasets. There are three exper-
451 iments using different separator weights, but none of
452 them has achieved better results. Further experiment-
453 ing is to be done to examine if there is any parameters
454 settings combination that leads to better final results.

455 5. Conclusions

456 This paper investigates the usage of generative net-
457 works to automatically augment training data for speech
458 separation systems. This could replace manually de-
459 signed data augmentation methods and make the speech
460 separation system more robust.



Figure 9. Results from finding the best separator model from trained GAN with parameter $w_{\text{sep}} = 0.7$. The first chart shows the SI-SNR means of evaluated separators. The other charts show the results on generated augmented mixtures of the separator models from the each selected epoch

495 References

- 496 [1] Yi Luo and Nima Mesgarani. Conv-tasnet: Sur- 507
 497 passing ideal time–frequency magnitude mask- 508
 498 ing for speech separation. *IEEE/ACM transac- 509
 499 tions on audio, speech, and language processing*, 27(8):1256–1266, 2019. 510
 500
 501 [2] Hervé Abdi and Lynne J Williams. Principal 513
 502 component analysis. *Wiley interdisciplinary re- 514
 503 views: computational statistics*, 2(4):433–459, 2010. 515
 504
 505 [3] James V Stone. Independent component analysis: 516
 506 an introduction. *Trends in cognitive sciences*, 6(2):59–64, 2002. 517
 518
 [4] Fahimeh Bahmaninezhad, Jian Wu, Rongzhi Gu, 508
 Shi-Xiong Zhang, Yong Xu, Meng Yu, and Dong 509
 Yu. A comprehensive study of speech separation: 510
 spectrogram vs waveform separation. *arXiv 511
 preprint arXiv:1905.07497*, 2019. 512
 [5] Dong Yu, Morten Kolbæk, Zheng-Hua Tan, and 513
 Jesper Jensen. Permutation invariant training 514
 of deep models for speaker-independent multi- 515
 talker speech separation. In *2017 IEEE Interna- 516
 tional Conference on Acoustics, Speech and Sig- 517
 nal Processing (ICASSP)*, pages 241–245. IEEE, 518

- 519 2017.
- 520 [6] Tom Ko, Vijayaditya Peddinti, Daniel Povey, and
521 Sanjeev Khudanpur. Audio augmentation for
522 speech recognition. In *Sixteenth annual confer-*
523 *ence of the international speech communication*
524 *association*, 2015.
- 525 [7] Divya Saxena and Jiannong Cao. Generative
526 adversarial networks (gans) challenges, solutions,
527 and future directions. *ACM Computing Surveys*
528 (*CSUR*), 54(3):1–42, 2021.
- 529 [8] J. R. Hershey, Z. Chen, J. Le Roux, and S. Watan-
530 abe. Deep clustering: Discriminative embed-
531 dings for segmentation and separation. In *2016*
532 *IEEE International Conference on Acoustics,*
533 *Speech and Signal Processing (ICASSP)*, pages
534 31–35, 2016.
- 535 [9] Gordon Wichern, Joe Antognini, Michael Flynn,
536 Licheng Richard Zhu, Emmett McQuinn, Dwight
537 Crow, Ethan Manilow, and Jonathan Le Roux.
538 Wham!: Extending speech separation to noisy
539 environments. In *Proc. Interspeech*, September
540 2019.
- 541 [10] Jian Wu. A PyTorch implementation of the Con-
542 v-tasnet: Surpassing Ideal Time-Frequency Mask-
543 ing for Speech Separation, November 2019.
- 544 [11] Steven R Young, Derek C Rose, Thomas P
545 Karnowski, Seung-Hwan Lim, and Robert M Pat-
546 ton. Optimizing deep learning hyper-parameters
547 through an evolutionary algorithm. In *Proceed-*
548 *ings of the workshop on machine learning in high-*
549 *performance computing environments*, pages 1–5,
550 2015.
- 551 [12] Franck Dernoncourt and Ji Young Lee. Optimiz-
552 ing neural network hyperparameters with gaus-
553 sian processes for dialog act classification. In
554 *2016 IEEE Spoken Language Technology Work-*
555 *shop (SLT)*, pages 406–413. IEEE, 2016.