

# **Composing Dynamic Soundscapes Using Neural Networks and Sentimental Input**

**CS350 Dissertation Project Final Report**

University of Warwick, BSc Data Science

## **Author**

Harrison Wilde (u1600779)

## **Supervisor**

Professor Graham Cormode

# Table of Contents

<b>Table of Contents</b> . . . . .	<b>i</b>
<b>List of Figures</b> . . . . .	<b>iii</b>
<b>List of Tables</b> . . . . .	<b>iv</b>
<b>Acknowledgements</b> . . . . .	<b>v</b>
<b>Abstract</b> . . . . .	<b>vi</b>
<b>1 Introduction</b> . . . . .	<b>1</b>
1.1 Motivation . . . . .	2
1.1.1 Brief History of Computational and Algorithmic Approaches to Com- position and Musical Applications . . . . .	2
1.2 Research Objectives . . . . .	3
1.3 Eventual Objectives and Project Requirements . . . . .	4
<b>2 Background</b> . . . . .	<b>4</b>
2.1 Related Work . . . . .	4
2.1.1 Competitive Existing Solutions . . . . .	4
2.1.2 Transcription and Encoding Techniques . . . . .	4
2.1.3 Data . . . . .	6
2.2 Composition . . . . .	6
2.2.1 Markovian Approach . . . . .	6
2.2.2 Neural Network Approach . . . . .	7
2.3 User Input and Interface . . . . .	9
2.4 Articles of Interest in the Field . . . . .	9
<b>3 Methodology</b> . . . . .	<b>10</b>
3.1 Research Methodology . . . . .	10
3.2 Project Management . . . . .	10
3.3 Development Methodology . . . . .	11
3.4 Ethical Considerations . . . . .	11
<b>4 Technical Component</b> . . . . .	<b>12</b>
4.1 MIDI Transcription . . . . .	12

4.1.1	Implementation . . . . .	13
4.2	Sentiment Analysis . . . . .	13
4.2.1	Implementation . . . . .	15
4.3	Recurrent Neural Networks and Their Variants . . . . .	16
4.3.1	Recurrent Neural Network Justifications and Definitions . . . . .	16
4.3.2	Long Short-Term Memory Recurrent Units . . . . .	19
4.3.3	Gated Recurrent Units . . . . .	19
4.4	Dilation . . . . .	20
4.5	Bi-Axial Architectures . . . . .	21
4.6	The Model . . . . .	24
4.7	Design . . . . .	24
4.8	Implementation . . . . .	24
4.9	Theory . . . . .	24
4.10	Testing . . . . .	24
<b>5</b>	<b>Evaluation . . . . .</b>	<b>24</b>
5.1	Conclusions . . . . .	24
5.2	Internal Comparison of Work . . . . .	24
5.3	Contextualised Comparisons with Existing Solutions . . . . .	24
5.4	Qualitative Surveying Assessment . . . . .	25
5.5	Future Work . . . . .	25
5.5.1	Improved Data Collection and Corpus Creation . . . . .	25
5.5.2	Sentimental Input from Images . . . . .	25
5.5.3	Alternative Architectures . . . . .	25
5.5.4	Performance and Interactivity . . . . .	25
5.5.5	Synthesiser Parameters . . . . .	26
5.6	Author's Assessment of the Project . . . . .	26
	<b>References . . . . .</b>	<b>27</b>

## List of Figures

1	Two dimensional representation of the Circumplex model with labels plotted on the resulting plane . . . . .	14
2	Recurrent connections between the hidden layers of a network . . . . .	17
3	Identical to the previous figure in meaning but with the recurrent connections unrolled along the time axis . . . . .	18
4	A single Long Short-Term Memory Unit taking input at a single time step .	19
5	A single Gated Recurrent Unit taking input at a single time step . . . . .	19
6	Different dilation factors (represented by the value of D) illustrated on a 2-dimensional input . . . . .	20
7	An example of a three-layer DRNN with dilation factors 1, 2, and 4 . . . . .	21
8	Diagram showing the circle of fifths . . . . .	22
9	Example of a transposition from F major to B flat major and D major . . .	22

## List of Tables

### Acknowledgements

My personal tutor Professor Bärbel Finkenstädt and trusted advisor and friend Dr Matthew Leeke alongside my supervisor Professor Graham Cormode have shown consistent support on a personal, professional and academic level throughout my time at the University; for this I owe them great thanks.

The support of my family and friends has been invaluable throughout this process; especially my mother Katie who has always been unwavering in her support of me in the pursuit of my aspirations.

I would also like to thank all of the artists and musicians who gave feedback on my work or were featured in the training corpus; without their tireless work the world would be severely lacking in value and beauty. Most notably, the works of [Ryuichi Sakamoto](#), [Gigi Masin](#), [Jónsi](#) and [Alex Somers](#), [Jamison Isaak](#) and [Brock Van Wey](#) have shaped my years as an undergraduate and offered a great deal of creative and intellectual stimulation as well as instilling a great appreciation for the world and its inhabitants.

## **Abstract**

This dissertation focusses on the use of deep sequential learning techniques to imitate the very human artistic process of musical composition. Numerous techniques and approaches were applied mainly surrounding developments in Recurrent Neural Networks from the past two decades. The project highlights significant progress made by this research and similar papers as well as the work there is still to be done.

(Is talking about RNNs etc. in the abstract too technical, some of the writing sessions seem to think so but I would disagree given it is a key aspect of the project)

The devised model and architecture effectively captures musical structure and harmony in order to compose pieces which subscribe to the general style of the data used to train the model. Moreover, the chosen training data may have its sentiment analysed and attached as a feature during training to allow for retroactive tuning of the model through passing mood parameter values at the time of generation which then influence the mood of the outputs.

The model is versatile and performant, to the point of it being potentially on par with the current cutting edge attempts at making progress in this field; it utilises recent developments to combat some of the issues in previous attempts whilst improving on the training time required to reach equivalent results.

Testing indicates that the compositional engine can be trained with any genre of music and replicate it reasonably well. Future work would focus on improving the sentiment analysis part and improving integration with the model to perhaps focus on certain key defined features of the music.

# 1 Introduction

It is written with regards to the undertaken project exploring the possibilities of musical composition based upon perceived mood of an input and a library of examples with which a compositional engine is to be trained. The main focusses of the project remain on composition as this is the most interesting and challenging component; meaningful contributions to this area would be the most desirable outcome of the project. Given an input from the user, whether it be an image or some text, the main goal is to have a system which will output music to match the perceived sentiment of this input.

This will hopefully lead to the creation of an interactive web application for a user to upload something and receive a composition in return which ‘matches’ whatever inspiration they provide to the compositional engine. The specification document laid out a number of questions to be answered regarding how this composition should be achieved, a task which has multiple components spanning:

- **A means of transcription.** This must be decided upon to encode input and output data for the compositional engine. The main choices to consider were either raw audio data or MIDI data. This decision was a critical one as not only does it influence the choice of approach for composition, but it also dictates the output of the system: training a model with MIDI data means that it is somewhat limited to responding with MIDI itself. This has led to the consideration of one of the **extensions mentioned below**.
- **A means of composing musical pieces based on sentimental input data from a user alongside training data;** this question is heavily influenced by the answer to the one posed above. The main approaches to consider were a more classical and simple approach involving Markov Chains and stochastic state spaces, in comparison with a very modern and active research subject in composition using Neural Networks.

The project has broadened in scope in some ways with respect to the initial specification, whilst narrowing in others. Namely, constraints on the style of output have been relaxed due to potential issues with creating effective training sets; it remains unclear as to whether ambient is one of the hardest genres of music to compose (a lack of rhythmic elements being a major factor in this theory as they are likely to be the easiest elements to generate) or one of the easiest (less emphasis on continuity and coherence to be maintained throughout allows for more abstract output from the compositional engine which is a possibility).



## 1.1 Motivation

Since recurrent neural networks rose to prominence in the 1990s. Deep learning has shown great promise in the task of completely automatic musical composition. This is a complex task exploring some challenging computational and philosophical questions in terms of how the “rules” of music might be correctly encapsulated in order to compose musical pieces to eventually compete with what has always been an almost exclusively human process.

There would be contention in explicitly defining rules for something as complex and subjective as music; this is where deep learning’s relative impartiality comes in and has already proved effective on a number of previously inaccessible tasks where patterns were not thought to be present, or at least were difficult to capture due to the tasks often being very human in nature.

It is here that the motivation for the project lies. Others have been attempting for decades and this project hopes to contribute to that path of exploration, assessing and utilising a variety of techniques in the hope of replicating at least part of the process of composition. In order to do this assumptions must be made in terms of which parts of human behaviour the model is to imitate.

Ambient music would be used primarily because BLAH

### 1.1.1 Brief History of Computational and Algorithmic Approaches to Composition and Musical Applications

Algorithmic composition has developed rapidly in recent years alongside the increase in popularity of deep learning techniques. However, the task and potential solutions pre-date deep learning entirely:

- Markov chains were first formalised in the 1900s and were used to string together segments of score or individual musical notes based on a set of transition probabilities to probabilistically generate sequences of music following conditioning. This conditioning could take place using pre-existing musical score and naturally the outputs tend to follow the inputs closely in terms of style. Iannis Xenakis was a big fan REFERENCE
- Recurrent Neural Networks attempt to extend beyond the main limitations of Markov Chains regarding their restriction to only reproducing subsequences of the original pieces that are used to condition them. These techniques first rose to prominence in the 1980s. The first attempts were usually limited by their lack of coherence beyond

a small neighbourhood of notes; the outputs would often lack any real structure or explode into chaos or vanish into silence.

- In the early 2000s, some improvements on RNNs which aimed to solve some of the aforementioned issues were first applied to composition. The first attempt was made by Doug Eck in 2002 [1] and utilised LSTMs.
- Doug now leads the Magenta team at Google Brain [2] who have created a myriad of models mainly focussing on assisted accompaniment for musicians or improvisation with a user. They have continued to apply LSTMs to these problems as well as variations on the Transformer (a sequence model based on self-attention) pioneered by Huang et al. in 2018 [3].
- The field has fragmented in recent years as teams such as Magenta focus more heavily on interaction with a user. Another fragment is focussing on raw audio over MIDI or character representations as has been the standard for decades. The raw audio approach has huge requirements in terms of data and training time making it still somewhat inaccessible for the time being. Although notable research has been carried out by Google again at their Deepmind front via WaveNet [4]. Other, less creatively focused applications have begun to take hold as well especially within the field of speech synthesis.

There is further discussion to this history if the reader is inclined [4], [5]. The above points cover the main and most relevant milestones and aims to provide some historical context to the starting point of this project.

## 1.2 Research Objectives

Generative probabilistic models such as recurrent neural networks and their many improvements and variations were the main subject of exploration through research. As well as some perhaps simpler interpretations of the above name like Markov chains and fields. Development of a novel architecture which can compete with existing solutions and better them in as many areas as possible is desirable.

The method is inspired by numerous existing works and recent advances in many fields which reveal that deep learning is increasingly effective when the correct data is present. Formulating and utilising an appropriate data set is therefore of the utmost importance.

The goal of this project is to build a versatile model whose outputs are at best indistinguishable from that of a human and at worst impressive in terms of their structure and harmonic

coherence.

## **1.3 Eventual Objectives and Project Requirements**

The main goal

# **2 Background**

## **2.1 Related Work**

### **2.1.1 Competitive Existing Solutions**

In terms of a symbolic approach (something that was not considered due to lack of training data, using musical scoring in terms of sheet music, lettering representing notes), one of the most impressive was undertaken by Sturm et al. from the KTH Royal Institute of Technology and Kingston University in producing a full album with trained musicians and having it reviewed without revealing the nature of its composition. This is a very effective but resource intensive means of qualitative assessment as well as an effective example of a Turing Test for computational composition tasks.

Following on from the questions posed in the Specification, a great deal of research has been carried out in order to better define the associated objectives. Resolutions to most of these questions have been reached following testing or further research, and implementations have been attempted where appropriate in order to properly compare the options outlined previously. Balance in terms of resources required and complexity of the task were the main considerations throughout this process, leading to a realistic but challenging set of goals for the future.

### **2.1.2 Transcription and Encoding Techniques**

The first objective in the composition section of the Specification stated that a decision must be made regarding the format of data used to train the eventual compositional engine. The goal is to build a large library of this data with enough engineered features for the model to successfully produce meaningful output. It has become clear that the use of raw audio data involves a much larger computational overhead making it an infeasible approach given the

time-scale of this project. This decision will be discussed at greater length [in the appropriate section](#); in terms of encoding input data, the decision is mainly based upon the large body of existing research behind MIDI transcription and the intuitiveness of features which arise from relatively simple data such as this. Characteristics such as tempo can be extracted easily and used to train a model on the fundamental ‘rules’ of music, from which the model can then learn to improvise and compose within certain constraints.

Ableton (a company at the forefront of digital audio technology) employee Anna Wszeborska spoke at EuroPython 2016 on the applications of Python in encoding from raw audio data to MIDI [\[6\]](#). Indeed, an implementation to this end has been partially completed, producing some interesting results. Google Brain’s Magenta was heralded in the specification as one of the more promising existing tools: a TensorFlow sub-package focussing on musical composition, transcription and other artistic endeavours. The ‘Onset Frames’ model [\[7\]](#) from Magenta’s library of pre-trained models uses Neural Networks to transcribe piano pieces to MIDI, and also works reasonably well in transcribing other forms of music.

Cross-validation between this model’s output and a method of encoding raw audio to MIDI implemented in Python will be attempted as both approaches have shown themselves to be successful; often one is better than the other depending on the structure of the source. ‘Onset Frames’ is clearly optimal in the situations it was designed for, whilst implementations attempted in Python give mixed results that are more consistent across different styles of music. The biggest challenges encountered in this part of the project are complex sequences of chords and polyphony in the input from different instruments and vocals. Both of the above techniques can deal with these challenges, but with a slightly lower accuracy than quick successions of monophonic notes for example.

Another option which is yet to be considered is the use of Convolutional Neural Networks. These have been applied extensively in the past to the task of musical transcription [\[8\]–\[10\]](#) with state-of-the-art results at their resolutions. It is likely that as the project progresses, work on applying this existing research and exploration into the possibilities of other models will continue and be refined to best match the desired output of training data to be used as inputs for the compositional engine. This process exemplifies one of the largest challenges and complexities in this project: using machine learning to *generate* the data used to then compose music in a similar way introduces a lot of uncertainty and dependence between these two components.

### 2.1.3 Data

All of the approaches to transcription are impractical to a certain extent following on from the concerns regarding time requirements and uncertainty made above; it is for this reason that some large MIDI data libraries such as the Lakh dataset ??? and Metacreation’s corpus [11] could still be fallen back on if time becomes too constrictive. The existence of such datasets allows for work on the compositional engine to begin without delay, but ideally one of the approaches above will eventually lead to the creation of a similar dataset to be utilised in more intricate feature engineering as well as being a more relevant dataset to the genres of music which comprise the initial focus of this project. The MAESTRO Data Set of MIDI corresponding to piano transcriptions [12] is used by Magenta to train their Onset Frames model and so could also form part of the training data for the project.

MIDI is easy to work with and clear in its structure, especially when compared to the complexities of raw audio waveform data. However, there is still the matter of actually using this data to train a model; for this TensorFlow already has a technique for converting MIDI to so-called ‘Note Sequences’ [13] which define all the characteristics of a MIDI file in a Pythonic format. This conversion is painless and allows for more standard feature engineering and manipulation techniques to be done inside of Python using familiar packages.

## 2.2 Composition

### 2.2.1 Markovian Approach

This approach involved implementing a learned state-space based upon interpolations of a series of MIDI files. Markov chains could then be implemented to traverse this state-space and sequentially generate new MIDI, building compositions from probabilistic sequences of notes. The idea was to have an ‘improvisational’ algorithm which could stochastically traverse common progressions and chords, incorporating the ability to switch between these sequences and build a more complex overall piece. Initial work on this approach provided a clear-cut first step to the project as it did not require a huge amount of prior work due to familiarity with the mechanisms and concepts behind it, as well as the existence of other attempts at Markovian composition [14].

The ease of this approach also summarises its main disadvantage: there is a noticeable lack of complexity and novelty in the pieces composed. By its nature, most of the states which the composer can traverse are directly influenced by the input data and thus often lead to

sections which are identical to one of the inputs. This could be attributed to a lack of training data, though using more lead to increasingly small transition probabilities and a descent into near randomness, losing a lot of the musicality present in previous outputs along the way.

It is unlikely that this approach will be revisited and although some output was generated, it was not particularly impressive and certainly pales in comparison to some of the witnessed outputs generated by other techniques. It is for this reason that the second of the big objective questions in the Specification could be answered, concluding that Neural Networks should be used for composition rather than Markov Chains.

### 2.2.2 Neural Network Approach

These are recurrent neural networks which most importantly for my application have the property of time invariance, in that at a given time step the networks activations and learned properties can all be considered relative to previous time states. The outputs at one time step become inputs for the next and we can use this to sequentially train and generate music without worrying about specific locations or times within the network. To offer a counter example that might make this clearer I would imagine a network composed of a fixed sequence of layers which are visited once and then output to the next layer, dealing with absolute positions would make this kind of network unsuitable for musical composition.

Neural Networks are undoubtedly a very active area of research; a lot of which is relevant or could even be directly applied to this project. For example, “How we Made Music Using Neural Networks” [15] references an article by Andrej Karpathy (Tesla’s Director of AI); both of these pieces together formulate a good introduction to the technologies to be implemented for a large portion of the remainder of this project. The Karpathy article [16] showcases some of the capabilities of Recurrent Neural Networks working on generating code, text and images. Music is considered to be a more challenging feat for these systems, which is perhaps intuitive given its continuous and often chaotic nature. The former of the two articles discusses a short exploration into this challenge, a challenge which this project will go further in trying to tackle.

Google’s Magenta project was mentioned as one of the main leads for composition. Magenta’s collection of pre-trained models is testament to the promise of this platform; more specifically there are pre-trained models available which produce improvisation, polyphony and interpolation of input pieces [17], [18]. The aim of this project is to build and train a model sitting somewhere between these existing ones, one which is capable of generating inspired sequences of chords and notes and then recurrently feeding these generations back into itself

in order to emulate improvisation.

Based on Magenta’s current capabilities, an improvisational RNN [19] combined with an LSTM net to improvise off interpolations between existing sources would introduce enough variance and novelty to achieve the goals set for the compositional engine. Preliminary experiments and tests with the models are satisfactory for this application; there are a plethora of provided Jupyter Notebooks to test out the models and interaction is also possible through a command line interface. RNNs are likely to suffer with a lack of global structure, but work well in terms of identifying characteristics of an input and continuing to improvise, in this case. An LSTM would be able to maintain awareness of musical features such as bars, phrases and tempo which should lead to a more acceptable musical output. Hidden layer activations could be used in an RNN to ‘remember’ what it is doing and where it may be in a musical phrase. However, it is yet to be seen how effective this memory will be in practice.

Considering the work done by Magenta’s team and prior explorations into the field by other researchers leads to the conclusion that Neural Networks show more promise for composition than the Markovian approach. This choice links back into the choices made for the format of data to be used to train the models. There are Neural Network based systems which generate MIDI such as Magenta and DeepJazz [20], in contrast to some which instead generate raw audio waveforms such as WaveNet, GRUV and SampleRNN [21]–[23]. As mentioned previously, using raw audio would be infeasible in a project of this scale and so a lot of the tools and options for composition with Neural Networks could immediately be discarded. The outputs from these other approaches were often more abstract as they were not limited to standard musical notation enforced by MIDI, for example notes would often merge into one another without an initial point of attack.

After carrying out research into different Neural Network structures and use-cases, it can be concluded that a Recurrent Neural Network or LSTM would be the most appropriate implementation for this task due to the allowance for different lengths of input and output compared to something like a Convolutional Neural Network which has fixed input and output sizes. Some very promising work was carried out in studies on music generation from MIDI datasets [24], [25] which shows the potential of RNN’s for this task. A similar project called DeepJazz was produced in a hackathon in a matter of hours and also gave very promising results, again using an RNN.

## 2.3 User Input and Interface

This component of the project is relatively simple and thus has not been given a great deal of consideration at this stage. A Neural Network approach to image analysis similar to the one described above would allow for the definition of features and characteristics required for the model to gain an understanding of sentiment in music. This could be achieved through labelling and classification of data into different ‘moods’ or similar groupings which could eventually lead to it matching an input’s perceived mood with the type of output it creates. The analysis of images for features is a much simpler task than music; this area of research is also very active, and many examples of such analysis can be found on the internet. It is likely that a lot of the theory and some of the implemented components from the compositional engine could even be reused in image analysis; although here a CNN would likely be more appropriate than an RNN due to the stationary nature of image data.

As mentioned briefly in the report, an image as input is not the only planned way for a user to interact with the compositional engine. It should also be possible for a user to provide textual input or manually set parameters which will again influence the compositional engines output.

## 2.4 Articles of Interest in the Field

It was clear in the initial stages of this project that the field of computational musical composition is one of intense and rapid development, with a lot of the papers cited being from 2015 and later. Numerous sources of information have been investigated and have proved conducive in terms of inspiration moving forward with the project. For example, the research and updates to Google Brain Team’s Magenta Blog [2] have been of great interest throughout the project’s development so far as there is a lot of objective overlap.

The general discussion in this blog alongside that found in articles on the internet [4], [26] helped influence the decision to pursue more complex Neural Network based approaches to composition. Both of these cite Markovian techniques as a *starting point* for this area; something which has quite safely been surpassed in every respect at this point. An article discussing comparisons of various deep learning tools for music generation [27] was also informative, highlighting Magenta as a tool of great promise.

Another deep learning tool encountered is GRUV [28], which was mentioned earlier as being an approach trained on raw audio data. It has been used to successfully reproduce some very recognisable snippets of music such as the ‘Amen Break’ (part of a funk record which



has been sampled countless times in popular music) from a set of training data, as well as producing individual instrument sounds [29]. The reproduction of individual instruments could lead to an alternative approach to synthesis and production; layering generated sounds on top of a MIDI composition would be an interesting challenge if the different instruments could somehow be associated with sentiment.

## 3 Methodology

### 3.1 Research Methodology

The first stages of this project in particular had stringent research requirements in order to formulate an effective approach. There was a large body of work to investigate, much of which consists of research from the past couple of years meaning much of it is still somewhat theoretical or lacking in implementation options. This imposed a requirement of careful thought in terms of which avenues might be worth exploring and which were feasible with the given resources and time-frame.

Due to the scale of the task taken on, and the technical requirements, it was decided that first formulating some research questions would be the best approach.

### 3.2 Project Management

Agile in places, refer to the initial posed questions and answer them to iteratively come up with a novel and effective approach supported by research and assessment.

PROJECT DEVELOPMENT this is more about the questions posed and how the project crystallised over time.

Meetings for advice and aid in resolving some of the key issues in the project have taken place roughly once or twice per week with the project's supervisor. This is planned to continue, with the meetings likely to increase in frequency due to a much lower external workload in the next academic term.

With most of the preliminary and supporting research complete, as well as the main few directional questions answered, development may continue in earnest. Below is an updated Gantt chart to show how the project's trajectory and schedule has changed since the initial specification, as well as showing confirmation of the completion of some of the defined

tasks. There are some slight changes following the issues or decisions discussed earlier in the document, but for the most part the schedule remains unchanged.

### 3.3 Development Methodology

DEVELOPMENT actual programming was iterative combined with agile I guess, rushed to MVP in order to gain better idea of foreign concepts.

Cues were taken from previous development experience alongside agile methodologies to try and accommodate for a flexible, research-driven development process. Deep learning lends itself to trialling ideas and performing evaluation due to the black box nature of some of its components; the task itself does not have a “fixed” optimal solution implying that the development process would have to accommodate for changing requirements and focusses on improvement rather than reaching some pre-determined threshold.

### 3.4 Ethical Considerations

There are no major ethical considerations in this project. All of the potential issues are briefly discussed here.

All citations have been checked and are included where necessary to ensure credit is given where existing work has been reused or had a significant influence on the course of the project. The code used for sentiment analysis was based upon the DeepSent project as already mentioned which is present on GitHub under an open source Apache 2.0 license via its creator Mu Chen. The rest of the code included in the technical appendix is referenced where appropriate or entirely original.

A small survey was devised and carried out in order to test a subject’s ability to differentiate between human composed pieces and those of the compositional engine. This survey was carried out with willing participants over the internet by sending them a mix of unmarked audio files all synthesised using the same virtual piano. The user is simply asked to categorise them into human and non-human composition folders before sending them back. This ensures the study was unbiased and also offered the chance to ask for feedback as to the quality of the pieces once a participant had carried out the first task. The results were recorded over a period of 2 weeks and are presented in the evaluation section.

This survey was discussed with the project’s supervisor but was deemed to be satisfactory in terms of any ethical considerations; no further justification or recourse was necessary to

ensure the efficacy of it as a means of evaluation.

## 4 Technical Component

### 4.1 MIDI Transcription

In order to determine which MIDI transcriptions were best, a number of factors were considered. Some preliminary testing was done by first designing a sequence of MIDI notes which could then be played internally and exported as an audio file. These MIDI notes acted as a ground truth. Within these test files, chords and complex structures were included as well as different levels of white noise applied to them. From here different techniques could be tried in order to decide which would be best:

- WaoN ??? is a transcription tool which converts `.wav` audio files to `.midi` files. It carries out frequency-domain analysis using Fast Fourier Transforms which are computationally intensive but are flexible and accurate when compared to simpler autocorrelational techniques [klapuri2004automatic; gerhard2003pitch] meaning they are often applied to more complex tasks such as pulling out polyphonic features in music as was the goal of this part of the project. This technique requires no training data and so can be quickly applied to a large array of audio; it is not restricted to a specific genre according to its creator and even has means to try and remove drum sounds when transcribing.
- Magenta’s ‘Onset Frames’ model [7] is considered cutting edge and utilises a convolutional recurrent neural network architecture to predicts pitch events both in a framewise and onset context. I.e. it first predicts where notes may begin and then uses these predictions to influence predicted frames where a certain pitch is present. Where these predictions are in agreement (an onset is predicted for that pitch within the predicted pitch frame) a note is presumed to be present and can be transcribed to MIDI. Training data and scope again rears its head as an issue with this approach, as is mentioned in the referenced paper.

Based on the transcriptions both create it is possible to compare them qualitatively. Magenta’s approach is decisive in its choice of notes and often more confident when it comes to timing due to its training instilling a greater sense of tempo and structure in its transcriptions. WaoN is delicate and plays with wide dynamic range, but often incorporates more false positives in terms of ghost notes (notes that should not be present but are determined to be through frequency clashes, overtones etc.). Both seemed to be viable and offered slightly

different versions of the training data to be used in training the compositional engine, and as both were viable it was decided that both should be used. Magenta’s model shows great promise and it is likely that with the correct training data this approach would be far superior to WaoN in fixed domain transcriptions (i.e. training the model on a dataset of ground truth transcriptions of a certain genre before using it to transcribe more of that genre, so that it might understand the ‘rules’ of transcribing genres other than classical more effectively). I have no doubt that in the next few years applications such as WaoN will be largely eschewed in favour of modern machine learning approaches similar to what has happened to the field of Computer Vision over the past decade or so.

#### 4.1.1 Implementation

Most of the code relevant to this section is present in the `dataset.py` and `conversion.py` files included with this document.

## 4.2 Sentiment Analysis

In order to satisfy the requirements of the project regarding sentiment, it was necessary to first extract the sentiment from each component of the training data corpus and then formulate a means of attaching this sentiment to the representation of the training data to be fed into the model during training.

Much of the technical work regarding extraction of sentiment from audio follows the work done by Mu Chen and their DeepSent project [30]. This code was used as a basis upon which a vector of mood values could be generated automatically for each element of a training data set and then associated with the accompanying MIDI files to be used in training. The sentiment is measured with respect to an arousal-valence emotional model known as the Circumplex model [31], with values being determined for both axes of arousal and valence.

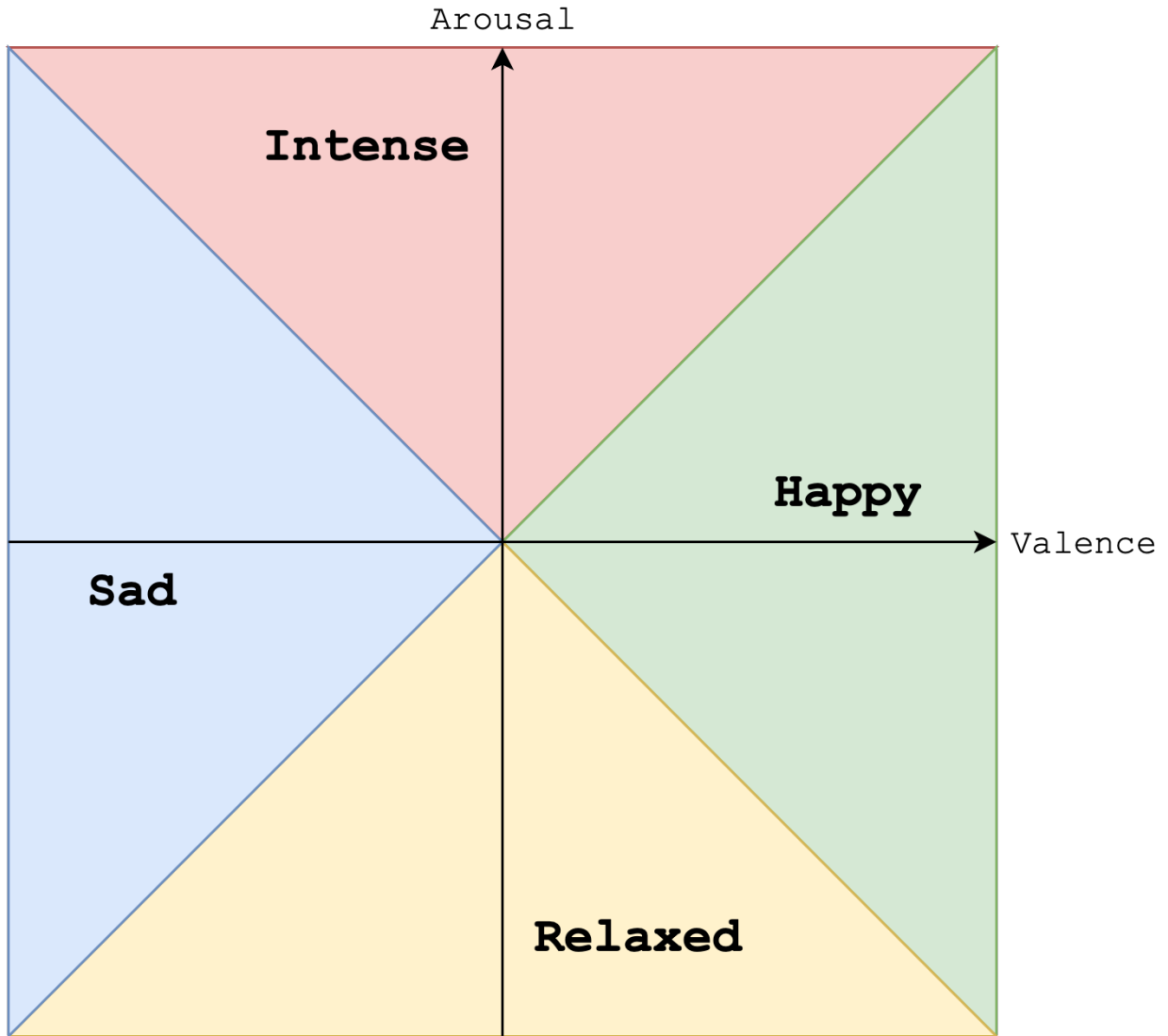


Figure 1: Two dimensional representation of the Circumplex model with labels plotted on the resulting plane

In this context, arousal refers to the perceived intensity of a piece; ranging from a relaxed or somewhat lethargic feeling to a more intense and excited stimulation. Valence is a spectrum of the affect on a listeners level of happiness or sadness which may be incited by the piece. Values for both axes were estimated and ratios calculated to fulfil three ratio values regarding each axis. The ratios are calculated as shown below following predictions made by the models:

$$\begin{aligned}
S_A &= \{\text{Scores assigned to each input by the arousal regressor model}\} \\
\text{Arousal Intense Ratio} &= \frac{\sum_{S_A} \mathbb{1}_{\text{Score} > 2}}{|S_A|} \\
\text{Arousal Relaxing Ratio} &= \frac{\sum_{S_A} \mathbb{1}_{\text{Score} < 1}}{|S_A|} \\
\text{Arousal Mid Ratio} &= 1 - (\text{Arousal Intense Ratio} + \text{Arousal Relaxing Ratio})
\end{aligned}$$

$$\begin{aligned}
S_V &= \{\text{Scores assigned to each input by the valence regressor model}\} \\
\text{Valence Happy Ratio} &= \frac{\sum_{S_V} \mathbb{1}_{\text{Score} > 2}}{|S_V|} \\
\text{Valence Sad Ratio} &= \frac{\sum_{S_V} \mathbb{1}_{\text{Score} < 1}}{|S_V|} \\
\text{Valence Neutral Ratio} &= 1 - (\text{Valence Happy Ratio} + \text{Valence Sad Ratio})
\end{aligned}$$

In both cases it can be seen that ratios are calculated through taking the top and bottom few scores and calculating the ratio of each above and below boundaries which were formulated through testing and with respect to the original work. It is worth noting that all scores fall between a range of 0 and around 5, but extreme values are clipped to 3 as they do not appear and are often spurious.

#### 4.2.1 Implementation

Most of the code relevant to this section is included in the `mood.py` file included with this document.

Each piece to be analysed is trimmed of its start and end (which are likely to be sparse and potentially irregular in terms of their content when compared to the main body of a piece occurring around its midpoint) to leave the middle 70% of the data. This raw data is then windowed into 5-second long frames of samples, with a step of 0.5 seconds between the starting points of each (such that there is overlap between frames). Within each of these there is then a further decomposition into 25ms sub-frames that are transformed into an array of Mel-frequency cepstral coefficients (MFCCs). The matrices representing the coefficients for each sub-frame are then flattened and fed into a 3-layer neural network regressor. The neural networks associate patterns common within certain moods based upon what was learnt from numerous large datasets used for training these models.

These ratios are then packed into a 6 item vector and associated with each of the training MIDI files to be loaded during training. Thresholding as a means of segmentation was implemented by setting all ratios greater than 50 to 1 and those below 50 to 0. This offered computational improvements but slightly decreased the effectiveness of the mood transfer to training pieces in the opinion of the author and some survey participants. Due to the subjectivity of this part of the output, it was hard to assess which approach was more effective in capturing mood and so the simpler, less computationally intense thresholding option was chosen. This resulted in a six element *binary* vector representing mood associated with each item in the training corpus.

## 4.3 Recurrent Neural Networks and Their Variants

### 4.3.1 Recurrent Neural Network Justifications and Definitions

Perhaps the most common form of artificial neural network is the *feedforward* network which utilise connected layers of neurons and activation functions to approximate different functions through the learning of parameters and weights. Their main limitation in this context is that they require their inputs to be of a fixed dimension and thus are not well suited to dynamic data of varying size such as a sequence  $(x_i)_{1 \leq i \leq T}$  of time stepped note states representing a composition where  $T$  is the length of the piece.

It is assumed that this type of sequential data has some system of dependence on its prior and potentially future elements and so it would not be appropriate to simply input each element of a temporal sequence into a feedforward network during training. It can clearly be seen by observing musical data and considering knowledge of musical harmony (songs are often in a certain and consistent key which may be inferred by chords and notes present throughout the piece) that the input at each time step is likely to be dependent on other time steps.

This class of situations led to the development of **Recurrent Neural Networks** which introduce recurrent connections between layers over a temporal dimension allowing the network to exhibit dynamic behaviour in this dimension.

It is now possible to provide some notation specific to this project regarding RNNs. A sequence of inputs is defined as above, with the vector corresponding to the sequence at time  $i$  denoted as  $x_i \in \{0, 1\}^{D_{\text{input}}}$ .

At the start of the process all weights and activations are initialised to some value. At each time

step following this the new activations for each hidden layer are calculated using a combination of the current time steps input and the previous activations. This process highlights the possibility of unrolling of an RNN into a feedforward network or DAG representation as shown in the figure below.

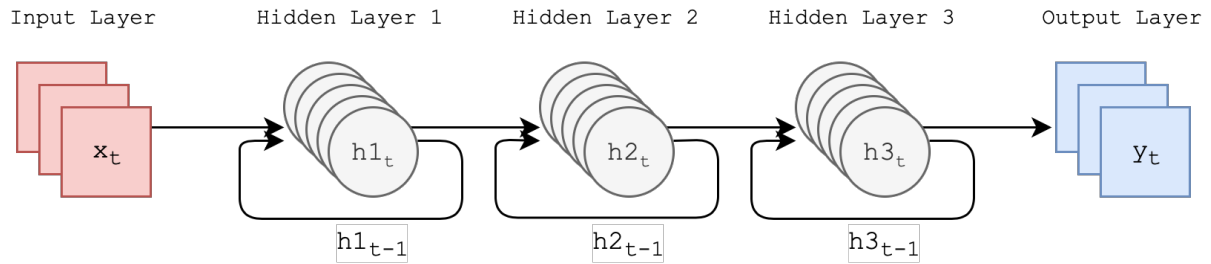


Figure 2: Recurrent connections between the hidden layers of a network



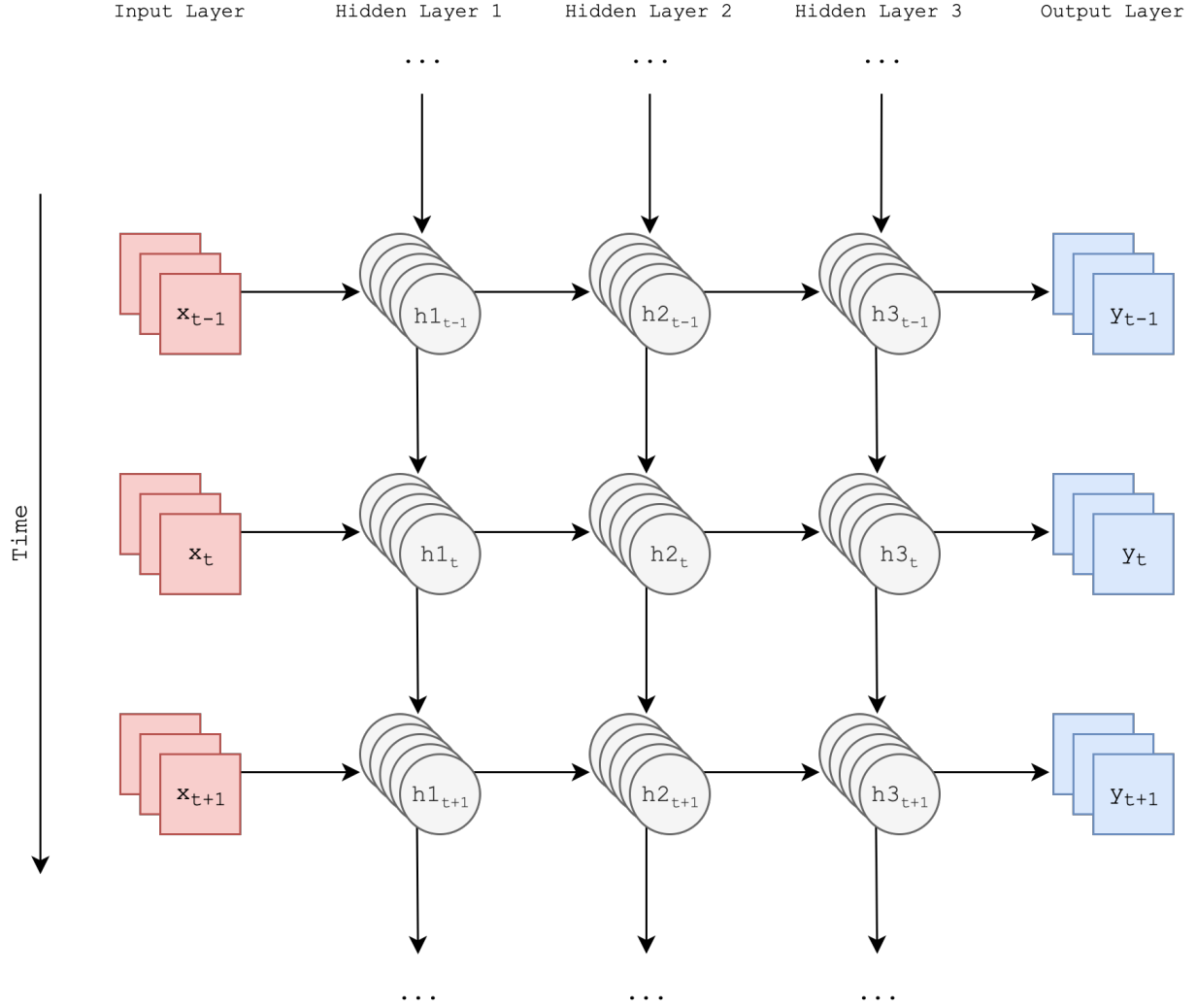


Figure 3: Identical to the previous figure in meaning but with the recurrent connections unrolled along the time axis

The most important property of this class of network is their *time invariance*, in that at a given time step the networks activations and learned properties can all be considered relative to previous time steps. The activations at one time step influence the next. The ability to recurrently input into the network is the clear reason for its usage in this context over traditional feedforward networks.

Despite this time invariance and flexibility in terms of the dimension of their inputs, they lack long term coherence meaning they often fail when complex dependencies are built up and different temporal structures must be captured. They are susceptible to exploding or

vanishing gradients over time due to repeated and recurrent backpropagation, i.e. tiny values will often compound and be multiplied together leading to the network getting stuck or significantly slowing its learning. These issues have inspired a number of variants aiming to mitigate these problems and form the basis of most current approaches to musical composition. They are discussed in the sections below.

### 4.3.2 Long Short-Term Memory Recurrent Units

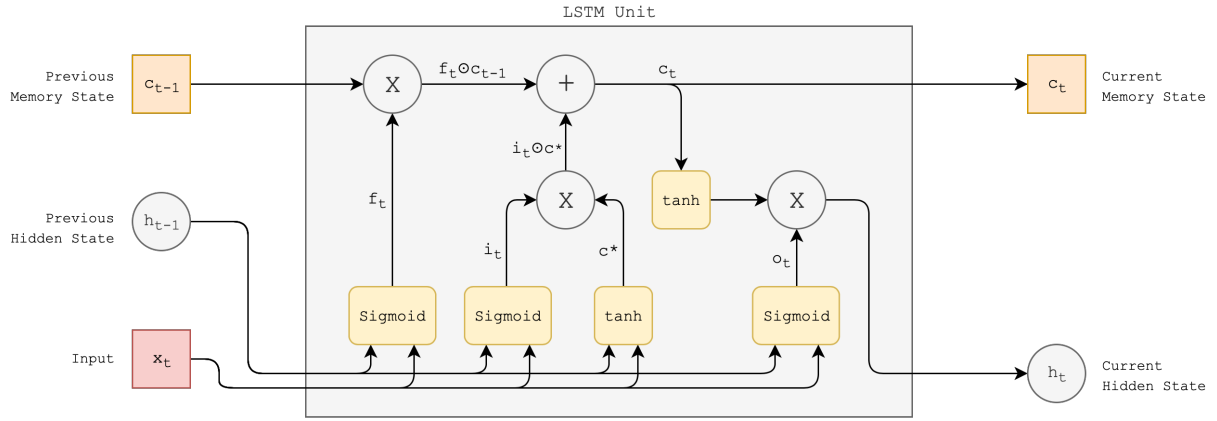


Figure 4: A single Long Short-Term Memory Unit taking input at a single time step

### 4.3.3 Gated Recurrent Units

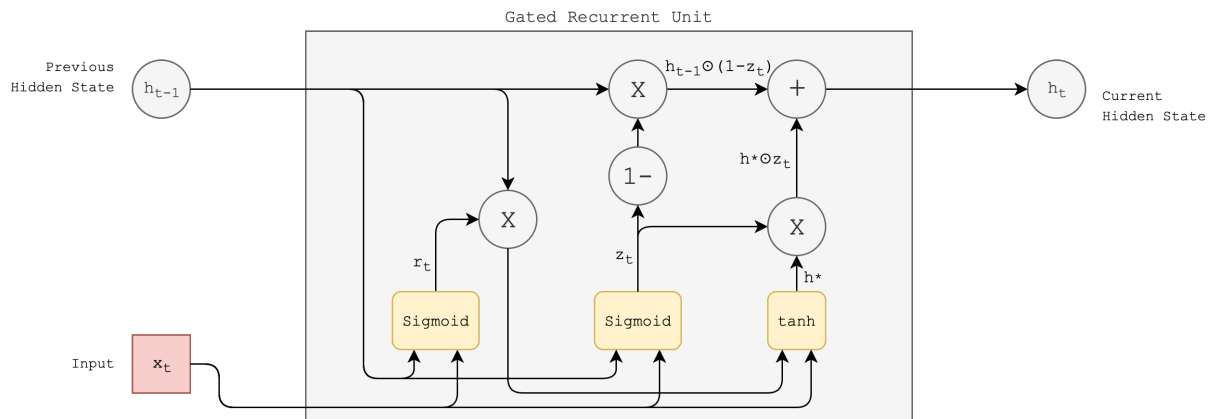


Figure 5: A single Gated Recurrent Unit taking input at a single time step

## 4.4 Dilation

The concept of dilation was first proposed in the context of convolutional neural networks [32] for image analysis and semantic segmentation; this work was discovered during research for a **potential extension** of this project. The main concept is to aggregate information at different contextual scalings without losing resolution by exploding a kernel's considered neighbourhood around a central pixel / element. This is done to increase the likelihood of discovering pattern structures at different resolutions within an input as well as increasing the area of an image or input which can be considered through a small number of steps.

It offers an alternative to other techniques often relying on down-sampling which sacrifice resolution rather than considering different contexts at full resolution as is the case with dilation.

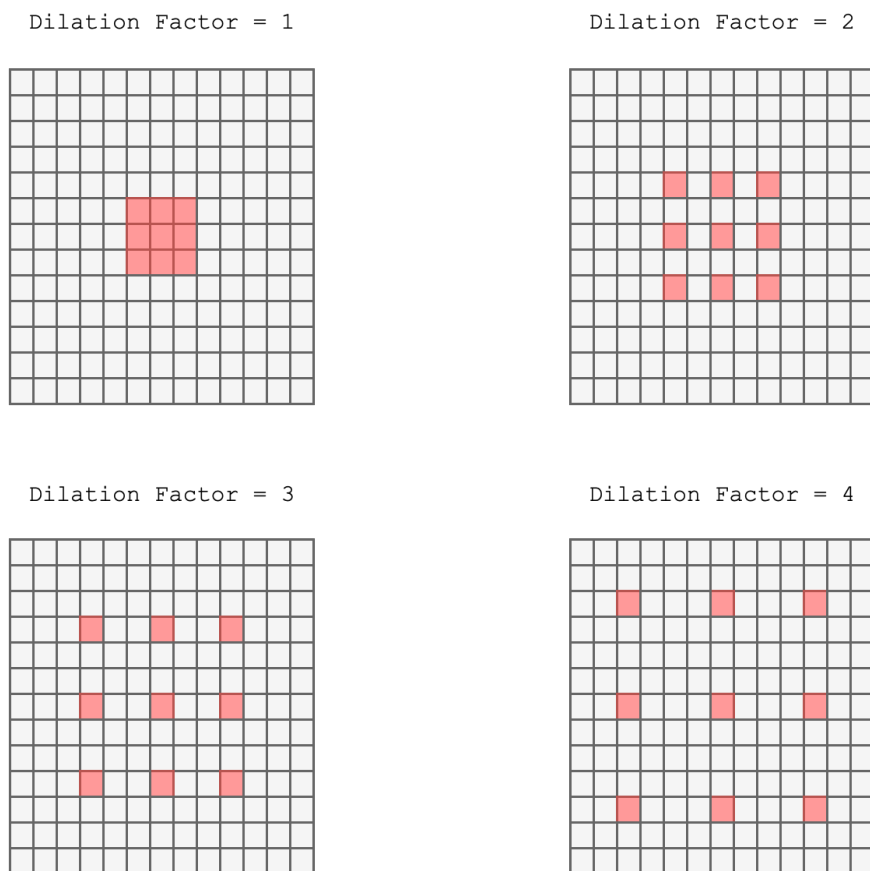


Figure 6: Different dilation factors (represented by the value of  $D$ ) illustrated on a 2-dimensional input

This technique has shown to be very effective in aiding dense prediction problems (predicting

labels for each pixel in an image, or with respect to this project’s proposed equivalence of predicting on or off states for each note on the piano over a series of time steps). Deepmind researchers applied a similar approach within their own convolutional network for WaveNet which is the first documented use of dilation in the musical composition domain [21]. Their results found that again the addition of dilation conceptually increased the accuracy and effectiveness of their model.

A paper was published linking this concept with recurrent neural networks in 2017 by Chang et al. [33]. This paper forms the basis for the justification of using dilation in the model present in this project; this project’s compositional model is potentially the first use of dilated recurrent neural networks in the musical composition domain.

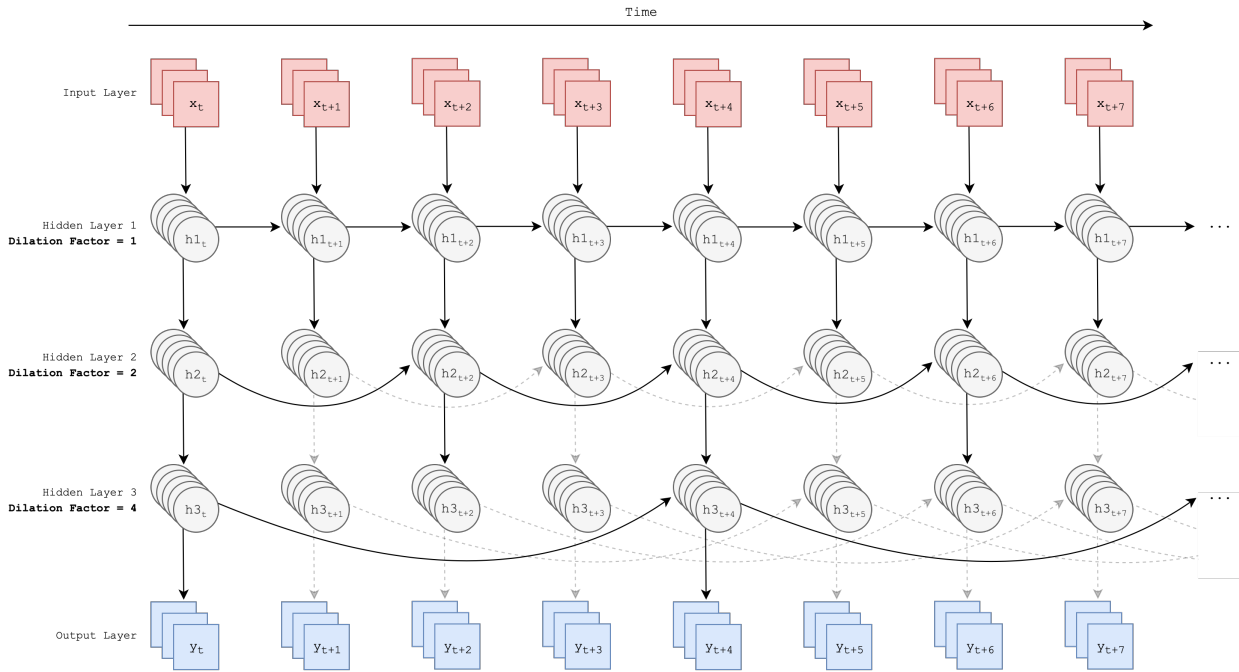


Figure 7: An example of a three-layer DRNN with dilation factors 1, 2, and 4

## 4.5 Bi-Axial Architectures

Consideration of musical theory is the main inspiration for the following section. There is a large body of documentation and research surrounding musical theory which a reader might find helpful in adding context to this paper; the essentials are included for convenience.

In general, compositions are written with respect to a **key** which roughly determines the

scales upon which harmonics and chords for a piece are constructed. Each key is simply a transformation or rather a transposition of another through some number of shifts up or down.

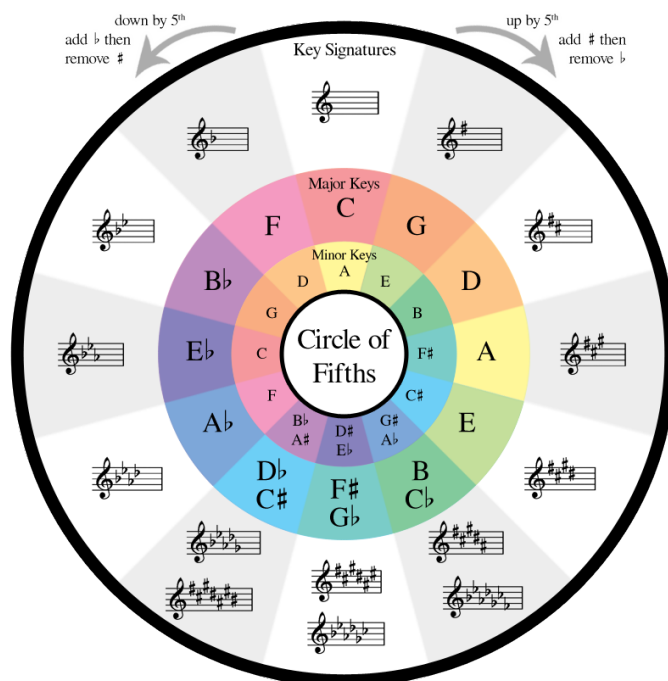


Figure 8: Diagram showing the circle of fifths

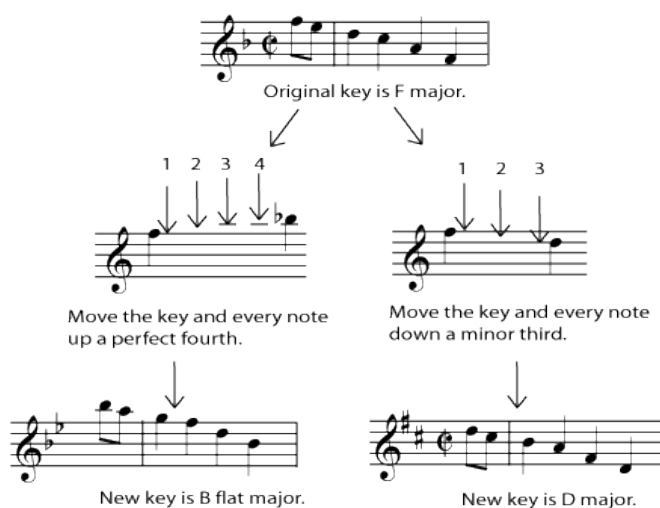


Figure 9: Example of a transposition from F major to B flat major and D major

This highlights a potential issue with many pre-existing attempts at composition using neural networks and limitations of some distribution estimators such as Restricted Boltzmann

Machines and Neural Autoregressive Distribution Estimators. All of the harmonics present in music (aside from some intentional dissonances) are entirely relative by nature. For example, if we represent all possible played MIDI notes as a binary vector where 1 represents a played note and 0 represents the lack of a note at this location. A major triad chord can be represented as:

...0100010010...

Regardless of its position in the input sequence. It is known that the relative position of notes represented by some subsequence of an input matching the locations above would result in a major triad of some key. This highlights the property of *harmonic invariance* in music

Time invariance

5. Image processing has some parallels to my approach as well. I wanted to introduce more concepts of relativity into my model over absolute structures beyond just the time thing; many existing solutions include networks where each note corresponds to a specific node. However, music is more of a relative set of rules in that pieces can be transposed up or down in key and you can always construct chords and melodies relative to that key. So I wanted my network to be able to capture this understanding rather than just being able to copy certain chords it learns; this way it can perform in a more varied and realistic way. Image processing has a similar idea already in which a kernel is used to take a neighbourhood of pixels in an image and carry out some transformation on them. My idea was to do this with an octave above and below each note, feeding this neighbourhood into its own network which would then focus only on the relative position of notes to learn harmonic structure without worrying which note was actually being played.

## **4.6 The Model**

## **4.7 Design**

## **4.8 Implementation**

## **4.9 Theory**

## **4.10 Testing**

Functional testing, unit testing, integration testing ? testing strategies

# **5 Evaluation**

## **5.1 Conclusions**

A favourable configuration was found though further investigation would be required. The function aimed for is a complex one and it is still somewhat unclear how many neurons might work best in modelling the full potential. Here resources of time and computational power must be considered.

## **5.2 Internal Comparison of Work**

Compare different layer configurations, pyramid and number of units and also GRU and LSTM.

## **5.3 Contextualised Comparisons with Existing Solutions**

Insert the tables of comparisons with existing solutions

## 5.4 Qualitative Surveying Assessment

## 5.5 Future Work

### 5.5.1 Improved Data Collection and Corpus Creation

As has already been alluded to, perhaps the biggest issue faced throughout this project has been one of data. New architectures were developed and tested iteratively but all of them faced similar limitations due to a lack of data. The model itself proved to be close to state-of-the-art through its training performance and quantitative evaluation. However, many of the pieces still seemed to leave something to be desired; especially when the created ambient corpus was used

### 5.5.2 Sentimental Input from Images

### 5.5.3 Alternative Architectures

Look into SRUs would be cool, potential improvement over GRUs, moving averages present potential relevance to music as a format.

### 5.5.4 Performance and Interactivity

This project proves that additional features can effectively be incorporated into the model. Indeed, a mood may be input at the time of generation in order to influence the model's output. A natural extension of this would be to provide additional ways to interact with the model at the point of generation through giving a user the chance to manually inspire the model. This could potentially be as intuitive as allowing a user to play a short number of notes or chords and have these set the initial weightings for the network, which would presumably go on to compose a whole piece from these starting conditions. This would certainly be possible given a greater corpus of training data to allow for the model to have a greater understanding of shorter inputs; could allow for musicians to use the model creatively in terms of inputting short ideas into it and exploring its responses.

Many of the Magenta team's efforts since the launch of interactive TensorFlow.js have incorporated ideas such as these; it is relatively simple to hook up a MIDI interface into the model such that it could receive input in a way similar to how some of the Magenta models do.



In order to further build upon the user experience of this project, a simple web app could be built to contribute to a larger population size for the **surveying method** described in the evaluation section. The web app would allow for a user to input a mood either by setting the Circumplex parameters manually or potentially tying in the use of an uploaded image or even NLU to analyse the mood of a sentence and then feed this into the model to generate an appropriate piece. The hosting of a pre-trained model would be relatively simple to achieve and again was not deemed to be of sufficient value to justify the time and resource cost during this project as it was not key to the initially laid out requirements. Despite this, it could potentially gain traction online and provide a much larger sample size for automated collection of qualitative assessment data regarding the model.

#### 5.5.5 Synthesiser Parameters

Other researchers have successfully trained a model to change parameters of a synthesiser during live performance based on a musician's input and learned preferences [34]. Moving forward this would be an excellent way to add another dimension to the project's outputted music; existing studies into the sentiments behind different timbres and urgency or latency introduced by manipulations of note attack-delay-sustain-release parameters could form a basis of some quantitative assessment of the results. The input could again influence the output through affecting the choice of parameter programming for the synthesiser used to play the compositions the system produces.

One of the limitations in choosing to train a model on MIDI as mentioned earlier is that the output is also somewhat restricted to this format. At which point choices could be made in terms of how this MIDI is synthesised and played back to the user. This could range from something as simple as instrumental selection, to the tuning of parameters of a chosen synthesis engine as described here.

### 5.6 Author's Assessment of the Project

The work discussed indicates the level of technical achievement this project encompasses. Much of the work is of a high level for a third year project and involved significant time investment to learn additional advanced material with which to accomplish the goals of the project.

The use of neural networks and deep learning is adherent to the description of my degree; probabilistic sequence generation is certainly relevant in almost all aspects to the statistical

and computational nature of the BSc in Data Science. The project illustrates domain knowledge spanning multiple fields relevant to data science and involved numerous challenging components from a research, theoretical, technical and engineering perspective.

It is hoped that this work forms a relatively robust and all-encompassing summary of not only the outputs of the project but the work involved in achieving it. Other academics and students will hopefully find value in this work and be able to continue to contribute to this fast-moving field. Musical composition as mentioned is a sufficiently complex task to fully leverage some of the more advanced deep learning and sequential modelling techniques which are currently of interest in academia; this work shows their potential applications as well as hopefully contributing to their further development.

Why should this project be considered an achievement?

The incorporation of mood is perhaps the initial limitation an observer may encounter in terms of the initial goals of the project.

## References

- [1] D. Eck and J. Schmidhuber, “Finding temporal structure in music: Blues improvisation with lstm recurrent networks,” in *Proceedings of the 12th ieee workshop on neural networks for signal processing*, 2002, pp. 747–756.
- [2] Google Brain Team, “TensorFlow Magenta.”  
<https://github.com/tensorflow/magenta>.
- [3] C.-Z. A. Huang *et al.*, “An improved relative self-attention mechanism for transformer with application to music generation,” *arXiv preprint arXiv:1809.04281*, 2018.
- [4] K. McDonald, “Neural nets for generating music.”  
<https://medium.com/artists-and-machine-intelligence/neural-nets-for-generating-music-f46dffac21c0>.
- [5] Y. Bayle, “Deep learning for music chronicle.”  
<https://github.com/ybayle/awesome-deep-learning-music>.
- [6] A. Wszeborowska, “Music Transcription with Python.”  
<https://www.youtube.com/watch?v=9boJ-Ai6QFM&feature=youtu.be>.
- [7] Google Brain Team, “Magenta - Polyphonic Piano Transcription.”  
<https://magenta.tensorflow.org/onsets-frames>.

- [8] M. Bereket and K. Shi, “An ai approach to automatic natural music transcription.”
- [9] mbereket, “Music Transcription Repository.”  
<https://github.com/mbereket/music-transcription>.
- [10] S. Sigtia, E. Benetos, and S. Dixon, “An end-to-end neural network for polyphonic piano music transcription,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 24, no. 5, pp. 927–939, 2016.
- [11] M. Lab, “Metacreation lab.”  
<http://metacreation.net/corpus-1/>.
- [12] C. Hawthorne *et al.*, “Enabling factorized piano music modeling and generation with the maestro dataset,” *arXiv preprint arXiv:1810.12247*. 2018.
- [13] Google Brain Team, “NoteSequence Conversion Guide.”  
<https://github.com/tensorflow/magenta/blob/master/magenta/scripts/README.md>.
- [14] anbud, “Markov composer (java application).”  
<https://github.com/anbud/MarkovComposer>.
- [15] A. Tavgen, “How we made music using neural networks.”  
<https://medium.com/@ATavgen/how-we-made-music-using-neural-networks-449a62b8a332>.
- [16] A. Karpathy, “The unreasonable effectiveness of recurrent neural networks.”  
<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>.
- [17] Google Brain Team, “Magenta - Music VAE Model.”  
[https://github.com/tensorflow/magenta/tree/master/magenta/models/music\\_vae](https://github.com/tensorflow/magenta/tree/master/magenta/models/music_vae).
- [18] Google Brain Team, “Magenta - Polyphony RNN Model.”  
[https://github.com/tensorflow/magenta/tree/master/magenta/models/polyphony\\_rnn](https://github.com/tensorflow/magenta/tree/master/magenta/models/polyphony_rnn).
- [19] Google Brain Team, “Magenta - Improvisational RNN Model.”  
[https://github.com/tensorflow/magenta/tree/master/magenta/models/improv\\_rnn](https://github.com/tensorflow/magenta/tree/master/magenta/models/improv_rnn).
- [20] J.-S. Kim, “DeepJazz.”  
<https://github.com/jisungk/deepjazz>.
- [21] A. van den Oord *et al.*, “WaveNet: A generative model for raw audio.” 2016 [Online]. Available: <http://arxiv.org/abs/1609.03499>
- [22] A. Nayebe and M. Vitelli, “GRUV : Algorithmic music generation using recurrent neural networks,” 2015.

- [23] S. Mehri *et al.*, “SampleRNN: An unconditional end-to-end neural audio generation model.” 2016 [Online]. Available: <http://arxiv.org/abs/1612.07837>
- [24] M. Hilscher and N. Shahroudi, “Music generation from midi datasets.”
- [25] L. Wyse, “Real-valued parametric conditioning of an rnn for interactive sound synthesis,” *arXiv preprint arXiv:1805.10808*, 2018.
- [26] M. Kofler, “Deep Learning with TensorFlow.”  
<https://towardsdatascience.com/deep-learning-with-tensorflow-part-3-music-and-text-generation-8a3fbfdc5>
- [27] F. Brinkkemper, “Analysing Six Deep Learning Tools for Music Generation.”  
<http://www.asimovinstitute.org/analyzing-deep-learning-tools-music/>.
- [28] A. Nayebi and M. Vitelli, “Gruv: Algorithmic music generation using recurrent neural networks,” *Course CS224D: Deep Learning for Natural Language Processing (Stanford)*, 2015.
- [29] Fiala, “Generating Audio with Deep Learning.”  
<http://fiala.uk/notes/deep-learning-and-sound-02-generating-audio>.
- [30] Y. Bayle, “Deep learning for music chronicle.”  
<https://github.com/muchen2/DeepSent>.
- [31] J. A. Russell, “A circumplex model of affect.” *Journal of personality and social psychology*, vol. 39, no. 6, p. 1161, 1980.
- [32] F. Yu and V. Koltun, “Multi-scale context aggregation by dilated convolutions,” *arXiv preprint arXiv:1511.07122*, 2015.
- [33] S. Chang *et al.*, “Dilated recurrent neural networks,” *arXiv preprint arXiv:1710.02224*, 2017.
- [34] N. Sommer and A. Ralescu, “Towards a machine learning based control of musical synthesizers in real-time live performance,” in *Proc. Of the 25th modern artificial intelligence and cognitive science conf., spokane, washington, usa*, 2014, pp. 61–67.