
Dilated Convolutions for Music Generation

Kaleb Morris, Hannah Leou
Department of Computer Science
Stanford University
CS 230 Spring 2018
SUNetIDs: ka.lebm, h.leou

Abstract

1 Art and science are often considered to be polar opposite disciplines. Many
2 debate that the creativity and natural fluidity behind music composition is the
3 sole product of human generation. They argue that technology, is far too cold
4 and rigid a field to produce the freeform and melodic arrangements. We would
5 like to challenge this notion and investigate the intersection of these two fields
6 through computer generated music. One may be skeptical of a computer's ability
7 to generate something so subjective as music, after all how would a computer know
8 if it was producing something that sounds pleasing? While the quality of musical
9 piece is subjective, we see that pieces from the classical genre follow fairly rigid
10 principles of music theory and composition. Unlike pop or other modern genres,
11 we see that classical melodies are highly structured and thus are perfect for training
12 a model that can intelligently generate music. Within our project, we sought to
13 learn the principles and recurrences encoded in classical music, utilizing them to
14 generate pieces through a Dilated Convolutional Neural Network.

15 1 Task Definition

16 For our project, we are using a Dilated Convolutional Neural Network to sequentially generate music.
17 Given MIDI files, we extract the primary melodic voice from the track. We then encode "active note"
18 vectors for each time step in our training sample compositions. For our model, we emulate Google
19 DeepMind's WaveNet for raw audio generation. Our model follows the general Wavenet architecture
20 of a Dilated Convolutional Neural Network, but is optimized for the task of music generation rather
21 than speech generation. Furthermore, our model is restructured such that it is designed to train on
22 and generate MIDI files. Our model produces output by means of generating an "active note" vector
23 for each time step. In our post-processing step, we transform these encoded vectors into MIDI events
24 and conjoin the successive MIDI events to generate a complete MIDI file.

25 As music is very subjective, it's difficult to quantify the results of each model to guide the development
26 of our final model. We therefore constructed a survey to determine the quality of a generated piece. In
27 this survey, we focus on evaluating the flow, rhythm and repetitiveness of the piece. We also gauged
28 listeners on how musically correct and classical the piece sounds. While the last two metrics are good
29 general guideposts, they are unreliable due to their subjectivity.

30 2 Infrastructure

31 We created a class called midiGenerator, that given the training data set outputs a generated MIDI
32 file. Our first goal in designing the class was to build an interface so that we can directly compare
33 the results of training the Dilated Convolutional Neural Network on different training sets. Our
34 second goal in designing the class was to build an interface so that we can experiment with different
35 tempos (ticks per beat) to generate the most sonically pleasing musical compositions. When running

36 the midiGenerator.py file, you can easily specify the desired model training set and output tempo
37 following the usage outlined below:

38 **Usage:** [# MIDI events to generate][path to training examples][ticks per beat]

39 **path to training examples** - The algorithm will use only the MIDI files in the specified
40 directory path to train the Diluted Convolutional Neural Network

41 **ticks per beat** - The algorithm will use the given ticks per beat to set the generated MIDI
42 file's tempo in the Meta Messages of the MIDI File

43 3 Approach

44 We begin by reading in a MIDI file, extracting the primary melodic voice, and then for each time
45 step, encode the active notes as a vector. We stack these "active note" vectors on top of one another
46 to create a matrix representation of the MIDI training file. After creating our matrix, we train our
47 model using the sequence of encoded "active note" vectors. We frame our training in the context of
48 an input x and a target y , where x is a series of "active note" vectors and y is another set of "active
49 note" vectors with overlapping time-steps.

50 3.1 Data Encoding

51 For our "active note" vector approach of data encoding, it was necessary to pair down the robust
52 MIDI files. MIDI files, especially those of classical compositions typically consist of many melodic
53 and harmonic voices represented by MIDI tracks. To single out the primary melodic voice in the
54 MIDI file, we visualized the MIDI files using GarageBand. From a simplified MIDI file, we encoded
55 the active notes at each time step with a modified "one-hot" vector. There are 128 note classes in
56 MIDI files and each of these classes is represented by a value of 1 (ON) or a value of 0 (OFF).

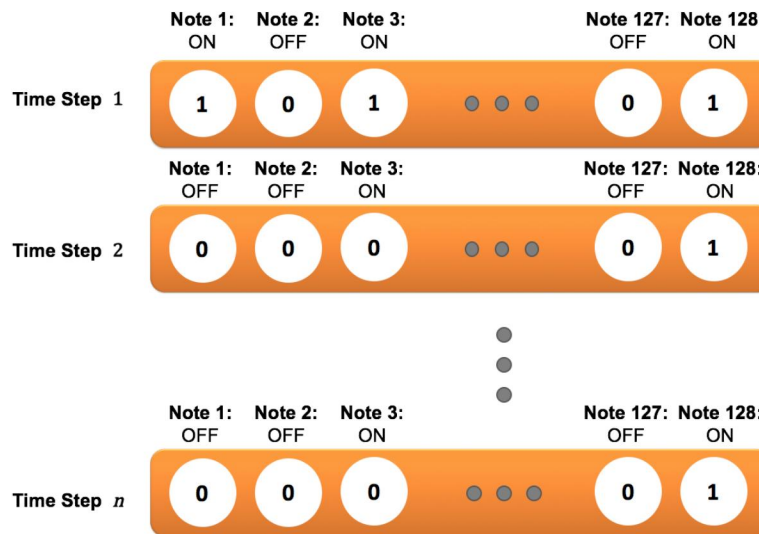


Figure 1: Visualization of the matrix composed of "active note" vectors encoded for 128 note classes

57 For a training piece, we stacked the sequential "active note vectors". The first row of the matrix is an
58 "active note" vector that corresponds to the first time step of the piece, the second row of the matrix
59 is an "active note" vector that corresponds to the second time step of the piece, and so on.

60 3.2 Model

61 We modeled our own architecture after the Google DeepMind's WaveNet, a dilated dilated convolution
62 (convolution where the filter is applied over an area larger than its length by skipping a constant
63 number of inputs)

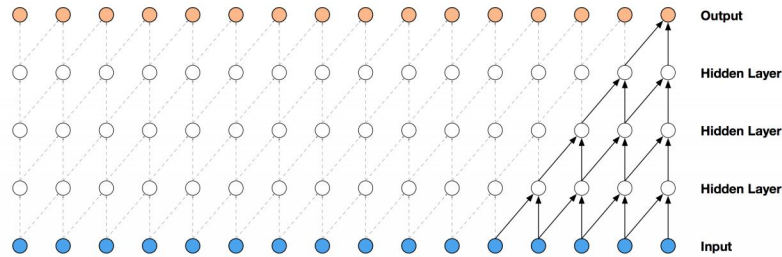


Figure 2: Figure 1: Visualization of a stack of causal convolutional layers.

64 The model we created takes as input MIDI files and generates as output MIDI files. In order to
65 generate MIDI files rather than raw audio output, we replaced the final activation function in the
66 original WaveNet with a sigmoid function in order to map the values for 128 note classes into a range
67 between $[0, 1]$. We fine tuned the threshold such that note classes with values above the threshold
68 would be switched off in the generated encoding to a value of 1 and classes with note values below
69 the threshold would be switched on in the generated encoding.

70 *Dilations:* The dilated convolutional layers incorporated to our model greatly improved the
71 receptive field. Our initial assumption that a larger receptive field would help generate better
72 music led us to add many dilated convolutional layers to our model. However, we observed
73 that additional layers minimally brought down training loss and significantly brought up
74 training time. We valued the ability to iterate on and redesign our model, so we chose to
75 incorporate a modest 5 stacked dilated convolutional layers to our model.

76 *Activation Functions:* For the activation functions within our network, we tested a number of
77 different activation functions. The activation that we settled on for the majority of our layers
78 was an ELU function, as this appeared to be giving us the best results. The final activation
79 function we utilized for our model was the sigmoid function. We utilized the sigmoid function
80 because our model was outputting "active note" vectors with note class values varying from
81 low negative values to low positive values. We chose to use the sigmoid function in our
82 final layer to recenter our data such that all class values ranged from 0 to 1 so that we could
83 transform the final vector to be representative of a note, chord, or silence.

84 *Learning Rate:* We tried a few values for our learning rate, but it appeared as though the
85 default of 0.001 in PyTorch was best suited to our task. Increasing the learning rate from
86 0.001 resulted in our loss being greatly increased, while decreasing the learning rate from
87 0.001 resulted in similar loss with much greater time to train.

88 4 Literature Review

89 Music generation is a very well-explored problem in the realm of artificial intelligence. There
90 is a diverse variety of models which have been used to model music composition. Among the
91 most popular models for music generation are Recurrent Neural Networks and Markov Chains.

92
93 Recurrent Neural Networks (RNN) are particularly well equipped for music genera-
94 tion because they are not constrained by direction like an archetypal directed Neural Network.
95 The RNN is a superior model because it has a degree of memory with each layer receiving
96 input from the previous layer and input from the previous time-step. In order to expand the
97 memory of this model, many music generation models employ Long Short-Term Memory
98 (LSTM) nodes which incorporate a memory cell which passes inputs from one layer's time
99 step down to multiple layers' preceding time-steps. Google's deep learning music project,
100 Magenta, utilizes an RNN and two LSTM's to generate a single melodic voice. Music

101 composition is a highly methodical process which contains a plethora of recurring patterns.
102 Music generation can be artificially initiated by breaking down musical pieces into common
103 subsequences and reordering the subsequences in innovative ways. The Markov Chain Model
104 is particularly apt for this task as it functions under the assumption that there is a hidden
105 structure or underlying relation between successive notes. Typically, Markov Models for
106 music generation will take into account time, pitch, and duration when crafting the state space.
107 While Markov Models are advantageous in the sense that they can be quickly trained, they
108 come with the major disadvantage of having "memorilessness" and inability to pick up on
109 dependencies between hidden layers.

110
111 For our project, we chose to model music generation using a Dilated Convolutional
112 Neural Network. We chose this model for two reasons: characteristics of classical music and
113 MIDI input/output. The classical music genre is characterized by its regimented, constrained
114 structure and contains a good deal of recurring note sequences. This underlying structure of
115 classical-style music can be exploited by Dilated Convolutional Neural Network. Moreover,
116 the encodings for MIDI files is very befitting for a Dilated Convolutional Neural Network.
117 MIDI files are encoded using messages with each message containing note, note velocity,
118 and note duration. These messages can be reinterpreted as a modified version of a "one hot"
119 vector (per time step) in which all active note classes are populated with a value of 1 and all
120 inactive note classes are populated with a value of 0.

121
122 We chose to utilize a Dilated CNN because of the incredible breakthroughs Google
123 DeepMind has made with its WaveNet, a deep generative model of raw audio waveforms.
124 WaveNet is "a fully convolutional neural network, where the convolutional layers have various
125 dilation factors that allow its receptive field to grow exponentially with depth and cover
126 thousands of timesteps." DeepMind's model, originally developed for applications in speech
127 generation is now being expanded to applications in music generation through the DeepSound
128 project.

129
130 Currently, the WaveNet model is trained on raw audio input and generates raw audio
131 output. In our project, we decided to modify and restructure the WaveNet architecture such
132 that the model is trained on MIDI files as input and generates MIDI files as output. This way,
133 the CNN will be able to capitalize on learning from focused information concentrated in MIDI
134 files rather than noisy data found in raw audio. Ultimately, the CNN we are constructing
135 will have to do significantly less work to generate equally natural music to that of a CNN
136 operating on raw audio, ideally making music generation more feasible for those without
137 significant computational power.

138 **5 Error Analysis**

139 Although the nature of the subject matter is highly subjective, we can measure error via our
140 population survey results, explained below.

141 **5.1 Population Survey**

142 The best measure of quality of a musical composition is the human ear. Therefore, we surveyed
143 individuals and asked them to score the training music compositions and generated .mid files
144 produced in five categories. We surveyed 30 students.

- 145 1. Score the flow of the song (1-5)
- 146 2. Score the rhythm/beat of the song (1-5)
- 147 3. How classical does this song sound? (1-5)
- 148 4. How repetitive is this song? (1-5)
- 149 5. How musically correct does this song sound? (1-5)

150 Using this model, we received the following scores throughout the different iterations on our
151 model trained on composer-specific data sets.

152 Throughout the process, our team managed to more closely bridge the gap between the
153 aspects of human-composed music and machine-composed music. Despite stark differences
154 in musical structure, which is primarily what allowed our survey responders to distinguish
155 human-composed music from machine-composed music, certain aspects generated by our
156 model more closely resembled those of real music.

Project Stage	Flow	Rhythm	Classical	Repetitive	Correct
Training on ONLY Mozart	3.2	2.0	3.3	1.0	1.8
Training on ONLY Haydn	3.1	2.1	3.0	1.1	2.1
Training on ONLY Beethoven	3.0	1.9	3.2	1.3	2.3
Training on ALL composers	2.6	2.2	2.9	.9	1.7

157 5.2 Flow and Rhythm

158 We found that flow and rhythm scores were fairly equal for the model when being trained
159 on exclusively Mozart, Haydn, or Beethoven. We believe that this is due to the fact that
160 composers have specific musical styling that translates to unique flow and rhythm patterns
161 our models were able to pick up on. Therefore, we observed that although all composers fall
162 under the classical genre, the model, when trained on the aggregate of Mozart, Haydn, and
163 Beethoven compositions performed significantly worse due to the discrepancies in musical
164 styling.

165 5.3 Repetitiveness and Correctness

166 We found that none of our trained models generated pieces that were characterized by a lot
167 of repetition. We believe that our model may have generated too large a receptive field with
168 the multiple dilated convolutional layers. Perhaps if we generated longer pieces a repetitive
169 quality may have been more apparent. Ultimately, we attribute this to the fact that the model
170 was trained on much longer classical pieces than we generated.

171 The correctness of our generated pieces was better among the models trained on one exclusive
172 classical artist. We believe that this is valid because each specific artist has an underlying structure
173 and style for their compositions. It is logical that generated pieces from the "Beethoven
174 ONLY" trained model would generate music resembling the Beethoven's actual compositions.
175 Further, a generated piece from the amalgamation of classical composers would struggle to
176 generalize and pinpoint specific patterns and structures among varied artists.

177 5.4 Classic Element

178 Our results seemed to somewhat capture the essence of classical music. However, we would
179 have liked to see our results more closely resemble the style of the input that our network was
180 trained on. This likely occurred due to the limitations of our model in our attempt to simplify
181 the music generative process. In the future, we would add more complexity to the modeling
182 of MIDI notes. Particularly, we could consider the velocity or duration of the notes more
183 accurately.

184 6 Conclusion

185 The music generated through our model did not sound as natural as we had initially thought it
186 would, but certain aspects of the music definitely improved upon altering our model. Notably,
187 the model's capability to generate more complex pieces consisting of notes, chords, and
188 silences came about after altering how we were processing and feeding data into our model.
189 In future iterations, we would train on a much larger dataset, and we would test with different
190 architectures, specifically deepening our network and possibly adding residual and skip
191 connections to improve our model's memory capability.

192 **References**

- 193 1. Oord, Aaron van den, et al. "Wavenet: A generative model for raw audio." arXiv preprint
194 arXiv:1609.03499 (2016).
- 195 2. Brinkkemper, Frank. "Analyzing Six Deep Learning Tools for Music Generation." The Asimov
196 Institute, 7 Oct. 2016, www.asimovinstitute.org/analyzing-deep-learning-tools-music/.
- 197 3. Johnson, Daniel. "Polyphonic Music Generation Using Tied Parallel Networks." Polyphonic Music
198 Generation Using Tied Parallel Networks, 12 Dec. 2017, www.cs.hmc.edu/~ddjohnson/tied-parallel/.
- 199 4. Merwe, A. and Shulze, W., "Music Generation with Markov Models," in IEEE MultiMedia, vol. 18,
200 no. 3, pp. 78-85, March 2011.
- 201 5. Steinsaltz, D., Wessel, D. "The Markov Melody Engine: Generating Random Two-Step Markov
202 Chains." Department of Statistics, California Berkeley, <http://www.steinsaltz.me.uk/papers/melody.pdf>.