

Master Thesis

Gaspard Ulysse Fragnière

August 2022

1 State-of-the-art

The problem of Acoustical Source Localization (ASL) is an important problem. It has many applications such as smart assistants (e.g. Google Home, Alexa, ...), industrial applications, **TODO: add more?**. Traditionally this problem is tackled with methods based on the physics of sound propagation (e.g. TDoA, beamforming) or with statistical inference (e.g. Sparse Bayesian Learning).

The recent success of Deep Learning (DL) based methods in other fields of research (e.g. Computer Vision) led to believe that Deep Neural Networks (DNN) based approaches could provide state-of-the-art results in solving the ASL problem. Castellini et al. (2021), Kujawski et al. (2019), Lee et al. (2021), Ma and Liu (2019), Pinto et al. (2021) and Xu et al. (2021) propose state-of-the-art DL-based methods for Source Characterization.

A common issue faced while implementing DL-based methods is that significant quantities of well-structured data are required. In the literature, the data has been obtained using the following approaches:

- **Real Measurement:** To create the different samples of such a dataset, sounds emitted with a loudspeaker or human voices are recorded in a real acoustic environment. Even though such a method allows for the creation of perfectly realistic samples, it does not come without any issues. Indeed, it is very tedious and time-consuming to record in different environments. Additionally, all the environments for measurement need to physically exist, which limits the quantity of possible samples. Moreover, to build a high-quality data set, expensive equipment is required to have an accurate ground truth (i.e. precisely identify the location of the sources). In the literature, He et al. (2018) and Ferguson et al. (2018) have used such an approach.
- **Synthetic Data:** The sounds used are artificial (i.e. white noise, sine wave). The room acoustic is also simulated. Indeed, the dry sound is convolved with a simulated Room Impulse Response (RIR) to mimic the effect of room acoustics (e.g. reverberation). Compared to real measurement, this approach allows samples in more diverse environments. Indeed, RIR for rooms of arbitrary size, different source positions as well as different dry signals can be used for the training. The issue with such a method is the important amount of time and storage required for the creation of the datasets. E.g. Chakrabarty and Habets (2017), Perotin et al. (2018) and Adavanne et al. (2018) created their datasets in this way.
- **Semi-synthetic data:** The creation of such a dataset is similar to the creation of synthetic datasets. The difference lies in the fact that the dry sound source used and the RIR are measured and not simulated. Then, the samples of such a dataset are generated by convolving

dry sounds with RIR. This method is not the best suited, since it is very time-consuming to generate a data set with enough samples for training a DL-based algorithm. Indeed, measuring all the RIR lead to the issues faced with real measurement. Takeda and Komatani (2016) use such an approach for to obtain their data.

Moreover it is to be noted that none of these methods are suitable for online data generation. Indeed, any of the above mentionned method do not allow for creating random sample while training DL-based algorithm. To use such datasets for training, they need to be fully created (and stored) before any training can occur.

1.1 DL-based data generation

In the past years, DL-based approaches have shown to be able to learn and realistically reproduce very complicated data structures (e.g. generation of pictures of faces in the field of Computer Vision). Those breakthroughs lead to believe that similar data generation methods could be used to fix the above-mentionned issues (e.g. offline training, lack of variance in the different samples, ...).

Moreover, it is relevant to note that the data used for source characterization in Castellini et al. (2021), Lee et al. (2021), Ma and Liu (2019), Xu et al. (2021) is the Cross Power Spectra (CPS), i.e. a direct representation of the signals received in the array of microphone. Indeed those approaches do not use direct recording of microphone input but instead features extracted from the raw data. This is crucial because it means that recording, simulating or generating raw microphone data is no longer necessary, if features (e.g. CPS) could be generated directly. We therefore need to identify what acoustic quantities:

- have already been generated using a DL approach
- are potential feature for a Source Characterization Algorithm.

1.1.1 Generation of Signal

Neekhara et al. (2019), Kumar et al. (2019), Engel et al. (2019) use Generative Adversarial Network (GAN) to generate realistic audio waveform. Neekhara et al. (2019) and Kumar et al. (2019) specifically focus on the generation of audio waveform conditioned on a spectrogram (cGAN). On the other hand, Engel et al. (2019) design a GAN to generate realistic audio waveform of single music notes played by an instrument. The data generated in those approaches is single-channel data, but maybe it could be extended to multi-channel to simulate the different signals recorded in an array of microphone. It is relevant to note that the GAN designed by Neekhara et al. (2019) is the one implemented in Vargas et al. (2021) in order to compare the accuracy of a network for single source DoA estimation when trained with different sound classes.

1.1.2 Generation of Impulse Response

Papayiannis et al. (2019) introduce a GAN approach to generate artificial Acoustical Impulse Response (AIR) of different environment in order to generate data for a NN used for classification of acoustic environment.

Ratnarajah et al. (2021) proposes a fast method (NN-based) for generating Room Impulse Response (RIR). The input parameters of the networks used for creating the IR are the desired dimensions of the rectangular room, listener position, speaker position and reverberation time (T_{60}).

TODO: is it worth mentionning both papers ?

This is relevant for a problem at hand because if we are to be able to generate impulse responses with known source and listener position, we could simply convolve them with the dry source sounds. This way, we could generate raw microphone signal and use them to train a DL-based algorithm for source characterization.

1.1.3 Generation of potential NN feature

Bianco et al. (2020) proposes an approach to generate another acoustic feature: the phase of the relative transfer function (RTF) between two microphones. In this paper a Variational Auto Encoder (VAE) is designed to simultaneously generate phases of RTF and classifying them by their Direction of Arrival (DoA).

Gerstoft et al. (2020) use a GAN to generate Sample Cross Spectra Matrices (CSM). for a given DoA. In their approach, the GAN is trained with data only coming from one DoA, making it unable to generate sample for different DoA. This approach could be extended by creating a conditional Generative Adversarial Network (cGAN) taking as input the DoA. Such a GAN would receive a DoA as input and use it to produce a CSM corresponding to the received DoA.

1.1.4 Other possible approaches to generate the data

In Hübner et al. (2021) introduce a low complexity model-based method for generating samples of microphones phases. This method proposed is not based on DL. Indeed, it is based on a statistical noise model, a deterministic direct-path model for the point source, and a statistical model. The claim of this paper is that the low complexity of the proposed model makes it suited for online training data generation.

Vera-Diaz et al. (2021) introduce a CNN for denoising (i.e. removing the effects of reverberation and multipath effects) on the Generalization Cross Correlation (GCC) matrix of an array of microphone. More specifically than a CNN, the network used has a encoder-decoder structure. This means that a possible approach to create the data we want, would be to attempt to invert network proposed. With this we could realistically add noise to GCC matrices and hence making it suitable for training.

2 Fundamentals

2.1 GAN

Goodfellow et al. (2020) introduce a new approach to solve the problem of generative model. The goal is to learn the probability distribution that was used to generate samples, by observing them.

The method introduced in this paper is called Generative Adversarial Network (GAN). The idea is to create a game (in the game theory sense) where two networks compete against each other. The first network is called a discriminator and its goal is to determine real from fake samples. The other network is a generator with the aim to produce data realistic enough that the discriminator cannot determine that it is fake.

The generator takes as input a random vector and use it generate a sample. This random vector is referred to as latent variable. The vector space of latent variable is called latent space. After training, the generator should be a mapping from the latent space to the data space. The latent space is a representation of smaller dimension of the data space.

The discriminator takes as input a real or generated sample and predicts its authenticity. The discriminator is a simple binary classifier. The real sample come from the datasets and the fake are output of the generator.

Both models are trained simultaneously. First the generator create a batch of fake sample. This batch is then fed to the discriminator, alongside a batch of real samples. For each of them, the discriminator makes a prediction about their authenticity. The discriminator then gets updated based on how accurate it was at classifying samples and the generator based on how many times fake sample were able to fool the discriminator. We can see here that the training of both networks is supervised, on the contrary of typical generative models.

In a gaming theory framework, both networks are competing in a zero-sum game. This means that if one network performs well at its task and gets rewarded by little weights update, the other must have performed poorly and hence gets penalized by heavy weights updated. E.g. if the discriminator was successful at classifying all samples, it means that the generator had not been able to fool the discriminator by producing any realistic fake samples.

TODO: include more maths and/or more graphs?

2.2 DCGAN

Radford et al. (2015) introduce GAN, but unlike in Goodfellow et al. (2020), the architecture of the discriminator and generator is not achieved with regular perceptron, but with Convolutional layer. We first remind how a regular perceptron layer works. For an input vector $\mathbf{x} \in \mathbb{R}^k$ and output vector $\mathbf{y} \in \mathbb{R}^l$, a perceptron layer consists of two parts:

- an activation $\mathbf{a} = \mathbf{W}\mathbf{x} + \mathbf{b}$
- a non-linearity $y = f(\mathbf{x}; \theta) = \sigma(\mathbf{a})$

where $\mathbf{W} \in \mathbb{R}^{l \times k}$ are the weights of the perceptron and $\mathbf{b} \in \mathbb{R}^l$ its bias. $\sigma(\cdot)$ is called the activation function and is the source of non linearity in neural networks. An example of such activation function is the Rectified Linear Unit (ReLU) where:

$$\sigma_{\text{ReLU}}(\mathbf{a}) = [\max(a_0, 0), \dots, \max(a_{l-1}, 0)]^T \quad (1)$$

We can now introduce the architecture of a convolutional layer. Without loss of generality, we assume that our network receive as input an square image with a single channel, i.e. $\mathbf{H}_0 \in \mathbb{R}^{d_0 \times (m \times m)}$ with $d_0 = 1$, transforms it with one square kernel, i.e. $\mathbf{K} \in \mathbb{R}^{s \times (r \times r)}$ with $s = 1$ to output another square image with a single channel, i.e. $\mathbf{H}_1 \in \mathbb{R}^{d_1 \times (n \times n)}$, with $d_1 = 1$

The relationship between input image $\mathbf{H}_0 \in \mathbb{R}^{m \times m}$, output image $\mathbf{H}_1 \in \mathbb{R}^{n \times n}$ and kernel $\mathbf{K} \in \mathbb{R}^{r \times r}$ is defined as:

- an activation $\mathbf{A} = \mathbf{H}_0 * \mathbf{K}$
- a non-linearity $\mathbf{H}_{1(i,j)} = \sigma(\mathbf{A}_{(i,j)})$

Radford et al. (2015) proposes a set of constraint that makes the use of convolutional layers possible in GAN setting.

2.3 WGAN

Arjovsky et al. (2017) introduce a new type of Generative Adversarial Network, namely the Wasserstein GAN (WGAN). The claim is that WGAN improves the stability in learning and get rid of typical problem of the traditional GAN approach such as Mode Collapse or Convergence failure.

TODO: here need to explain mode collapse ?

More specifically, Arjovsky et al. (2017) provides the following insights:

TODO: Here needs to state the problem that the Earth Mover distance (distance between probability) try to solve.

- Analyses how the Earth-Mover (EM) distance, also known as Wasserstein distance behaves compared to other distance between probability distribution (e.g. Kullback-Leibler distance)
- Define a GAN that minimizes the an approximation of the EM distance, namely the WGAN.
- Show that unlike traditional GANs, WGANs do not need to maintain a balance when training the discriminator and generator. Indeed in regular GAN approach, it was crucial to avoid the discriminator to become too good before the generator, since this would prevent the generator to learn any distribution.

2.3.1 The Earth-Mover or Wasserstein distance

The goal of WGAN, remains the same as GAN, namely solving the problem of generative model. More precisely, this mean approximating the probability distribution P_r of some data by a distribution P_θ . For this reason, it is necessary to have metrics to quantify distance between two probability distributions P_r and P_m . Example of such distances are the Kullback-Leibler divergence or the Jensen-Shannon divergence. In WGAN, the distance used is the Earth-Mover (EM) distance or Wasserstein-1, defined as

$$W(\mathbb{P}_r, \mathbb{P}_m) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_m)} \mathbb{E}_{(x,y) \sim \gamma} [|x - y|] \quad (2)$$

Where $\Pi(\mathbb{P}_r, \mathbb{P}_m)$ is the set of all joint distribution $\gamma(x, y)$ whose marginals are respectively \mathbb{P}_r and \mathbb{P}_m . Informally, $\gamma(x, y)$ shows how much "mass" must be carried to transform \mathbb{P}_r into \mathbb{P}_m . The EM distance is then "the cost" of the optimal "transport".

TODO: here explain why the earth-mover is better than the the other probaility distances -> use <https://www.alexirpan.com/2017/02/22/wasserstein-gan.html>

Unfortunately the EM distance is intractable, due to the infimum part in its equation. But using the Kantorovich-Rubinstein duality, it can be reformulated as:

$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_r} [f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta} [f(x)] \quad (3)$$

Where the supremum is over 1-Lipschitz functions. It is important to note that if we replace $\|f\|_L \leq 1$ by $\|f\|_L \leq K$, i.e. consider also the K-Lipschitz functions, then we obtain $K \cdot W(\mathbb{P}_r, \mathbb{P}_\theta)$. Hence, for a family of functions $\{f_w\}_{w \in \mathcal{W}}$ (all functions being K-Lipschitz), we can consider solving the optimization problem:

$$\max_{w \in \mathcal{W}} \mathbb{E}_{x \sim \mathbb{P}_r} [f_w(x)] - \mathbb{E}_{z \sim p(z)} [f_w(g_\theta(z))] \quad (4)$$

Note: we can approximate the solution of the above mentioned problem with a Neural Network with weights \mathcal{W} . \mathcal{W} need to be compact to assure that the function f_w are K-Lipschitz. Therefore in order to have \mathcal{W} being compact, Arjovsky et al. (2017) proposes to simply clip the weights, such that they lay in a small box $[-c, c]^l$, e.g. with $c = 0.01$

2.3.2 Necessary changes to turn a GAN into a WGAN

Implementation of a WGAN requires a few changes from implementation of a regular GAN, i.e.

- Use a linear activation function in the output layer of the "discriminator" model (instead of sigmoid). The "discriminator" then becomes a critic that quantify the realness of a sample, instead of discriminating between real or fake.
- Use Wasserstein loss to train the critic and generator models that promote larger difference between scores for real and generated images.
- Constrain critic model weights to a limited range after each mini batch update (e.g. $[-0.01, 0.01]$). As seen above, this is to ensure, that the function estimated in the for approximating the Wasserstein distance are K-lipschitz, a necessary condition.
- Update the critic model more times than the generator each iteration (e.g. 5). Contrary as in a GAN, this is not a issue. Indeed, switching to the EM distance, allow for more stability when training the two networks in the WGAN. Moreover, the fact that the EM distance is continuous and differentiable means that we should train the critic until optimality.
- Use the RMSProp version of gradient descent with small learning rate and no momentum (e.g. 0.00005).

2.4 WGAN-GP

As we have seen above, WGAN improve the stability in training the critic (\approx discriminator). But it is still subject to poor sample generation, convergence failure or mode collapse. This is due to the weight clipping happening while training the critic. In order to remedy to this, Gulrajani et al. (2017) propose to replace weight clipping by the introduction of a penalization of the norm of gradient of the critic with respect to its input. The issue is that trying to orient the critic to 1-Lipschitz function by weight clipping, biases the critic for too simple function. Gulrajani et al. (2017) observes that implementing the Lipschitz constraint using weights clipping leads to either exploding or vanishing gradient, unless the threshold c used for the clipping is carefully fine-tuned.

2.4.1 The gradient penalty

In order to enforce the Lipschitz constraint, Gulrajani et al. (2017) proposes to add a penalty term to the loss function. The loss function then becomes:

$$L = L' + P \quad (5)$$

where:

- Original loss function:

$$L' = \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g}[D(\tilde{\mathbf{x}})] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r}[D(\mathbf{x})] \quad (6)$$

- Penalty:

$$P = \lambda \mathbb{E}_{\hat{\mathbf{x}} \sim \mathbb{P}_{\hat{\mathbf{x}}}} [(\|\nabla_{\hat{\mathbf{x}}} D(\hat{\mathbf{x}}) - 1\|)^2] \quad (7)$$

The goal of the penalty is to enforce the 1-Lipschitz constraint. Indeed, by definition, a function is 1-Lipschitz if and only if its gradient norm is smaller or equal to 1 everywhere. It can be easily seen that the penalty is here to enforce this constraint. In order to make such a penalty tractable, a soft version of the penalty is considered, where the constraint is only enforced on the gradient norm of a few random samples $\hat{\mathbf{x}} \sim \mathbb{P}_{\hat{\mathbf{x}}}$.

The sampling distribution $\mathbb{P}_{\hat{\mathbf{x}}}$ is defined by sampling uniformly on a line between a pair of points respectively sampled from \mathbb{P}_r and \mathbb{P}_g . This was proven experimentally to give sufficiently good results.

In Gulrajani et al. (2017), the penalty coefficient λ was set always to 10 in all experiments done. No batch normalization was used in Gulrajani et al. (2017). They claim that batch normalization shifts the discriminator problem from trying to match a single input to a single output from trying to match a batch input to a batch output. This makes the penalty invalid, since the penalization is performed with each input individually and batch normalization introduce correlation between samples. Instead of batch normalization, Gulrajani et al. (2017) recommends layer normalizations.

2.5 On GAN/WGAN/WGAN-GP Convergence

quote: "The loss of WGAN does have a convergence point: 0. We should arrive at this point when the generator is capable of generating samples so good that no Lipschitz continuous discriminator can distinguish real from generated samples."

In fact it is a major selling point of WGAN that the loss should steadily converge in a way that informs you whether the training is making progress or not. With traditional GANs, pretty much the only way of telling if the generated samples are improving is via visual inspection, and you stop the training when the visual quality of the samples is satisfying."

from Arjovsky et al. (2017) **quote:** "To our knowledge, this is the first time in GAN literature that such a property is shown, where the loss of the GAN shows properties of convergence."

-> Arjovsky et al. (2017) shows that in a WGAN, the loss function of the generator is directly correlated with the quality of the sample produced. This is not true in a regular GAN and is actually one of the main selling point of a WGAN over GAN.

From forum (<https://stats.stackexchange.com/questions/505696/when-is-my-wasserstein-gan-gp-overfitting>) : "Have you run your model on a simple dataset like MNIST to verify it's implemented correctly? In a WGAN-GP, the generator loss is typically not meaningful, and the discriminator loss is an approximation of the negative wasserstein distance between the generator distribution and real distribution. Of course, this only works if the discriminator is powerful enough, and you train it enough (usually several iterations per generator iteration). Otherwise, the discriminator loss is pretty meaningless. When it's working properly, the discriminator loss should start from 0, rapidly drop to some negative number, and then slowly work its way back to 0. (which means the generator and true distribution are getting closer)."

-> **TODO:** this need to be mentionned when showing the results of WGAN-GP

3 Our approach (TODO: title to change)

We decided that it made sense to try to generate the Cross Spectral Matrix (CSM), as done in Gerstoft et al. (2020) and extend his work to create a network to generate CSM, conditionally on Direction of Arrival (DoA). Indeed, such a network would allow us to have the online generation of labeled data required to train the network **TODO: find name of network for Source Characterization**. By providing a DoA (i.e. a label) to the network, we would generate the corresponding CSM data.

More specifically than the CSM, we thought it would make more sense to generate separately the eigenvalues and eigenvectors its eigendecomposition. A CSM $\hat{\mathbf{C}}$ can be decomposed as :

$$\hat{\mathbf{C}} = \mathbf{V} \Lambda \mathbf{V}^H \quad (8)$$

where $\mathbf{V} = [\mathbf{v}_1^T, \dots, \mathbf{v}_M^T] \in \mathbb{C}^{M \times M}$, \mathbf{v}_i being the i th eigenvector and where $\Lambda \in \mathbb{R}^{M \times M}$ is a diagonal matrix, where λ_{ii} is the i th eigenvalue, corresponding to the i th eigenvector.

Indeed, since we choose a Generative Adversarial Approach, the data will be generated using two networks: a generator and a discriminator/critic. Those two network are competing against each other: the goal of the generator is to produce data realistic enough so that discriminator can not tell it is fake. The goal of the discriminator is to tell whether a given input is real or has been generated. Both the generator and the discriminator have to be trained simultaneously until convergence. A typical issue occurring during the training, is that the discriminator becomes too good at discerning real from fake sample and hence the generator does not improve anymore.

Generating the eigenvalues and eigenvectors instead of the CSM is done in order to help the generator. This allow to normalize all the eigenvectors before feeding them to the discriminator, whether they are real or generated. The eigenvalues can also be scaled the biggest of them is equal to one. **TODO: develop on that**

3.1 Data

3.1.1 Synthetic data

The goal of this thesis was to generate data as realistic as possible, by learning the distribution of measured CSM. As mentionned in the introduction using only real measurement was not a feasible, since the number of measurement required for training is too big. For this reason, the chosen approach was to first pretrain the models using synthetic CSM and then fine tune our model using real measurement.

The synthetic measurement were sampled from the Wishart distribution. **TODO: developp here , mention added gaussian noise**

TODO: include here the eigenspectrum

3.1.2 Measurement

In order to generate realistic measurement data, measurement needed to be performed. An example of the measurement setup used is displayed in Fig.1. It is made of audio sources located in a $1\text{m} \times 1\text{m}$ plane located half a meter in front an array of 64 microphones. The microphones are organized in a circular pattern, with the central one acting as the reference microphone. This means that when the CSM are computed, they are normalize using the autopower of the reference microphone.

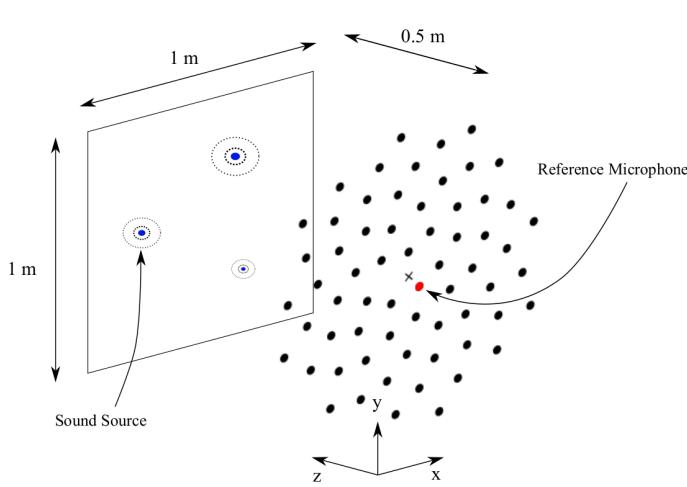


Figure 1: Example of the measurement setup used, with three audio sources

When performing the actual measurement, we only consider the case of a single audio source. The loudspeaker used in our measurement setup is displayed in Fig.2 and the microphone array in Fig.3.

TODO: what else should be specified about the used microphones/loudspeaker?

In Fig. **TODO: add figure**, we show a comparison of the normalized eigenvalues of a synthetic CSM compared with a measured CSM. In the synthetic eigenvalues, we observe first that all the eigenvalues except the first ones are close to zero. Such a string decay cannot be observed in the spectrum of the measured data. This can be explained by the fact that the synthetic CSM are closed to the perfection of a mathematical model and hence do not show presence of natural noise created by reverberation for instance. **TODO: information to check.**

3.2 Generation of Eigenvalues

3.2.1 WGAN-GP

In order to generate the eigenvalue, a WGAN-GP was built. The implementation was adapted from Nain (2020). The data to generate was the eigenvalues of the CSM. The dimension of the CSM are 64×64 , hence there are 64 eigenvalues to generate.

3.2.2 Networks Architecture

As mentionned in the fundamentals the basic structure is two competing networks, a generator and a discriminator or critic. Both architectures of the generator and critic are illustrated respectively in tables Tab.1 and Tab..2.

Unlike in the implementation of Nain (2020), both inner networks (generator and critic) do not have a convolutional but a perceptron structure. Indeed, a convolutional structure is relevant when the data to generate is an image. A convolutional layer is good at seeing patterns in image, by identifying structure in batches of neighbour pixels. We could have stacked pieces of our vector of

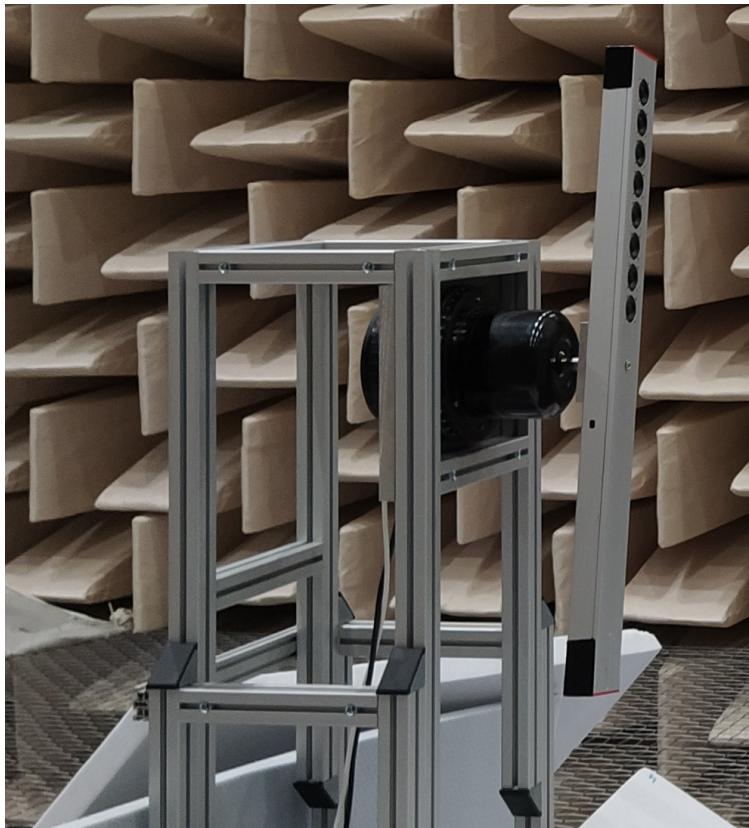


Figure 2: Device used as audio source to create the measurement

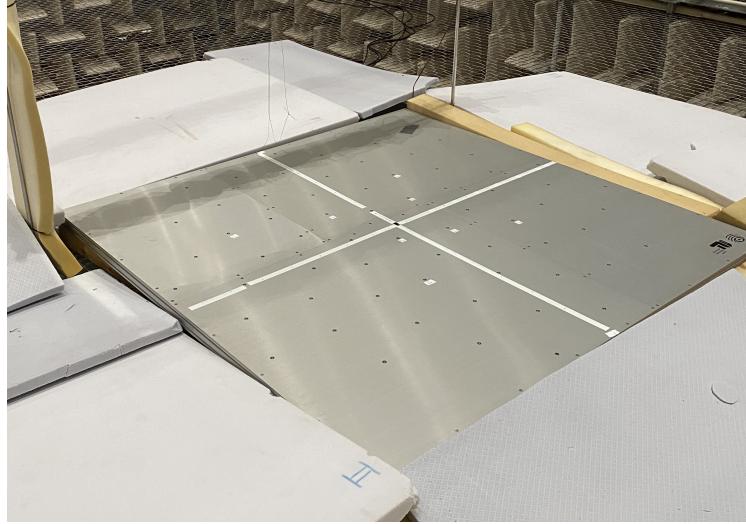


Figure 3: Picture of the array of microphone used to create the real measurement

eigenvalues of size 64 into an image of dimesion 8×8 but this would have resulted in an image where most neighbour pixel are not sharing any information about each other, hence making the convolutional operation irrelevant at detecting pattern, at least compared to images.

Another changes that was made compared to the implementation of Nain (2020), is that a ReLU activation function was added as the last layer of the generator. Indeed before this change our generator would generate eigenvalues slightly smaller than zero (of order -1e-10 **TODO: check actual value**). Moreover, we added 1e-10 to every eigenvalue produced to ensure that they are positive.

Finally, before being fed to the discriminator, real and generated eigenvalues are scaled such that the biggest eigenvalue is equal to one. Normalizing eigenvalues in this way is done to reduce the size of the sample space. This is illustrated in Fig.14.

3.2.3 Performance

A problem with regular GANs is that it is really hard to know when a good model has been found, from looking only at the loss function. The quality of a model is typically assessed by visually looking at the generated sample and deciding if they are realistic or not. But Arjovsky et al. (2017) shows that in a WGAN, the loss function of the generator is directly correlated with the quality of the sample produced. This is not true in a regular GAN and is actually one of the main selling point of a WGAN over GAN. With this knowledge, we can observe from the loss function of the generator showed in Fig.5 actually shows convergence.

Moreover when comparing real and generated eigenvalues sample in Fig.6, it can be observed that our WGAN-GP produces sufficiently realisitic results.

Layer	Output Shape	#Param.
InputLayer	(None, 128)	0
Dense	(None, 1024)	131072
Batch Normalization	(None, 1024)	4096
LeakyReLU	(None, 1024)	0
Reshape	(None, 4, 4, 64)	0
UpSampling2D	(None, 8, 8, 64)	0
Conv2D	(None, 8, 8, 128)	73728
BatchNormalization	(None, 8, 8, 128)	512
LeakyReLU	(None, 8, 8, 128)	0
UpSampling2D	(None, 16, 16, 128)	0
Conv2D	(None, 16, 16, 64)	73728
BatchNormalization	(None, 16, 16, 64)	256
LeakyReLU	(None, 16, 16, 64)	0
UpSampling2D	(None, 32, 32, 64)	0
Conv2D	(None, 8, 8, 1)	576
BatchNormalization	(None, 8, 8, 1)	4
Activation	(None, 8, 8, 1)	0

Table 1: Architecture of the generator used in the WGAN-GP to generate eigenvalues. Total params: 283,972, Trainable params: 281,538, Non-trainable params: 2,434

Layer	Output Shape	#Param.
InputLayer	(None, 8, 8, 1)	0
ZeroPadding2D	(None, 12, 12, 1)	0
Conv2D	(None, 6, 6, 64)	1664
LeakyReLU	(None, 6, 6, 64)	0
Conv2D	(None, 3, 3, 128)	204928
LeakyReLU	(None, 3, 3, 128)	0
Dropout	(None, 3, 3, 128)	0
Conv2D	(None, 2, 2, 256)	819456
LeakyReLU	(None, 2, 2, 256)	0
Dropout	(None, 2, 2, 256)	0
Conv2D	(None, 1, 1, 512)	3277312
LeakyReLU	(None, 1, 1, 512)	0
Flatten	(None, 512)	0
Dropout	(None, 512)	0
Dense	(None, 1)	513

Table 2: Architecture of the critic used in the WGAN-GP to generate eigenvalues. Total params: 4,303,873, Trainable params: 4,303,873, Non-trainable params: 0

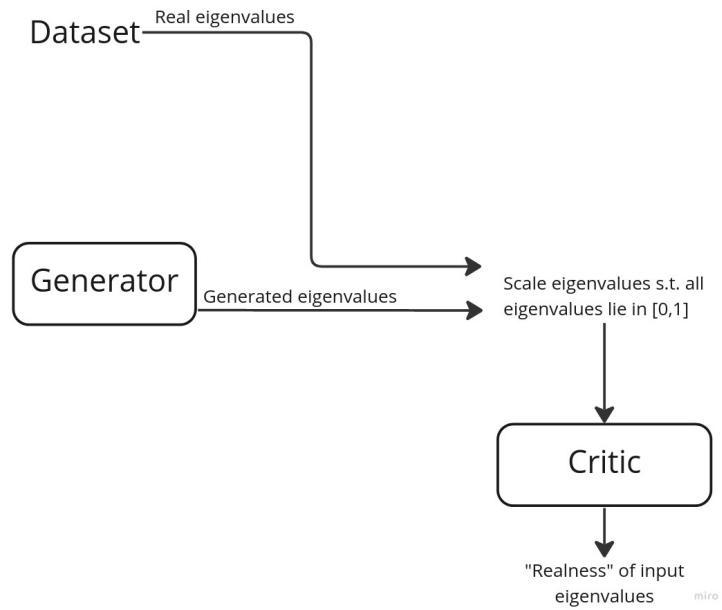


Figure 4: Full structure of the used implementation of the WGAN-GP to generate eigenvalues.

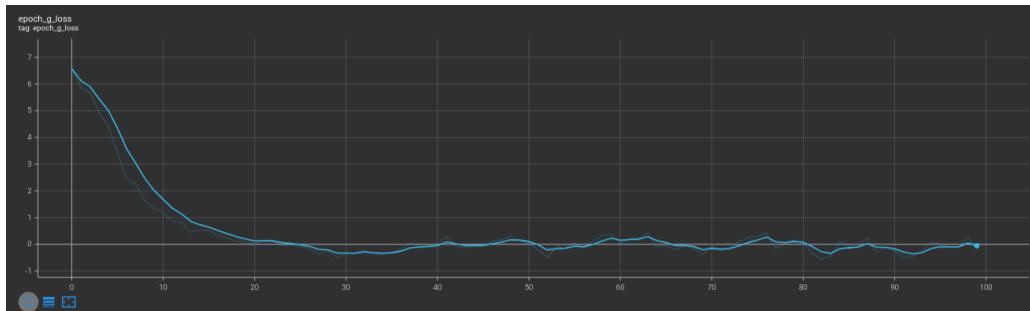


Figure 5: The loss function of the generator while training the WGAN-GP for eigenvalues generation.

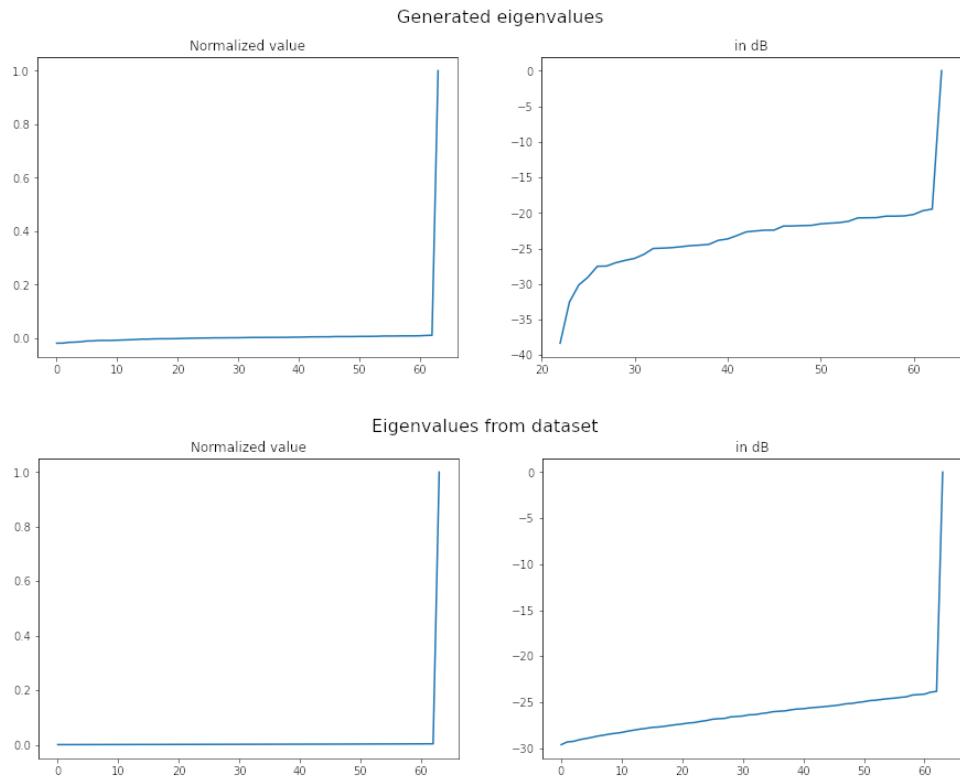


Figure 6: Top row shows eigenvalues generated by the WGAN-GP, respectively in normal values, followed by their representation in dB. The second shows the same for real eigenvalues.

3.2.4 Generating eigenvalues from their spectrum

As seen above, the WGAN-GP create convincing spectrum, until they are displayed as level. This led to idea of generating the spectrum from its levels values. We remind here how the level of the eigenvalues $[\lambda_0, \dots, \lambda_{63}] \in \Lambda$ are computed:

$$L_{\lambda_i} = 10 \log_{10}\left(\frac{\lambda_i}{\lambda_0}\right) \quad (9)$$

where λ_0 is the biggest eigenvalue. The WGAN-GP used to produce the eigenvalues needed to be adapted in order to produce the level of the spectrum. The architectures used for the critic and the generator are shown respectively in **TODO: add new architecture**.

TODO: finish writing this part when evals dB is working better

3.3 Generation of Eigenvectors

3.3.1 WGAN-GP

The eigenvectors of the cross spectral were generated using two different GANs. The idea behind this is not all vectors are sampled from the same distribution. Indeed there is only one main eigenvector (corresponding to a single DoA/source) and 63 "noise" vectors. Each eigenvector belongs to a different source and all are incoherent. If there would be multiple coherent source, all the energy and their position merges into a single eigenmode. The main eigenvector is the one corresponding to the highest eigenvalues. Formally, the set of all eigenvectors $\mathbf{V} \in \mathbb{C}^{64 \times 64}$ can be divided into two subsets $\mathbf{V}_{\text{main}} \in \mathbb{C}^{64 \times 1}$ and $\mathbf{V}_{\text{noise}} \in \mathbb{C}^{64 \times 63}$, with:

$$v \sim \mathbb{P}_{\text{main}} v \in \mathbf{V}_{\text{main}} \quad (10)$$

$$v \sim \mathbb{P}_{\text{noise}} v \in \mathbf{V}_{\text{noise}} \quad (11)$$

from mail: Yes, that makes sense. This means that the information of multiple DOAs is included in a single eigenvector. Thus, as long as we are dealing with incoherent sources, your remark fully applies.

Similarly as for generating the eigenvalues, the eigenvectors were generated using a WGAN-GP whose implementation was adapted from Nain (2020). There were 64 vectors to generate all of size 64. The vector were stacked in a matrix such that the data to generate was an image of dimension 64×64 .

3.3.2 Architecture

Again, we have two networks forming the WGAN-GP: a generator and a discriminator or critic. The architectures of both the generator and critic are illustrated respectively in the tables Tab.3 and Tab.4. Similarly as for the eigenvalues, the architectures used are directly adapted from Nain (2020) in order to produce data with desired dimensions.

In order to reduce the sample space, both the real and generated eigenvectors are normalized before being fed to the discriminator. The full WGAN-GP structure is illutrates in Fig.7.

Layer	Output Shape	#Param.
InputLayer	(None, 128)	0
Dense	(None, 16384)	2097152
Batch Normalization	(None, 16384)	65536
LeakyReLU	(None, 16384)	0
Reshape	(None, 8, 8, 256)	0
UpSampling2D	(None, 16, 16, 256)	0
Conv2D	(None, 16, 16, 128)	294912
Batch Normalization	(None, 16, 16, 128)	512
LeakyReLU	(None, 16, 16, 128)	0
UpSampling2D	(None, 32, 32, 128)	0
Conv2D	(None, 32, 32, 64)	73728
Batch Normalization	(None, 32, 32, 64)	256
LeakyReLU	(None, 32, 32, 64)	0
UpSampling2D	(None, 64, 64, 64)	0
Conv2D	(None, 64, 64, 1)	576
Batch Normalization	(None, 64, 64, 1)	4
Activation	(None, 64, 64, 1)	0

Table 3: Architecture of the generator used in the WGAN-GP to generate eigenvectors. Total params: 2,532,676, Trainable params: 2,499,522, Non-trainable params: 33,154

Layer	Output Shape	#Param.
InputLayer	(None, 64, 64, 2)	0
ZeroPadding	(None, 68, 68, 2)	0
Conv2D	(None, 34, 34, 64)	3264
LeakyReLU	(None, 34, 34, 64)	0
Conv2D	(None, 17, 17, 128)	204928
LeakyReLU	(None, 17, 17, 128)	0
Dropout	(None, 17, 17, 128)	0
Conv2D	(None, 9, 9, 256)	819456
LeakyReLU	(None, 9, 9, 256)	0
Dropout	(None, 9, 9, 256)	0
Conv2D	(None, 5, 5, 512)	3277312
LeakyReLU	(None, 5, 5, 512)	0
Flatten	(None, 12800)	0
Dropout	(None, 12800)	0
Dense	(None, 1)	12801

Table 4: Architecture of the critic used in the WGAN-GP to generate eigenvectors. Total params: 4,317,761, Trainable params: 4,317,761, Non-trainable params: 0

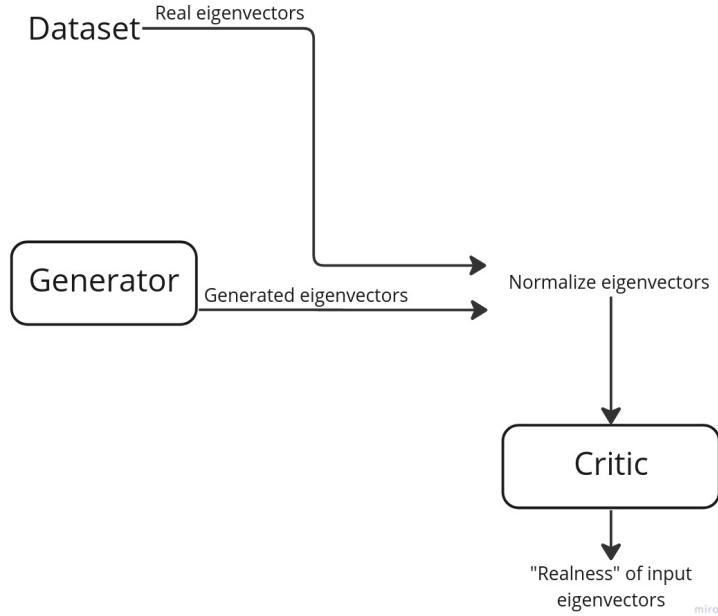


Figure 7: Full structure of the used implementation of the WGAN-GP to generate eigenvectors.

3.3.3 Performance

3.4 Data Augmentation

As seen before, we have some quite realistic looking generated eigenvalues. This led to the following idea for increasing the number of training samples (data augmentation). First, for a set of received CSM, decompose them in eigenvalues Λ and eigenvectors $\mathbf{V} = [\mathbf{v}_1^T, \dots, \mathbf{v}_M^T]$ with:

$$\hat{\mathbf{C}} = \mathbf{V}\Lambda\mathbf{V}^H \quad (12)$$

The set of all eigenvalues can then be used to train a WGAN-GP. Using this network, we can generate eigenvalues $\hat{\Lambda}$, which are realistic, up to an unknown scaling factor $c \in \mathbb{R}$.

Then using the generated eigenvalues $\hat{\Lambda}$ with real eigenvectors \mathbf{V} we can have semi-generated CSM:

$$\hat{\mathbf{C}}_{\text{augm.}} = \mathbf{V}\hat{\Lambda}\mathbf{V}^H \quad (13)$$

which are realistic up to the scaling factor c . **TODO: show that the scaling factor c does influence the beamforming.**

3.4.1 Comparison between augmented data and real data.

Fig.8 and Fig.9 show respectively the result of performing beamforming on a real CSM issued from the dataset and from a semi generated CSM. Both CSM have been recreated from the eigenvalue decomposition. It can be observed that both results are visually similar. A more careful observation

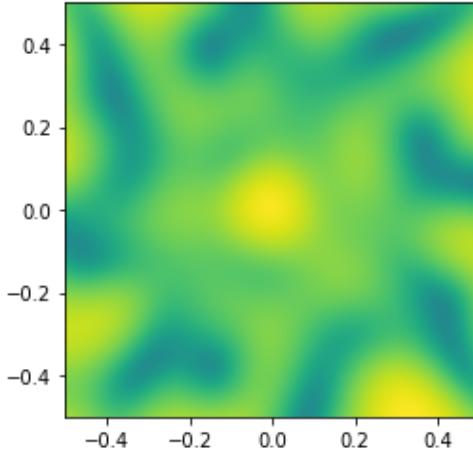


Figure 8: Results of beamforming from a real CSM

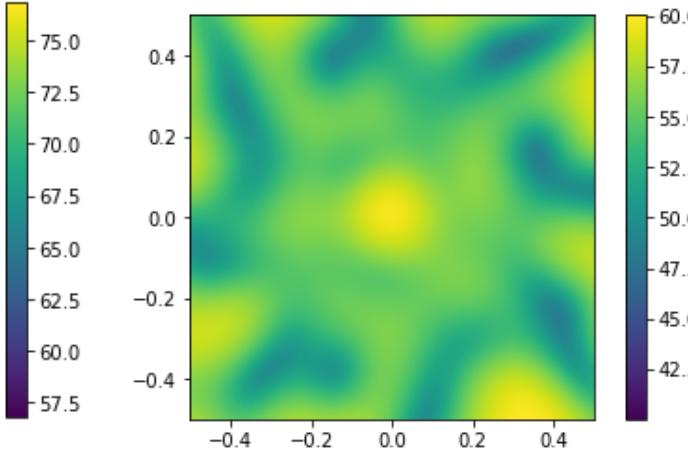


Figure 9: Results of beamforming from a CSM issued from augmented data

allows to notice that despite looking similar proportionally, the overall values of in the beamforming map resulting from real data are slightly higher. This is due to the scaling of the eigenvalues taking place in the generative process.

3.5 Generation of Cross-Correlation Matrix

4 Future work

TODO: here need to mention TransGAN, need to mention that the GAN could be made conditional

References

- Sharath Adavanne, Archontis Politis, and Tuomas Virtanen. Direction of arrival estimation for multiple sound sources using convolutional recurrent neural network. In *2018 26th European Signal Processing Conference (EUSIPCO)*, pages 1462–1466. IEEE, 2018.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pages 214–223. PMLR, 2017.
- Michael J Bianco, Sharon Gannot, and Peter Gerstoft. Semi-supervised source localization with deep generative modeling. In *2020 IEEE 30th International Workshop on Machine Learning for Signal Processing (MLSP)*, pages 1–6. IEEE, 2020.
- Paolo Castellini, Nicola Giulietti, Nicola Falcionelli, Aldo Franco Dragoni, and Paolo Chiariotti. A neural network based microphone array approach to grid-less noise source localization. *Applied Acoustics*, 177:107947, 2021.

- Soumitro Chakrabarty and Emanuël AP Habets. Broadband doa estimation using convolutional neural networks trained with noise signals. In *2017 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pages 136–140. IEEE, 2017.
- Jesse Engel, Kumar Krishna Agrawal, Shuo Chen, Ishaan Gulrajani, Chris Donahue, and Adam Roberts. Gansynth: Adversarial neural audio synthesis. *arXiv preprint arXiv:1902.08710*, 2019.
- Eric L Ferguson, Stefan B Williams, and Craig T Jin. Sound source localization in a multipath environment using convolutional neural networks. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2386–2390. IEEE, 2018.
- Peter Gerstoft, Herbert Groll, and Christoph F Mecklenbräuker. Parametric bootstrapping of array data with a generative adversarial network. In *2020 IEEE 11th Sensor Array and Multichannel Signal Processing Workshop (SAM)*, pages 1–5. IEEE, 2020.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- Ishaan Gulrajani, Faruk Ahmed, Martín Arjovsky, Vincent Dumoulin, and Aaron C.Courville. Improved training of wasserstein gans. *CoRR*, abs/1704.00028, 2017. URL <http://arxiv.org/abs/1704.00028>.
- Weipeng He, Petr Motlicek, and Jean-Marc Odobez. Deep neural networks for multiple speaker detection and localization. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 74–79. IEEE, 2018.
- Fabian Hübner, Wolfgang Mack, and Emanuël AP Habets. Efficient training data generation for phase-based doa estimation. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 456–460. IEEE, 2021.
- Adam Kujawski, Gert Herold, and Ennes Sarradj. A deep learning method for grid-free localization and quantification of sound sources. *The Journal of the Acoustical Society of America*, 146(3): EL225–EL231, 2019.
- Kundan Kumar, Rithesh Kumar, Thibault de Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose Sotelo, Alexandre de Brébisson, Yoshua Bengio, and Aaron C Courville. Melgan: Generative adversarial networks for conditional waveform synthesis. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL <https://proceedings.neurips.cc/paper/2019/file/6804c9bca0a615bdb9374d00a9fcba59-Paper.pdf>.
- Soo Young Lee, Jiho Chang, and Seungchul Lee. Deep learning-based method for multiple sound source localization with high resolution and accuracy. *Mechanical Systems and Signal Processing*, 161:107959, 2021.
- Wei Ma and Xun Liu. Phased microphone array for sound source localization with deep learning. *Aerospace Systems*, 2(2):71–81, 2019.
- Aakash Nain. Wgan-gp overriding model train step, May 2020. URL https://keras.io/examples/generative/wgan_gp/.

- Paarth Neekhara, Chris Donahue, Miller Puckette, Shlomo Dubnov, and Julian McAuley. Expediting tts synthesis with adversarial vocoding. *arXiv preprint arXiv:1904.07944*, 2019.
- Constantinos Papaiannis, Christine Evers, and Patrick A Naylor. Data augmentation of room classifiers using generative adversarial networks. *arXiv preprint arXiv:1901.03257*, 2019.
- Lauréline Perotin, Romain Serizel, Emmanuel Vincent, and Alexandre Guérin. Crnn-based joint azimuth and elevation localization with the ambisonics intensity vector. In *2018 16th International Workshop on Acoustic Signal Enhancement (IWAENC)*, pages 241–245. IEEE, 2018.
- Wagner Gonçalves Pinto, Michaël Bauerheim, and Hélène Parisot-Dupuis. Deconvoluting acoustic beamforming maps with a deep neural network. 2021.
- Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- Anton Ratnarajah, Shi-Xiong Zhang, Meng Yu, Zhenyu Tang, Dinesh Manocha, and Dong Yu. Fast-rir: Fast neural diffuse room impulse response generator. 10.48550. *arXiv preprint ARXIV.2110.04057*, 2021.
- Ryu Takeda and Kazunori Komatani. Sound source localization based on deep neural networks with directional activate function exploiting phase information. In *2016 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 405–409. IEEE, 2016.
- Elizabeth Vargas, James R Hopgood, Keith Brown, and Kartic Subr. On improved training of cnn for acoustic source localisation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:720–732, 2021.
- Juan Manuel Vera-Diaz, Daniel Pizarro, and Javier Macias-Guarasa. Acoustic source localization with deep generalized cross correlations. *Signal Processing*, 187:108169, 2021.
- Pengwei Xu, Elias JG Arcondoulis, and Yu Liu. Deep neural network models for acoustic source localization. In *Berlin Beamforming Conference*, 2021.