# Auto-adaptive Resonance Equalization using Dilated Residual Networks

Maarten Grachten[1] and Emmanuel Deruty[2] and Alexandre Tanguy[3]

July 24, 2018

[1]Contractor for Sony CSL Paris, France
[2]Sony CSL Paris, France
[3]Yascore, Paris, France

**Abstract**

In music and audio production, attenuation of spectral resonances is an important step towards a technically correct result. In this paper we present a two-component system to automate the task of resonance equalization. The first component is a dynamic equalizer that automatically detects resonances and offers to attenuate them by a user-specified factor. The second component is a deep neural network that predicts the optimal attenuation factor based on the windowed audio. The network is trained and validated on empirical data gathered from an experiment in which sound engineers choose their preferred attenuation factors for a set of tracks. We test two distinct network architectures for the predictive model and find that a dilated residual network operating directly on the audio signal is on a par with a network architecture that requires a prior audio feature extraction stage. Both architectures predict human-preferred resonance attenuation factors significantly better than a baseline approach.

## 1    Introduction and related work

Equalization is part of the audio mixing and mastering process. It is a redistribution of the energy of the signal in different frequency bands. The process has been traditionally performed by skilled sound engineers or musicians who determine the proper equalization given the character and peculiarities of the input audio. In recent years, methods have been developed for semi-automatic and automatic equalization. Such methods may be used by audio professionals to save time, or by recording enthusiasts lacking the skills required to use manual equalization tools effectively. These methods include automatic detection of frequency resonances [1], automatic equalization derived from expert practices [2], and automatic conformation to a target spectrum [3]. Equalization profiles may also be derived from semantic descriptors [4]. Appropriate equalization settings can be found through different means, for example by comparing the input source to previously equalized content [5], or by formulating equalization as an optimization problem in which inter-track masking is used as the cost function [6]. Some automated equalization functionalities are featured in

commercial products. Examples include the "learn" function in Neutron 2's equalizer[1], and the "Tame" function in SoundTheory's Gulfoss[2].

The use of machine learning to solve audio production related tasks is recent. Aside from an early approach using nearest neighbor inference to infer equalization [5], most applications of machine learning for automatic mixing date from the past few years. Automatic mixing tasks that have been addressed in this way include automatic reverbation [7], dynamic range compression [8], and demixing/remixing of tracks [9]. All three studies use neural networks, which are rapidly becoming a *de facto* standard method for machine learning. To our knowledge, there is no documented example of the use of neural networks for automatic equalization.

A particular method of equalization is the attenuation of resonating or salient frequencies, *i.e.* frequencies that are substantially louder than their neighbors [10]. The focus of this paper is the automation of such as process using a deep neural network. Salient frequencies may originate from different phenomena, such as the acoustic resonances of a physical instrument or an acoustic space. They may be considered a deficiency, in the sense that they may mask the content of other frequency regions. One particular difficulty in resonance attenuation is finding the right amount of attenuation. For example, too much attenuation may unmask noise that would otherwise remain unheard, or flatten the spectrum to the point of garbling the original audio.

Our method fully automates the resonance attenuation process. It includes, 1) a windowed, dynamic resonance attenuation process that works on 0.5s windows and can be controlled with a single parameter—the *attenuation factor*, 2) a deep neural network that predicts the attenuation factor from the input audio, making the process auto-adaptive [11].

For the training and validation of the predictive model we conduct a listening experiment determining optimal resonance attenuation factors for a set of tracks, as chosen by sound engineers. We describe and test two alternative modeling approaches, and find that a state-of-the-art convolutional network for image processing can be successfully adapted for audio processing. Experimental results show that this network architecture, which directly processes the raw audio signal, is on a par with a more traditional approach of training a neural network on a set of pre-computed audio descriptors.

The paper is organized as follows. Section 2 describes the resonance equalization process. The listening experiment is described in Section 3. The design, training, and evaluation of the predictive models is presented in Section 4, and conclusions are presented in Section 5.

## 2 Resonance Equalization

Traditionally, resonance attenuation has been a manual task in which a musician or sound engineer determines the resonating frequencies by ear or using a graphical tool, in order to reduce the energy of the signal in those frequencies by an appropriate amount [12]. In this section, we describe a procedure that identifies resonating frequencies autonomously, and reduces the energy in those frequencies by a factor that is controlled by the user. The procedure works

---

[1]`https://www.izotope.com/en/products/mix/neutron/neutron-advanced.html`
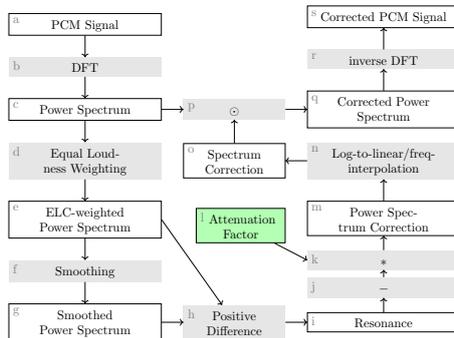[2]`https://www.soundtheory.com/home`

Figure 1: Resonance equalization block diagram; White and gray blocks represent data and processes respectively; The green block depicts the single user-controlled parameter; The symbols $\odot$, $*$ and $-$ represent elementwise vector/vector multiplication, elementwise scalar/vector multiplication, and unary negation respectively.

on overlapping audio windows that must be large enough to allow for spectral analysis at a high frequency resolution.

Figure 1 displays a block diagram of the resonance attenuation process, where each element is denoted by a letter. In the following, we will use these letters to refer to the corresponding elements in the diagram. First the audio signal is used to compute a power spectrum weighted by Equal-Loudness Contours (ELC) [13] at a fixed monitoring level of 80 phon (Figure 1, element d) in order to reflect the perceptual salience of the signal energy at different frequency bands. The value of 80 phon is chosen in relation to the procedure detailed in Section 3.3. The ELC-weighted power spectrum (e) consists of 400 log-scaled frequency bands.

Resonances (i) are determined by smoothing ELC-weighted power spectrum (e) to obtain (g) and computing the elementwise differences (e) minus (g), setting negative elements to zero (h). The negative of the resonances is then scaled by the user defined attenuation factor (l), transformed back to a linear scale and converted back to the shape of the original spectrum using interpolation (h). The result (o) is a vector of scaling factors (one for each DFT bin) that range between 0 and 1 for resonating frequencies and are 1 otherwise. Multiplying the original power spectrum (c) with the scaling factors gives the corrected power spectrum (q) from which the corrected audio signal (s) is recovered through the inverse DFT (r).

# 3 Listening experiment

A listening experiment was carried out to obtain ground truth in terms of optimal resonance attenuation values for a set of audio tracks. In the experimental design of the listening test it proved unpractical to ask subjects to set a varying resonance attenuation factor. Therefore, we chose relatively homogeneous sound fragments, and let the subjects choose a single attenuation factor for the whole fragment.

## 3.1 Participants

A group of 15 subjects was recruited for the experiment, in the area of Paris (France). All of them are recognized professionals in the industry. Nine subjects specialize in studio recording (Classic, Jazz, Pop, Rock, Movie Music, audio post-production), three are experts in live music, and three are composer/music producer. The subjects were between 24 and 42 years old, with an average age of 32 years. They were recruited and paid as if they were working on a commercial project.

## 3.2 Data

A set of 150 audio tracks was used for the listening experiment. The tracks are excerpts from longer pieces, with a mean duration of 46 seconds and a standard deviation of 16 seconds. All tracks were processed using Nugen AMB R128[3] so that they were aligned to the same median loudness. The set comprised pop and rock music, as well as film scores. Of this set, 131 tracks were unique recordings, while the remaining 19 tracks were variants of some of the unique 131 recordings, with differences in mixing. None of the tracks were previously mastered.

## 3.3 Procedure

The listening experiment took place in a recording studio, where participants listened to the audio tracks individually, using studio monitors, at a measured listening loudness of 80 dBC. This value was chosen as being representative of a normal listening loudness during audio production. The participants were presented with a web interface in which they could listen to each track with different degrees of resonance attenuation, ranging from 0 (no attenuation) to 1 in 17 steps. They could select their preferred degree of resonance attenuation, or alternatively decline to select any version, indicating that none of the versions sounded acceptable. The tracks were separated from each other by 10 seconds of pink noise surrounded by a short silence to give the participants a fixed reference. Sessions of 50 tracks were alternated with breaks.

## 3.4 Results and discussion

Basic statistics of the results per subject are given in Table 1. Subject 13 stands out because of the number of missing ratings (21 versus a median of 1 over all subjects). Subjects 1 and 15 have abnormally high rates of 0.0 ratings (72 and 58 respectively, versus a median of 16 over all subjects). Finally, Subject 7 stands out in terms of median rating (0.469 versus a median of 0.188 over all subjects). Figure 2 shows the distribution of ratings per subject.

To see how strongly the ratings are linearly related among subjects, we compute the Pearson correlation for each pair of subjects (Figure 3). Apart from Subjects 13 and 15 (and to a lesser degree Subject 1) who appear to have different rating patterns from the majority of the subjects, the figure shows weak to moderate positive correlations between all subjects. This suggests that

---

[3] https://nugenaudio.com/amb

Table 1: Rating statistics per subject. Outlying values are highlighted in bold (see text).

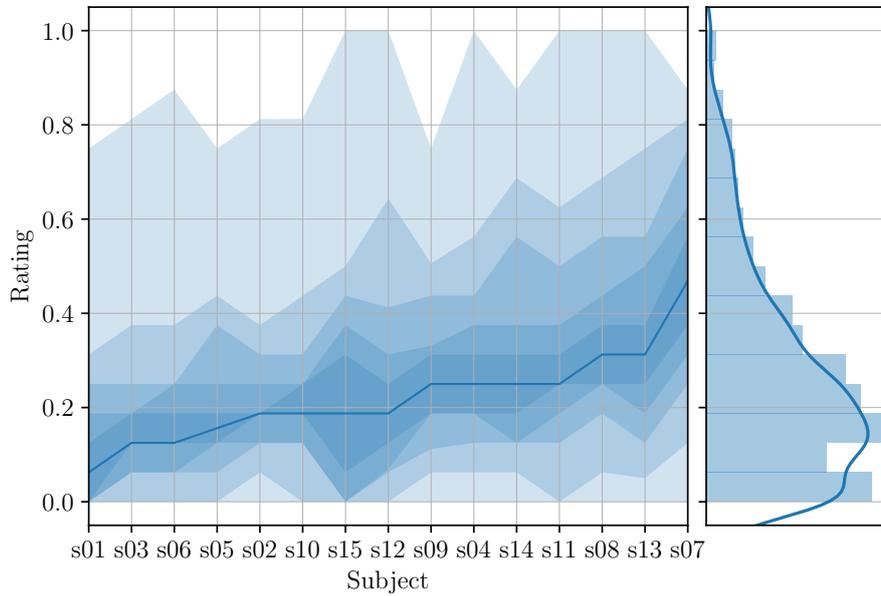| Subject | # No rating | # 0.0 | Min | Median | Max |
|---|---|---|---|---|---|
| s01 | 1 | **72** | 0.0 | 0.062 | 0.750 |
| s02 | 0 | 2 | 0.0 | 0.188 | 0.812 |
| s03 | 2 | 25 | 0.0 | 0.125 | 0.812 |
| s04 | 2 | 7 | 0.0 | 0.250 | 1.000 |
| s05 | 4 | 25 | 0.0 | 0.156 | 0.750 |
| s06 | 0 | 22 | 0.0 | 0.125 | 0.875 |
| s07 | 0 | 2 | 0.0 | **0.469** | 0.875 |
| s08 | 8 | 9 | 0.0 | 0.312 | 1.000 |
| s09 | 0 | 9 | 0.0 | 0.250 | 0.750 |
| s10 | 1 | 17 | 0.0 | 0.188 | 0.812 |
| s11 | 1 | 16 | 0.0 | 0.250 | 1.000 |
| s12 | 2 | 25 | 0.0 | 0.188 | 1.000 |
| s13 | **21** | 13 | 0.0 | 0.312 | 1.000 |
| s14 | 0 | 3 | 0.0 | 0.250 | 0.875 |
| s15 | 0 | **58** | 0.0 | 0.188 | 1.000 |



Figure 2: Rating percentiles (in steps of 10%) per subject. Darker areas correspond to more central percentile ranges, lighter areas to more peripheral ranges. The bold line in the left plot shows the median ratings.
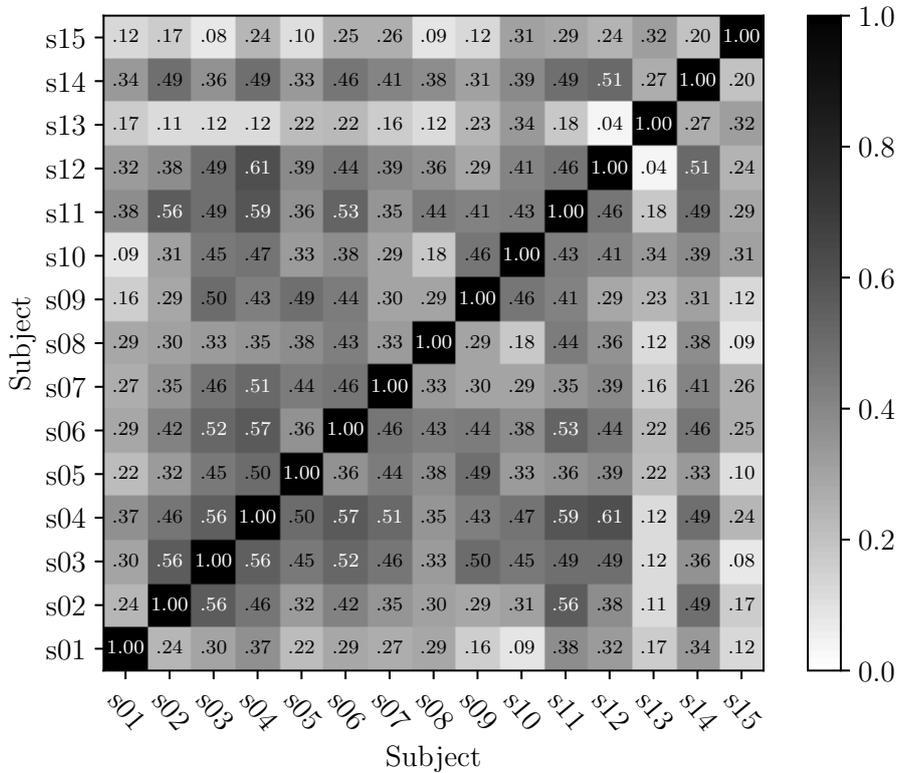
Figure 3: Correlation coefficients of ratings among subjects.

in spite of different preferred rating ranges, the subjects made their judgments according to common criteria.

# 4    Prediction of optimal attenuation factors

In this section we describe the design and experimental validation of two modeling approaches to predict optimal attenuation factors from audio. The models are ultimately intended to be used in a real-time plugin for audio workstations. Although the details the real-time aspects of the implementation are beyond the scope of this paper, it does guide some important design decisions for the modeling. Most importantly it implies a causal design in which the track cannot be analyzed as a whole in order to estimate the optimal attenuation factor. On the other hand, the audio latency upper-bound for real-time operation (maximum observed audio latency in real-time commercial plug-ins is 4096 samples) is too low for accurate prediction of the attenuation factor. This implies that the attenuation factor that will be applied at time $t$ will be estimated from a windowed part of the signal (immediately) before $t$. Whether the predicted attenuation factor is still approximately valid for the signal at time $t$ depends on the length of the window in relation to how quickly the resonance characteristics of the signal can change. Our point of view is to target phenomena whose time-scale

is around three seconds or longer—the same time scale as the *short-term loud-ness* in EBU R128 [14]. From this perspective we consider a window size of 0.5s a good trade-off, offering sufficient data for an informed prediction and at the same time being short enough to adapt to changes in resonance characteristics at the 3s time-scale.

We describe and test two alternative neural network approaches to the problem of predicting optimal attenuation factors. The first is a more traditional approach in which a feature set is computed from an audio window, from which the attenuation factor is predicted. The second approach skips the intermediate feature representation. Instead, it takes the stereo PCM signal directly as input to a neural network that predicts the attenuation factor.

## 4.1   Predicting based on audio features (FFN)

Performing regression or classification tasks on audio using feature descriptions of the audio has been the predominant approach for the past decades, and is based on the intuition that the prediction is determined by characteristics of the signal that can be defined explicitly and computed from the audio. These descriptors often capture spectral characteristics of the signal, but may also approximate perceptual characteristics, such as loudness. Many audio descriptors that have been proposed in the literature over time are implemented in a software library called Essentia [15]. The descriptors used in this study are listed in Table 2. All descriptors are available in the Essentia library, except *harmonics-to-noise ratio* [16] and *stereo width* (two descriptors, computed as the correlation between channels and absolute difference in RMS between channels, respectively), for which we used our own implementation.

The features are computed on shorter timescales (typically 1024 samples) than the 0.5s audio window for which our prediction will be made. Thus the feature computation stage returns a vector of values for each feature. We summarize each of these vectors by 7 statistics: the *mean, median, standard deviation, skew, kurtosis*, the $10^{th}$ *percentile*, and the $90^{th}$ *percentile*. This yields a total of 679 values per data instance, based on which a prediction must be made. A common preprocessing step for high-dimensional input spaces like this is dimension reduction by means of *principal component analysis* (PCA). Tests showed that PCA did not lead to any substantial improvements however, and was thus not included in the final experiment.

The network consists in a stack of linear layers (also called *dense*, or *fully connected*), each of which is followed by a *batch normalization* (BN) layer and a layer of *rectified linear* units (ReLU). The BN layer transforms the distribution of the output activations of the preceding linear layer to zero mean and unit variance by keeping track of mean and variance during the training of the model. The ReLU layer performs a non-linear transformation by setting negative output activations of the preceding layer to zero. The number of linear layers and their sizes are not fixed in advance but determined using a hyper-parameter optimization scheme (Section 4.3). A final linear layer is added after the after the last ReLU layer. This layer has a single output—the predicted resonance attenuation factor.

Table 2: List of audio descriptors used in the FFN.

| | |
|---|---|
| MFCC (13 values) | zero-crossing rate |
| GFCC (13 values) | spectral flatness dB |
| inharmonicity | high frequency coefficient |
| pitch | barkbands (30 values) |
| pitch salience | pitch instantaneous confidence |
| spectral complexity | silence rate (at 20/30/60dB) |
| spectral crest | odd-to-even harmonic energy ratio |
| spectral decrease | spectral energy band (4 values) |
| spectral energy | tristimulus (3 values) |
| spectral flux | spectral contrast (6 values) |
| spectral rms | spectral valley (6 values) |
| spectral rolloff | stereo width (2 values) |
| spectral strong peak | harmonics-to-noise ratio |

## 4.2 Predicting directly from audio: Dilated Residual Networks (DRN)

In this section we describe a convolutional neural network that takes slices of a stereo PCM signal of the audio as input and provides an estimate of the optimal attenuation factor. Note that even for a window of moderate size and sample rate this quickly leads to tens of thousands of samples to be taken as model inputs. As opposed to a feature vector however, the inputs are ordered along a meaningful dimension (time), in which patterns can be identified. A common way do deal with data exhibiting such topological properties (sound, images, video) in neural networks is to use *convolution*. This approach, which was pioneered in [17], exploits the fact that such data display local patterns that may occur at different locations in the data. The strength of convolutional networks is that they learn to recognize patterns independently of their absolute location, and at the same time the convolution operation is much more space efficient than the "fully-connected" matrix dot product that takes place in regular feed-forward neural networks, allowing for larger models. By stacking convolutional layers on top of each other it is possible to detect patterns of increasing size, and by interleaving the convolution operation with so-called *pooling* or *sub-sampling* operations, the patterns become somewhat invariant to local deformations.

However promising, the potential of traditional convolutional networks has been limited by a number of factors. Two of these limitations have been addressed by recent extensions of the traditional convolutional network approach, namely *dilated convolution*, and *residual networks*. We integrate both extensions in our convolutional network for predicting resonance attenuation factors, and discuss each of them briefly before we describe the global architecture of the model.

### 4.2.1 Dilated convolution

An approach often used with traditional convolution in order to create high-level feature representations of data is pooling using *max* or *average* aggregation functions. For instance, max-pooling sub-samples the input by selecting maximal

elements in a sliding window, typically using a sliding step (*stride*) equal to the size of the pooling window. Stacking convolution/pooling operations leads to features with increasing *receptive fields*, meaning that the features can describe patterns of increasing size. However, it comes at the cost of resolution loss: The relative position of features becomes less precise as their size increases.

On the contrary, dilated convolution achieves high-level features without loss of resolution. Rather than increasing stride, it increases the receptive field of the features by "dilating" the convolution kernels. A normal convolution of the kernel $k$ with the signal $s$ involves multiplying kernel elements with contiguous signal samples ($\tau$ is a discrete variable that increases in steps of 1):

$$(k * s)(t) = \sum_{\tau=-\infty}^{\infty} k(\tau) \, s(t - \tau) \tag{1}$$

In convolution with dilation factor $d \in \mathbb{Z}^+$ on the other hand the kernel elements are multiplied with signal samples that are equally spaced at $d$ samples:

$$(k *_d s)(t) = \sum_{\tau=-\infty}^{\infty} k(\tau) \, s(t - d\tau) \tag{2}$$

By stacking convolutional layers with increasing dilation factors the higher level filter kernels aggregate information over input ranges of exponentially increasing size, even if the size the kernels (in terms of parameters) does not increase, and the resolution remains intact. This approach has proven successful in image processing tasks such as semantic segmentation [18].

### 4.2.2  Residual blocks

Another issue with convolutional networks is that as they grow deeper in order to capture higher level patterns, it becomes harder to optimize the lower level convolutional layers. This is directly related to the fact that for low level feature activations to influence the output of the model, they must pass through multiple layers of convolutions. Sometimes however, it is desirable for low level features to be able to directly influence the output of the model, not just to figure as a building block for higher level features.

This observation has led to the proposal of the *residual block* as a substructure used in deep networks [19, 20], an adaptation of which is depicted in Figure 4. In this structure the information flows from input to output through two pathways in parallel. The left pathway involves a typical convolution layer with configurable parameters: the kernel size $k$, number of kernels $n$, and the dilation factor $d$. The convolution in the right pathway uses kernels of size one (the dilation factor is irrelevant in that case), and thus does not compute any features from the input. Instead, it outputs $n$ linear combinations of the input in order to make the input shape compatible for elementwise addition to the $n$ *feature maps* of the left pathway. Both pathways further include a batch normalization operation. After the elementwise sum of both pathways a rectified linear unit allows for a non-linear response.

The term "residual" refers to the fact that the left pathway only needs to account for the part of the output that cannot be accounted for by linear combinations of the input—the right pathway. In this way, increasing the number
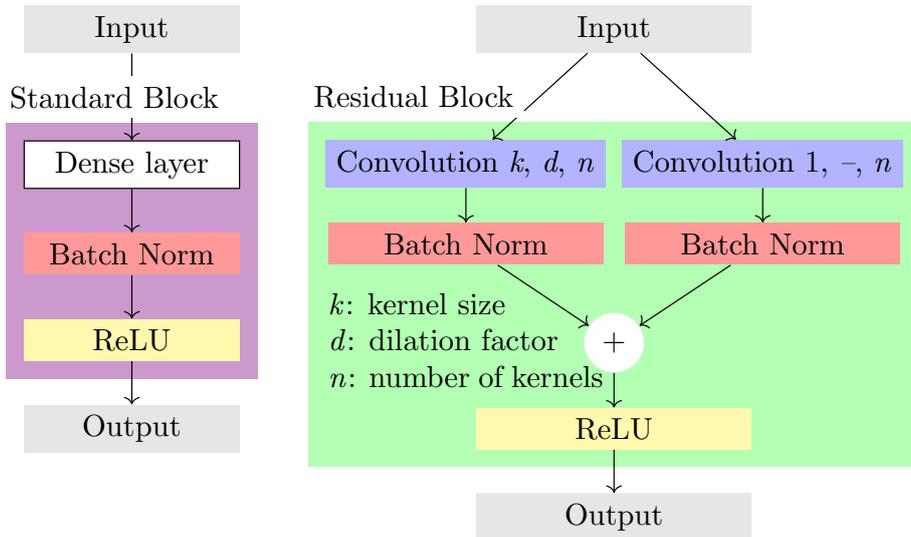
Figure 4: Building blocks for the FFN and DRN models. Left: Standard block composed of a dense linear layer followed by batch normalization and a rectified-linear layer (See Section 4.1); Right: Residual block (See Section 4.2.2).

of stacked convolution operations does not hamper the ability of the network to account for its output in terms of lower level features.

### 4.2.3 DRN architecture

Figure 4 shows the complete model consisting of multiple residual blocks. Note that the residual blocks maintain the original temporal dimension of the data, which amounts to a size of 11025 for 0.5s of audio sampled at 22050Hz. The temporal pooling operation reduces this number by down-sampling the output using window-wise averaging, and is followed by two dense layers with intermediate batch-normalization and non-linearity in order to produce an estimated resonance attenuation factor.

## 4.3 Experiments

In this section we describe the training and evaluation procedure of both model architectures described above. We use the human ratings of the 150 tracks gathered in the listening experiment to train and evaluate both architectures.

### 4.3.1 Procedure

**Evaluation Criterion** Predicting optimal resonance attenuation factors for given tracks is a regression problem and as such an obvious choice is to use the mean squared error of the predictions with respect to the optimal value (the *target*) as an objective to be minimized. However, given the variance in the ratings across subjects in the ground truth, it may be hard to determine a unique optimal value per track. Using the mean or median of the ratings per track as a target has the drawback that the mean squared error objective
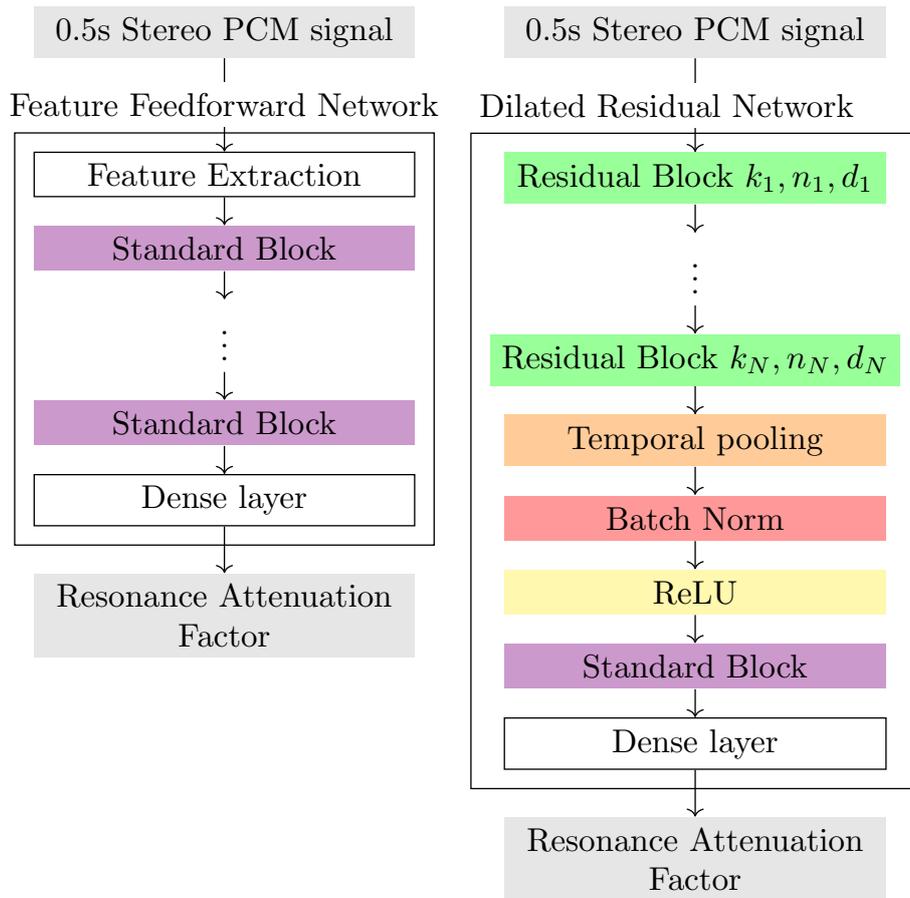
Figure 5: FFN and DRN architectures.

does not differentiate between tracks with different degrees of rater consensus. Ideally, we wish to impose a lower penalty on errors from the mean rating when the rater consensus is low. We do so by generalizing the mean squared error objective as follows. Rather than defining the objective function to be minimal only for a single value, we define it to be minimal whenever the prediction lies within a specified interval that varies from one data instance to another. For a given data instance consisting of an audio track $A \in \mathbf{A}$ and a set of ratings $\mathbf{F}$ we define the zero penalty interval as $[\ P_l(\mathbf{F}),\ P_h(\mathbf{F})\ ]$, where $P_l(\mathbf{F})$ and $P_h(\mathbf{F})$ denote the $l$-th and $h$-th percentiles of the ratings $\mathbf{F}$, respectively, with $l \leq h$. We refer to this objective as the *mean squared bounds errror* with bounds $l$, $h$, or $MSBE^{(l,\,h)}$. We use $l = 35$ and $h = 65$ throughout the experiments.

Formally, given a dataset $D$ consisting of pairs $(A, \mathbf{F})$, the $MSBE^{(l,\,h)}$ of a model $f : \mathbf{A} \to \mathbb{R}$ is defined as:

$$MSBE^{(l,\,h)}(f, D) = \frac{1}{|D|} \sum_{(A,\ \mathbf{F})\ \in\ D} L^{(l,\,h)}_{A,\mathbf{F}}(f), \tag{3}$$

where

$$
\begin{aligned}
L^{(l,\,h)}_{A,\mathbf{F}}(f) \ = \ & \Big( [\ f(A) - P_h(\mathbf{F})\ ]^+ + \\
& [\ f(A) - P_l(\mathbf{F})\ ]^- \Big)^2 .
\end{aligned}
\tag{4}
$$

The brackets $[\ \cdot\ ]^+$ and $[\ \cdot\ ]^-$ denote the positive and negative parts respectively.

**Hyper-parameter optimization**  We use Bayesian optimization to find the optimal hyper-parameters for each of the models, most importantly the depth of the networks and the hidden layer sizes. This is a heuristic to speed up the search for appropriate hyper-parameter values compared to an exhaustive grid search. The particular form of optimization we use is based on a gaussian process approximation of the loss as a function of the hyper-parameters. This approximation gives rise to the *upper confidence bound* [21], which estimates the expected loss for hyper-parameter settings that have not yet been tested, and is used as a guide to search the space of hyper-parameters [22].

Apart from the depth of the models and the hidden layer sizes, the optimization involved hyper-parameters to control the training procedure: the *learning rate*, and the thresholds for *early stopping*, and *learning rate reduction*.

**Cross-validation**  To perform the hyper-parameter optimization we use two partitions of the dataset into a test set (10 tracks), a validation set (10 tracks), and a train set (130 tracks). For each of the test tracks we compute the $MSBE^{(35,\,65)}$ loss on 100 randomly selected 0.5s frames. The criterion used to optimize the hyper-parameters is the average frame-wise loss across both test sets.

With the best hyper-parameters found for the FFN and DRN architectures, respectively, we perform a further five fold cross-validation. To this end, we use the same dataset, but exclude the 20 test tracks used for hyper-parameter optimization. We repeat the five fold cross validation five times using different

Table 3: Optimal configuration for the FFN and the DRN architectures as found by hyper-parameter optimization.

|  | FFN | DRN |
|---|---|---|
| Depth | 3 Std. Blocks | 10 Res. Blocks |
| Block Size (Low/Mid/High) | 500 / 250 / 250 | 100 / 100 / 300 |
| Temporal Pooling | – | 300 |
| Final Std. Block Size | – | 10 |

Table 4: Means, standard deviations, and the 95% confidence interval (CI) for the mean $MSBE^{(35,\,65)}$ per model.

| Model | Mean | Std. dev. | 95% CI Low | 95% CI High |
|---|---|---|---|---|
| Baseline | 0.237 | 0.103 | 0.194 | 0.280 |
| FFN | 0.159 | 0.082 | 0.124 | 0.194 |
| DRN | 0.154 | 0.080 | 0.121 | 0.188 |

random seeds to reduce the effect of partitioning of the data into folds and model parameter initializations on the result.

**Baseline**    We define a baseline approach as a reference for evaluating the FFN and DRN architectures. This approach consists in computing the mean resonance attenuation factor observed over all tracks in the training set, and using this value as a prediction for the test set, irrespective of the input audio.

### 4.3.2    Results and discussion

The optimal configuration for the FFN and the DRN architectures, as found by hyper-parameter optimization, are shown in Table 3. Figure 6 shows the results of these architectures on the repeated five fold cross validations. A one-way repeated measures ANOVA reveals a significant effect of model on $MSBE^{(35,\,65)}$ ($F_{2,72} = 6.55$, $p = 0.002$). A post-hoc Tukey HSD test at $\alpha = 0.05$ indicates that DRN and FFN differ significantly from the baseline. The effect size of DRN over baseline corresponds to Cohen's $d = 0.88$, whereas the FFN over baseline effect size is $d = 0.82$. The difference between DRN and FFN is not significant.

Table 3 shows that the FFN architecture works best when it is comparatively shallow (three standard blocks, the minimal depth tested), whereas the DRN architecture performs better when it is deep (10 residual blocks, the maximal depth tested). This trend is consistent in a review of the 10 best FFN and DRN architectures as found in the hyper-parameter optimization, omitted here for the sake of brevity. The layer sizes however do not show a similarly consistent trend, and vary considerably throughout the 10 best FFN and DRN architectures.

The fact that both modeling approaches show similar accuracies on the attenuation factor prediction is in line with a general trend that arises from the deep learning literature, showing that current techniques in deep learning are powerful enough to work with raw digitizations of information—such as sampled
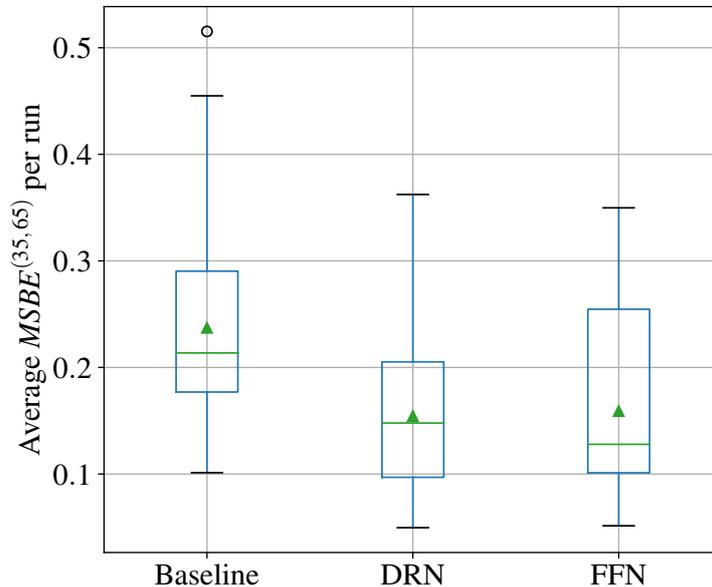
Figure 6: Average $MSBE^{(35,\,65)}$ per run for baseline, FFN, and DRN models; The horizontal lines in the boxes indicate the median, triangles the mean values per model.

waveforms—and still address tasks that require considerable abstraction from that low level representation, thus reducing the need for hand designed features.

At the same time, it must be noted that the roughly equivalent performance of the FFN and DRN measured here is at odds with a multitude of cases where end-to-end deep networks clearly outperform prior state-of-the-art methods that rely on a hand-designed feature extraction stage, especially in image processing [23]. For audio tasks such as automatic tagging however, end-to-end networks do not seem to have a strong advantage over spectrogram-based approaches [24], and require large training data sets in order to outperform spectrogram-based approaches in audio tagging tasks [25]. Similarly, in the study presented here the small size of the data set—especially in combination with inter-subject variance and the non-uniform distribution of the ratings—is a plausible explanation of why the DRN does not outperform the FFN.

# 5   Conclusion

In this paper we addressed the problem of automatic resonance equalization. We have proposed a method to attenuate automatically identified resonances by a user-controlled factor. Using this method we carried out a listening experiment in which sound engineers identify optimal attenuation factors for a set of audio tracks. The results, which show general consensus in ratings among subjects, were used to train and evaluate two types of predictive models to estimate optimal resonance attenuation factors based on the content of the audio tracks.

The results show that an intermediate stage of feature extraction is not strictly necessary for this task: a dilated residual network performs equally well when applied directly to the audio signal.

The proposed system is a fully auto-adaptive resonance equalization system in which the attenuation factor is chosen automatically by a deep neural network. To our knowledge, this system is the first documented self-adaptive equalizer based on neural networks. Future work includes a real-time implementation of the presented model as a real-time plugin that can be used in audio work stations.

# References

[1] J. Bitzer, J. LeBoeuf, "Automatic detection of salient frequencies," presented at the *Audio Engineering Society Convention 126* (2009).

[2] B. De Man, J. D. Reiss, "A knowledge-engineered autonomous mixing system," presented at the *Audio Engineering Society Convention 135* (2013).

[3] Z. Ma, J. D. Reiss, D. A. Black, "Implementation of an intelligent equalization tool using Yule-Walker for music mixing and mastering," presented at the *Audio Engineering Society Convention 134* (2013).

[4] M. B. Cartwright, B. Pardo, "Social-EQ: Crowdsourcing an Equalization Descriptor Map." presented at the *ISMIR*, pp. 395–400 (2013).

[5] D. Reed, "A Perceptual Assistant to Do Sound Equalization," presented at the *5th International Conference on Intelligent User Interfaces*, pp. 212–218 (2000).

[6] S. Hafezi, J. D. Reiss, "Autonomous multitrack equalization based on masking reduction," *Journal of the Audio Eng. Soc.*, vol. 63, no. 5, pp. 312–323 (2015).

[7] E. T. Chourdakis, J. D. Reiss, "A Machine-Learning Approach to Application of Intelligent Artificial Reverberation," *J. Audio Eng. Soc.*, vol. 65, no. 1/2, pp. 56–65 (2017).

[8] S. I. Mimilakis, K. Drossos, T. Virtanen, G. Schuller, "Deep Neural Networks for Dynamic Range Compression in Mastering Applications," presented at the *Audio Engineering Society Convention 140* (2016 May).

[9] S. I. Mimilakis, E. Cano, J. Abeßer, G. Schuller, "New Sonorities for Jazz Recordings: Separation and Mixing using Deep Neural Networks," presented at the *2nd AES Workshop on Intelligent Music Production* (2016).

[10] J. Bitzer, J. LeBoeuf, U. Simmer, "Evaluating perception of salient frequencies: Do mixing engineers hear the same thing?" presented at the *Audio Engineering Society Convention 124* (2008).

[11] J. Reiss, Ø. Brandtsegg, "Applications of cross-adaptive audio effects: automatic mixing, live performance and everything in between," *Frontiers in Digital Humanities*, vol. 5, p. 17 (2018).

[12] G. McCandless, D. McIntyre, *The Craft of Contemporary Commercial Music* (Routledge) (2017).

[13] International Standardization Organization, "ISO 226:2003: Acoustics– Normal equal-loudness-level contours," Geneva, Switzerland.

[14] EBU-R-128, "BU Tech 3341-2011, Practical Guidelines for Production and Implementation in Accordance with EBU-R-128," (2011).

[15] D. Bogdanov, N. Wack, E. Gomez, S. Gulati, P. Herrera, O. Mayor, G. Roma, J. Salamon, J. R. Zapata, X. Serra, "Essentia: an audio analysis library for music information retrieval," presented at the *14th International Society for Music Information Retrieval Conference* (2013 November 4-8).

[16] P. Boersma, "Accurate short-term analysis of the fundamental frequency and the harmonics-to-noise ratio of a sampled sound," presented at the *IFA 17*, pp. 97–110 (1993).

[17] Y. Lecun, L. Bottou, Y. Bengio, P. Haffner, "Gradient-based learning applied to document recognition," presented at the *IEEE*, vol. 6, pp. 2278–2324 (1998).

[18] F. Yu, V. Koltun, "Multi-Scale Context Aggregation by Dilated Convolutions," *CoRR*, vol. abs/1511.07122 (2015).

[19] K. He, X. Zhang, S. Ren, J. Sun, "Deep Residual Learning for Image Recognition," presented at the *IEEE Conference on Computer Vision and Pattern Recognition CVPR*, pp. 770–778 (2016).

[20] F. Yu, V. Koltun, T. A. Funkhouser, "Dilated Residual Networks," *CoRR*, vol. abs/1705.09914 (2017).

[21] T. Lai, H. Robbins, "Asymptotically Efficient Adaptive Allocation Rules," *Adv. Appl. Math.*, vol. 6, no. 1, pp. 4–22 (1985 Mar.).

[22] J. Snoek, H. Larochelle, R. P. Adams, "Practical bayesian optimization of machine learning algorithms," presented at the *Advances in neural information processing systems*, pp. 2951–2959 (2012).

[23] L. Zheng, Y. Yang, Q. Tian, "SIFT meets CNN: A decade survey of instance retrieval," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 5, pp. 1224–1244 (2018).

[24] S. Dieleman, B. Schrauwen, "End-to-end learning for music audio," presented at the *IEEE Int. Conf. on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6964–6968 (2014).

[25] J. Pons, O. Nieto, M. Prockup, E. M. Schmidt, A. F. Ehmann, X. Serra, "End-to-end learning for music audio tagging at scale," presented at the *Workshop on Machine Learning for Audio Signal Processing (NIPS)* (2017).