
Context-Based Music Composition using RNN Generative Model (Team Name: Playlist)

Apurva Pathak
A53097569
appathak@ucsd.edu

Kshitiz Gupta
A53104364
ksg005@ucsd.edu

Chaitanya Baratam
A53104872
cbaratam@ucsd.edu

Ritvik Jaiswal
A53090298
rjaiswal@eng.ucsd.edu

Puneeth BommiReddy
A53093725
pbommire@eng.ucsd.edu

Abstract

In this paper, we present a recurrent neural network that generates music compositions based on auxiliary information such as genre (Irish and Swedish). We represent the music in Abc notation, a text-based music notation system, so that a character-level recurrent neural network could be used for generating music composition. Our results are based on a corpus of over 2500 music files collected from Henrik Norbeck’s Abc Tunes. Our network, composed of LSTM units, produces music in Abc notation tailored to a genre. Qualitative evaluation demonstrates that the model learns the rhythm, keys and tempo of music from the training music files. Since our model generates music one character at a time, it also adds the metadata information into the generated music file like title, composer, etc. without any machinery explicitly dedicated to the purpose.

1 Introduction

Our work is motivated by the success of *BeerMind* [5] which uses a character-level Recurrent Neural Network (RNN) to synthesize beer reviews tailored to star rating or category using an auxiliary input fed to the network at each time step. Our goal is to use a similar model to generate different genres of music using the context as the auxiliary input. Through a character based representation of musical notes as training data, we generate a sequence of notes using varied RNN architectures and compare their performance. Further, we attempt to generate different genres of music using an auxiliary input that is fed to the network at each time step. We then visualize the learned parameters of our networks to understand what features the network has learnt.

Recurrent Neural Networks (RNNs) trained on the word-level i.e. with a word as input at each time-step suffer from the computational complexity that accompanies the very large vocabulary of legal words in any given context [5]. This is where character-level RNNs come in. They have been shown to generate coherent passages but suffer from being unable to decide which set of the previous inputs are relevant to the current output i.e. they are unable to effectively model the *context* [5][12]. They have been used to generate music with limited success. They fail to capture global structure, which most often manifests as generation of repetitive segments of text or music [1]. Being unable to decide when to stop generation are problems with this model. This can be combated by augmenting the vocabulary with suitable start and stop tags. Boulanger Lewandowski, Bengio, Vincent (2012) [3] talk about the difficulty in training RNNs and corresponding optimization techniques.

Character-level Long Short-Term Memory Networks (LSTMs) demonstrate the ability of RNNs to model sequences on multiple time scales simultaneously, i.e., they learn to form words, to form sentences, to generate paragraphs of appropriate length, etc. This makes them able to determine

the context effectively as compared to a vanilla RNN. Nayebe and Vitelli (2015) [4] conclude that LSTMs can successfully model long term dependencies present in musical sequences effectively. Eck and Schmidhuber (2002) [2] used an LSTM based model for music composition to successfully learn the global structure of a musical form and use that information to compose new pieces in that form. The model was used to generate new instances of a bebop-jazz variation of standard 12-bar blues.

However, we wish to add an auxiliary input to determine the kind of music it should generate. Since the LSTM could possibly *forget* the auxiliary input in a few steps. Lipton, Vikram and McAuley propose that this context deciding auxiliary input be concatenated along with the input all all time steps. While it might seem redundant to replicate the input at each sequence step, but by providing it, we eliminate pressure on the model to memorize it [5]. They have used such a model to generate beer reviews and it looks to be a plausible for the context based note-level music we are trying to generate.

However LSTMs are not without limitations. Boulanger, Lewandowski and Bengio, Vincent (2012) [3] talk about the difficulty in training RNNs and optimization techniques. The limitations of LSTMs include generating a single expected value, as opposed to a full conditional distribution as generated by RBMs (which allows for greater variety in generated music). They combine the time dependencies modelled by an RNN with an RBM and provide a training algorithm.

In this paper, we have used LSTM with dropout to generate music one character at a time. The music was represented as a sequence of characters in the ABC format. We have used two types of music - Irish and Swedish tunes. The context is modelled as a concatenated input. We then try to validate our results.

This section is followed by a discussion of the dataset and its representation. Following which, the architecture and mathematical details of the models used are discussed. Finally, we conclude the paper with a discussion of experiments, results and their analysis. Our quantitative and qualitative analysis shows that our model learnt the syntax of Abc notation and generated good quality music conditioned on the genre (Irish vs Swedish).

2 Dataset and Representation

2.1 Dataset

We have scraped music files from *Henrik Norbeck's Abc Tunes* [14]. It is a repository of 2180 Irish and 388 Swedish tunes in Abc notation (described in the 2.2) made available on the web as a free resource. In addition to music notes and rhythm, the dataset also contains the metadata information like title, composer, source, etc. We have used all the tunes for training in our experiments.

2.2 Data Representation

We represent the music in Abc notation [13]. Abc is a text-based music notation system designed to be comprehensible by both people and computers. Music notated in abc is written using characters - letter, digits and punctuation marks - on paper or in computer files. An abc tune itself consists of a tune header and a tune body, terminated by an empty line or the end of the file. The tune header is composed of several metadata information fields like title, composer, etc. The tune header should start with an X:(reference number) field followed by a T:(title) field and finish with a K:(key) field. The tune body, which contains the music code, follows immediately after. Figure 1 explains different sections of abc notation using a sample music in its abc notation from Henrik Norbeck's Abc Tunes [14].

3 Recurrent Neural Network Methodology

These differ from traditional feed-forward neural networks in the sense that the inputs and outputs are not independent. There exist feedback connections between nodes in a way that these connections form a directed cycle. RNNs have an internal memory which stores information about what has been calculated so far, thus giving them the capability of handling sequential data.

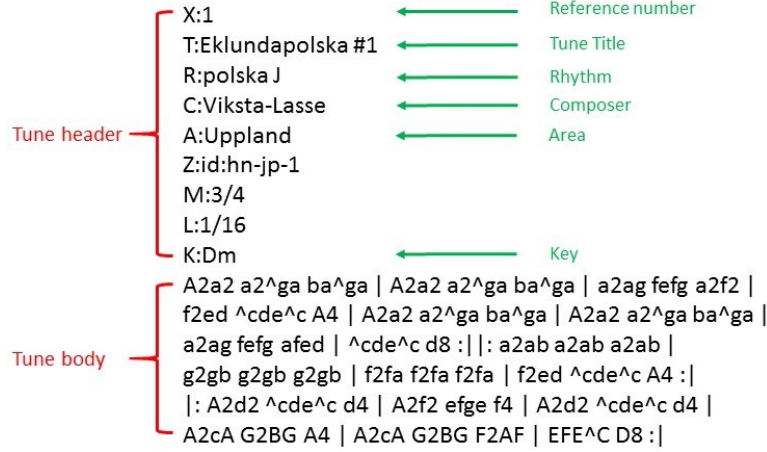


Figure 1: Abc notation of music file

The problem with RNNs is that of the vanishing/exploding gradient. When trained on an arbitrarily long sequence, if the weights are small the gradients shrink exponentially and if the weights are large they grow exponentially. This problem is overcome by using LSTMs (Long Short Term Memory), introduced by Hochreiter and Schmidhuber (1997) [6]. Each memory cell has an internal state s in which activation is preserved along a self-connected recurrent edge. Each cell also has sigmoidal input, output and forget gates. Information gets into the cell whenever its input (i) gate is on and it stays in the cell until its forget (f) gate is turned on. The information can be read from the cell by turning on its output (o) gate. Several layers of LSTMs can be stacked together (Graves- 2013) [7]. At step t , each LSTM layer $h_l^{(t)}$ receives as input the output from layer $h_{l-1}^{(t)}$ at step t and the output from layer $h_l^{(t-1)}$ at step $t-1$. As a base case, we take $h_0^{(t)} = x^{(t)}$ and $h_l^{(0)} = 0$. The following are the equations to calculate the forward pass through an LSTM:

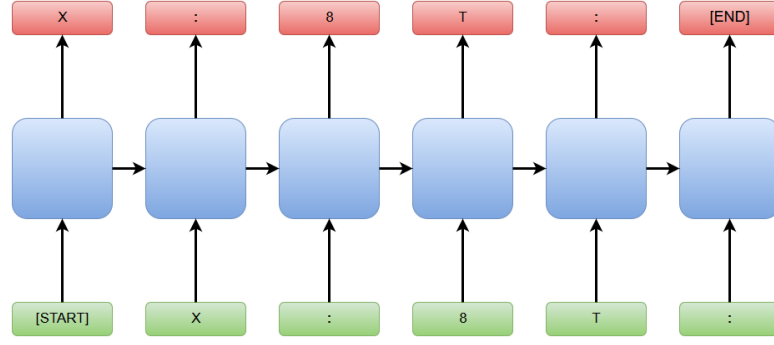
$$\begin{aligned}
 g_l^{(t)} &= \phi(W_l^{gx} h_{l-1}^{(t)} + W_l^{gh} h_l^{(t-1)} + b_l^g) \\
 i_l^{(t)} &= \sigma(W_l^{ix} h_{l-1}^{(t)} + W_l^{ih} h_l^{(t-1)} + b_l^i) \\
 f_l^{(t)} &= \sigma(W_l^{fx} h_{l-1}^{(t)} + W_l^{fh} h_l^{(t-1)} + b_l^f) \\
 o_l^{(t)} &= \sigma(W_l^{ox} h_{l-1}^{(t)} + W_l^{oh} h_l^{(t-1)} + b_l^o) \\
 s_l^{(t)} &= g_l^{(t)} \odot i_l^{(t)} + s_l^{(t-1)} \odot f_l^{(t)} \\
 h_l^{(t)} &= \phi(s_l^{(t)}) \odot o_l^{(t)}
 \end{aligned}$$

Here, σ denotes an element-wise sigmoid function, ϕ denotes element-wise \tanh , and \odot is element-wise product.

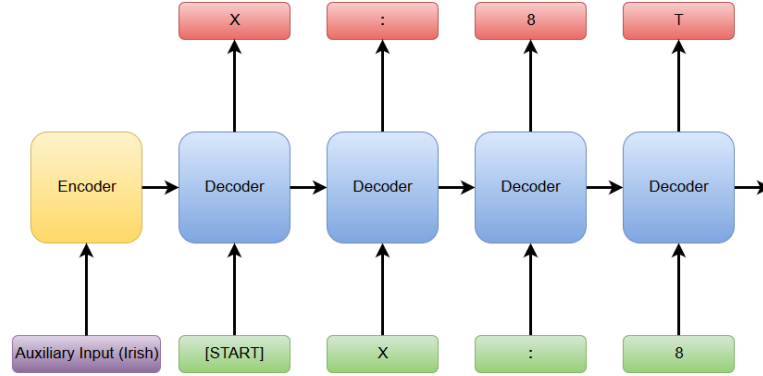
3.1 Generative Recurrent Neural Networks

Similar to the beer review paper [5], we build on the generative RNN model of Sutskever et al. (2011; 2014) [8] [9]. A generative RNN predicts the next token in a sequence (x^{t+1}) given the sequence up to that point (x^1, \dots, x^t). Thus, the output string is simply the input string shifted by one token. (Figure 2a). The output layer is fully connected with softmax activation. Cross entropy is the loss function used for training.

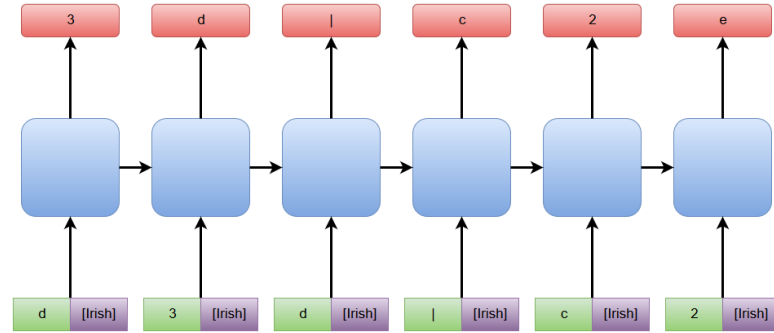
The model can then be used to generate arbitrarily long sequences given some starting token and state by passing the generated output as input in the subsequent step and repeating.



(a) Vanilla generative recurrent neural network



(b) RNN with LSTM units without dropout



(c) RNN with LSTM units with dropout

Figure 2: Plot of loss vs epochs

3.2 Concatenated Input Recurrent Neural Networks

In our experiment, we generate music based on context given in the form of an auxiliary input x_{aux} . In the model described by Sutskever et al. (2014) [9], and Karpathy and Fei-Fei (2015) [10], the auxiliary input is encoded and passed as an initial state to the decoder (Figure 2b). This auxiliary input has to be preserved across many sequence steps which is why it is successful in producing short image captions and the like, but impractical in generating long music notes.

To overcome this issue of having to remember the auxiliary input across each step, we simply concatenate x_{aux} with each character representation in the sequence (Figure 2c). At train time, x_{aux} is a feature of the training set; in the prediction phase, x_{aux} is fixed and concatenated with each character in the sequence to produce the next character.

3.3 Weight Transplantation

Models with auxiliary inputs are considerably harder to train as compared to a typical unsupervised model. To overcome this problem, the model is first trained without the auxiliary input. Once the weights have converged, they are then used in the concatenated input model. The extra set of weights from the input layer to the first hidden layer are initialized to zero. This is similar to the pre-training common within the computer vision community (Yosinski et al., 2014) [11]. Here, instead of new output weights, we train new input weights.

4 Experiment and Results

All the experiments were performed using Lasagne [15] built over Theano. Before training, we concatenated all the tunes represented in ABC format, delimiting each of the tunes with <START> and <END> tags, to produce one big file containing all the different melodies. Before generating music based on genre we wanted to check the feasibility of music generation using recurrent neural networks. We trained a vanilla RNN to see if music generation in this way is possible. It had single hidden layer with 100 nodes. We split the concatenated string into mini-batches of size 512, which are in turn split into segments of sequence length 25. Training is done in mini-batches and the weights are updated using AdaGrad with a learning rate of 0.1. We found that the model generated good quality music in Abc format after some epochs of training. This gave us the confidence to move forward with the idea of music generation using RNN. Samples of music generated using vanilla RNNs can be found at <https://goo.gl/Iw7IU2>.

Once we established the feasibility of generating music using RNNs, we trained an LSTM network with 2 hidden layers, with 512 hidden nodes in each. We experimented with various size of hidden units before finalizing the number of units to 512. Decreasing the number of hidden nodes lead to a higher error rate, meaning that the network was probably underfitting the data (Figure 3). Increasing the number of units beyond 512 lead to GPU memory issues.

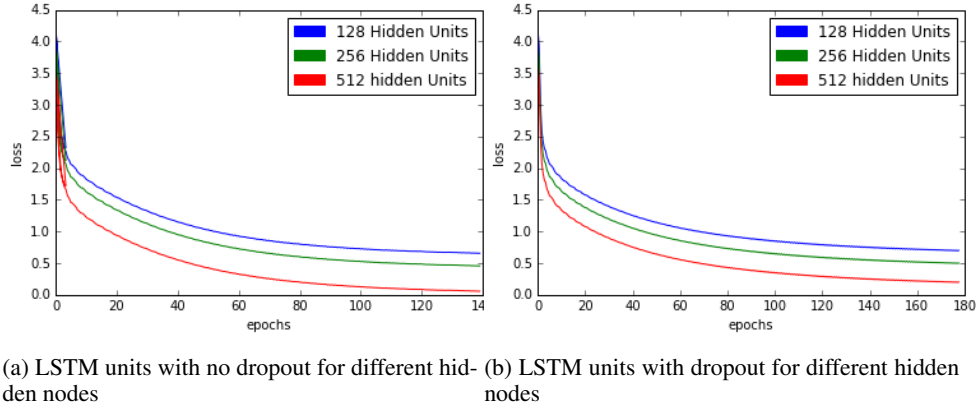


Figure 3: Loss plotted against epochs for different no. of hidden nodes

We also experimented by varying the mini-batch sizes. As the average length of a melody in the training data was found to be around 400, we observed that a small mini-batch size could not accommodate a complete melody with <START> and <END> tags. With a small mini-batch the network couldn't learn how to start or end a melody, leading to abrupt starts and finishes to the generated tunes. A mini-batch size greater than 512 led to memory errors on the GPU. So, we split the concatenated string into mini-batches of size 512 which was in turn split into segments of sequence length 64. To generate music of a particular genre, we provided a one-hot encoded auxiliary input representing Irish and Swedish melodies to the LSTM. Training was done in mini-batches and the weights were updated using AdaGrad with a learning rate of 0.1.

Next, we introduced dropout of 0.1 in LSTM. In the LSTM network with no dropout, training loss dropped quickly implying that the network is overfitting the training data (Figure 4). This led to the network generating music very similar to the training melodies, defeating the purpose of generating

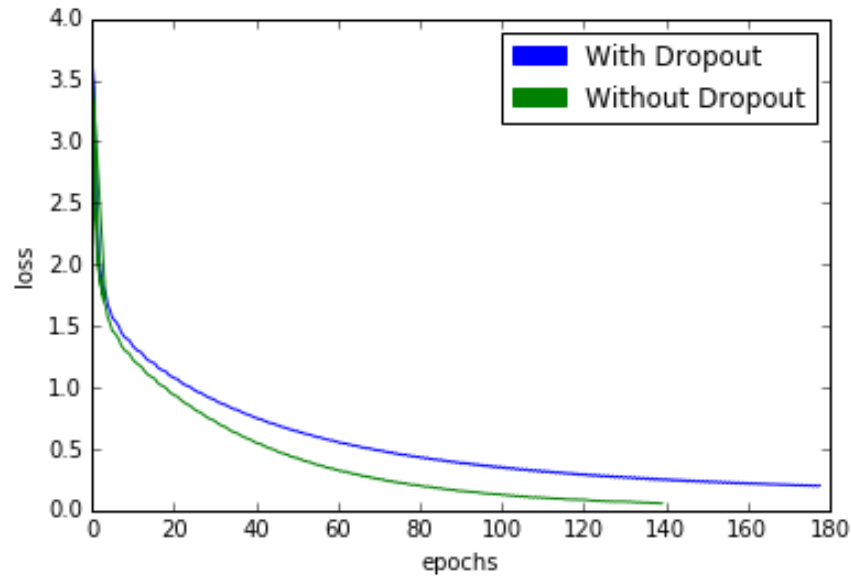


Figure 4: Plot of loss against epochs for LSTM with dropout and no dropout

new music. Introduction of dropout solved this problem and generates sounds different from the input.

Training the network with the auxiliary input leads to very slow learning. So, we first trained the plain character RNN model (i.e. without the auxiliary input), on the combined Irish and Swedish music datasets. We then transplanted the weights of the trained network into the concatenated input network for fine-tuning. The final music files generated using the model can be found at <https://goo.gl/Iw7IU2>.

X:18

T:Solle else Toll Ta er ools To ole an er oe er or on er ool en er ool en er or on er
oll an er on er or on er or on er oo so er on er or on er oor War oo sar ar on er or
ole Ar or on er or on er ool en er ols ar ool ia dor ian en er ole tar ion Ah ton

Z:id:hn-slide-

M:6/8

L:1/8

LK:1 Should be K:1

A2A A2A | A2A A2 | A2A A2A | A2A A2A | A2A A2A | A2A A2A | A2A A2A |
A2A A2A | A2A A2A | A2A A2A | A2A A2A | A2A A2A | A2A A2A | A2A A2A |
A2A A2A | A2A A2A | A2A A2A | A2A A2A | A2A A2A | A2A A2 | A2A A2 | A2
A2 A2 | A2A A2 | A2A A2 | A2A A2 | A2A A2 | A2 A2 | d2 d2 | d2 d2 | d2 d2 | e2
d2 | e2 d2 | d2 d2

Figure 5: Music generated after 10 epochs

In the generation phase, along with the auxiliary input, this network initially generated random text. After 10-20 epochs, the network started generating text which followed the ABC format with some syntax errors (Figure 5). At 50-60 epochs the network started generating syntactically correct ABC text which can be played in the Midi format. However, this output had the problem of repeating notes for long duration (Figure 6). At around 80 epochs, the network started generating ABC text

which followed the chord structure of a real melody (Figure 7). Though we could differentiate Swedish melodies from the Irish ones by listening to it, there was no concrete metric in the literature to quantitatively differentiate them. So, we compared the top 5 rhythms (*R* field in the tune header in Abc notation) generated by the network for each tune and found that these rhythms are the same in the generated tune and the training tune for both Irish and Swedish melodies. Figure 8 shows the top 5 rhythms of the training Swedish data which are consistent with the top 5 rhythms of the generated data. The same was observed for the Irish data which can be seen in Figure 8. This helped us verify that the network is actually generating music based on the genre provided as input. Also, we observed that 5 rhythms those were only present in Swedish training data and not on the Irish training data were observed in the generated Irish data in a very low fraction (Figure 10). Similarly, we found the 5 rhythms which were only present in Irish training data occurred in the generated Swedish text in very less number (Figure 11).

```
X:8
T:Pol's ar Rer Ror Roll Rar Rols an Rers on ser son sar son sols or
son son sor son son son son sor son sor son son son sor son
sors on sol sor son sons
Z:id:hn-sl-13
M:3/4
L:1/8
K:D
A2 A2 A2 | A2 A2 | A2 A2 | B2 A2 | B2 A2 | B2 A2 |
A2 A2 | B2 A2 | B2 A2 | B2 A2 | B2 A2 | B2 A2 | B2 A2 | B2 A2
| A2 A2 | B2 A2 | A2 A2 | A2 A2 | A2 A2 | B2 A2 |
A2 A2 | B2 A2 | B2 A2 | B2 A2 | B2 A2 | B2 A2 | A2 A2 | B2 A2
| A2 A2 | B2 A2 |
A2 A2 | B2 A2 | B2 A2 | B2 A2 | B2 A2 | B2 A2 | B2 A2 |
A2 A2 | B2 A2 | B2 A2 | A2 A2 | B2 A2 | A2 A2 | A2 A2 | A2 |
A2 A2 | A2 A2 |
A2 A2 | A2 A2 | A2 A2 | A2 A2 | A2 A2
```

Figure 6: Music generated after 50 epochs

```
X:2
T:Pol'sk of The
R:slide
D:Seeral Ossen
Z:id:hn-slide-6
M:6/8
L:1/8
K:Ador
A2 | G2 A2 A2 | B2 GE GE | E2 A2 A2 | B2 G2 BG |
A2 A2 c2 | B2 BA BA G2 | B2 BA/G/ GE |
D2 G2 BG | B2 AB/c/ BA | 1 G2 G2 :| 2 G2 G2 B2 ||
```

Figure 7: Music generated after 100 epochs

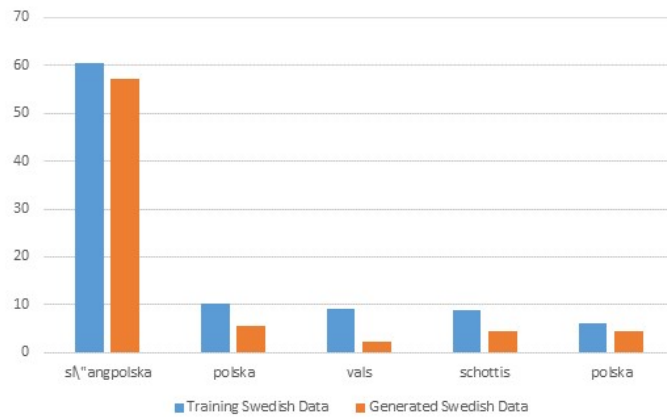


Figure 8: Top 5 Swedish rhythms on training and generated data

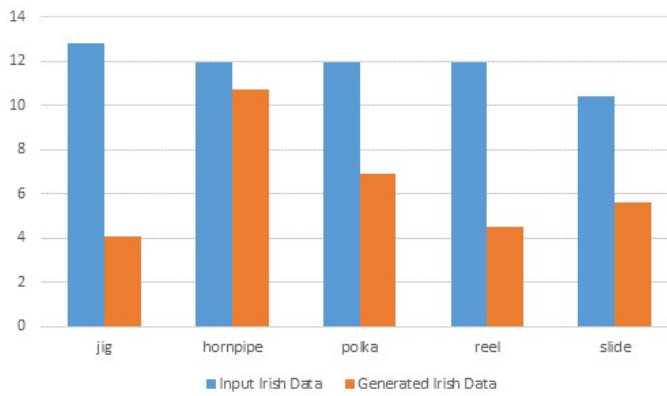


Figure 9: Top 5 Irish rhythms on training and generated data

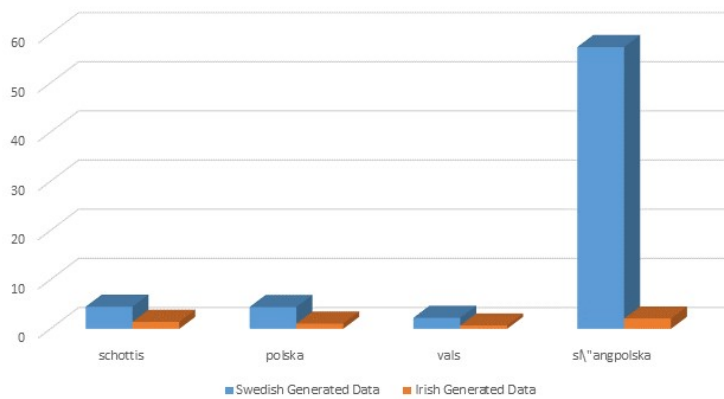


Figure 10: Rhythms not originally observed in Irish training data

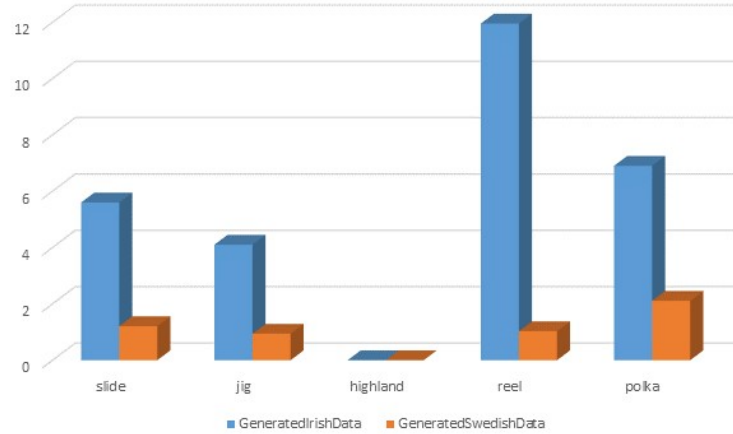


Figure 11: Rhythms not originally observed in Swedish training data

5 Conclusion

In this paper, we demonstrate the character-level recurrent neural network to generate relevant melodies conditioned on auxiliary input. This is also the first work, to our knowledge, to generate the music conditioned upon genre. Our qualitative and quantitative analysis shows that our model learnt the syntax of Abc notation and generated good quality melodies. Further, our analysis support that the model was able to generate different music conditioned on the auxiliary input. In this paper, the number of genres used were restricted to two due to the unavailability of proper dataset, however, the same model can be extended for any size of the genres.

Although, the results of the experiments seem promising but a good measure to distinguish between different genres of music is yet to be investigated. In addition to this, much work also remains to be done to quantitatively evaluate the quality of the music.

We will include the results of feature visualization at the final submission. We also plan to make a website to present a live demonstration of the project before the final submission.

6 Contributions

Apurva and Kshitiz wrote the code for LSTM generative model. Puneeth and Ritvik wrote the code for character level vanilla RNN. Chaitanya wrote the code to scrape data and to add dropout. Kshitiz and Apurva generated the plots from the results. Ritvik, Puneeth and Chaitanya wrote the most of the parts of the report and were helped by Apurva and Kshitiz.

References

- [1] Mozer, Michael C., (1994) Neural network music composition by prediction: Exploring the benefits of psychoacoustic constraints and multi-scale processing., *Connection Science*, pp. 247-280.
- [2] Eck, Douglas., (2002) A first look at music composition using lstm recurrent neural networks.
- [3] Boulanger-Lewandowski, Nicolas and Bengio, Yoshua and Vincent, Pascal., (2012) Modeling temporal dependencies in high-dimensional sequences: Application to polyphonic music generation and transcription., *arXiv preprint arXiv:1206.6392*
- [4] Nayeibi, Aran and Vitelli, Matt., (2015) GRUV: Algorithmic Music Generation using Recurrent Neural Networks.

- [5] Lipton, Zachary C and Vikram, Sharad and McAuley, Julian., (2015) Capturing Meaning in Product Reviews with Character-Level Generative Text Models., *arXiv preprint arXiv:1511.03683*
- [6] Hochreiter, Sepp, and Jrgen Schmidhuber. "Long short-term memory." *Neural computation* 9, no. 8 (1997): 1735-1780.
- [7] Graves, Alex. "Generating sequences with recurrent neural networks." *arXiv preprint arXiv:1308.0850* (2013).
- [8] Sutskever, Ilya, James Martens, and Geoffrey E. Hinton. "Generating text with recurrent neural networks." In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pp. 1017-1024. 2011.
- [9] Sutskever, Ilya, Oriol Vinyals, and Quoc V. Le. "Sequence to sequence learning with neural networks." In *Advances in neural information processing systems*, pp. 3104-3112. 2014.
- [10] Karpathy, Andrej, and Li Fei-Fei. "Deep visual-semantic alignments for generating image descriptions." In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3128-3137. 2015.
- [11] Yosinski, Jason, Jeff Clune, Yoshua Bengio, and Hod Lipson. "How transferable are features in deep neural networks?." In *Advances in Neural Information Processing Systems*, pp. 3320-3328. 2014.
- [12] Andrej Karpathy., The Unreasonable Effectiveness of Recurrent Neural Networks., <http://karpathy.github.io/2015/05/21/rnn-effectiveness/>
- [13] The abc music standard 2.1, <http://abcnotation.com/wiki/abc:standard:v2.1>.
- [14] Henrik Norbeck's Abc Tunes, <http://www.norbeck.nu/abc/>
- [15] Lasagne, <https://github.com/Lasagne/Lasagne>.