

The BopBot

Austin Reilly, CSE, Mathew Cierpial, CSE, Max Jaffe, EE, and Vee Upatising, CSE

Abstract—The most memorable songs often feature a compelling (vocal) melody. Musicians often struggle to write melodies that are catchy and unique. The BopBot can inspire this creative process by offering melodies based on what you play into the machine. After inputting a basic chord progression, the BopBot will generate and play back a cohesive and interesting sequence of notes using machine learning.

I. INTRODUCTION

With the introduction of the radio, musicians became able to share their musical creations to a much larger audience. Instead of exclusively sharing their music through physical media and in live performances, music could be spread to anyone with a radio. In doing this, a pressure was applied to musicians to create music that had broad appeal. Catchy music catches the attention of people quickly scanning through radio stations. Although the definition of “catchy” can be hard to define, it is usually indicative of one thing: a memorable melody that gets stuck in listeners heads played on top of a usually simple chord progression.

A. Significance

Where a chord progression is generally quite simple and follows a generic pattern, melodies can be highly complex, fast, and use a higher variety of notes. Fortunately, when a chord is played, a “key” is established which limits the number of melody notes that would sound align well with the current chord. Despite this, melodies can use varying timings and change octaves, making them even more difficult to create. A full understanding of music theory is very helpful in understanding what chord corresponds to what key and timings of notes; however, it takes many years of study to fully understand music theory, and thus everyday musicians often lack this knowledge. Because of this, writing a chart-topping melody can be a very daunting challenge.

B. Context and Existing Products

Using machine learning to generate music has been done many times in the past. In fact, much research has been done on the topic. Although music generation using machine learning has been done many times, it has not been done in the same way that the BopBot will do it.

AIVA is a subscription-based artificial intelligence program that generates emotional soundtrack music based on the selected genre or an existing MIDI file. One of the core differences between AIVA and the BopBot is the input and output. AIVA is intended to be used as a tool to replace hiring someone to write compelling themes for the user’s projects. Because of this, AIVA limits the user input to a genre or a track of music that it will be “inspired” by. The BopBot, on the other

hand, takes a chord progression as an input and outputs a melody that is overlaid on top of the chord progression. Additionally, the BopBot will be functional in a live setting using the MIDI input from a keyboard whereas AIVA cannot take live input.

Google Magenta is an open-source program that uses machine learning to append similar melodies to that which the user inputs. This is different than the BopBot in similar ways to AIVA. Firstly, Google Magenta takes a melody as input where the BopBot takes a chord progression as input. Secondly, Google Magenta is not intended to be used with live input.

C. Societal Impacts

The BopBot is intended to be used by everyday musicians. Because the BopBot will be in the form factor of a music foot pedal, as seen in Table 1, it will be able to be used in conjunction with the equipment and gear that musicians already own. Although the BopBot can be used as a tool to create melodies, it is intended to be used more as a creative-assistance tool. If a musician creates an interesting chord progression, they can use the BopBot to inspire their creative process of writing a melody. The BopBot can also be used as a tool in a jam session. In this way, it can act as another musician with his/her own unique ideas based upon the chord progression being jammed to. A concern of ours regarding the BopBot is that it will be used as a crutch for newer musicians to create melodies. Although the BopBot can be used in this way, it may instill bad habits on learning musicians.

D. Requirements Analysis and Specifications

1. Can fit onto a musician’s pedal board. (Typically 170x138mm)
2. Completes melody generation in an amount of time that is conducive for live music. Less than 5 seconds is acceptable.
3. 4 musical genres to choose from for models: Rock, Blues, Classical, Pop.
4. Powered by a typical music pedal power supply (9V, 1700mA max).
5. Neural network design must have less than 23,000 neurons in each of the 3 layers (fewer if using 4 layers) in order to generate a melody within 5 seconds on a 1 GHz processor.
6. Each neural network must be smaller than 200 MB such that the microprocessor has enough data memory (SDRAM).

Table 1: System Specifications

II. DESIGN

A. Overview

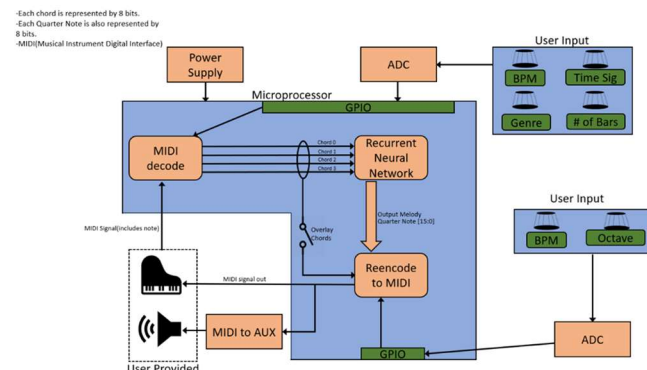


Figure 1: Block Diagram

To solve this problem, we are using a recurrent neural network. We decided to use this type of machine learning

network for several reasons which will be explained in detail in the section that goes over block 2. We are using a Raspberry Pi as our development platform and moving to a Raspberry Pi Compute Module for our final PCB. We decided to use the Raspberry Pi because it suits our needs well. We needed a system capable of running an operating system so we could easily compile the code that is necessary to generate our melodies via machine learning. We also needed a fairly powerful processor to ensure that we completed melody generation in a reasonable amount of time (5 seconds).

The first major block in the block diagram is the MIDI decode stage. Roughly speaking, this is the “recording stage” where the user input from the keyboard is turned into a format that can be used with our machine learning models to generate a melody. As you can see in the block diagram the decoded inputs (4 chords) are sent to the next major block in our diagram, the recurrent neural network. After the second block has generated a melody, the output is then sent to the third and final major block in our design, the MIDI encode stage. This can be thought of as the “playback” stage where the output from the machine learning is reformatted into MIDI so it can be output to a speaker. Besides the three main stages, there are several peripherals which allow for the user to alter different settings concerning the recording and playing back parameters.

B. Block 1 - MIDI Decode

The first stage to happen in our project is the MIDI Decode stage. This stage can be further broken down into two separate steps. The first of which is sampling the user input at the correct time steps. Human beings struggle significantly with keeping consistent time. Even trained drummers, or lifelong professional musicians struggle to keep consistent and accurate time. Below is a plot [1] of the tempo variation for Led Zeppelin’s classic song “Stairway to Heaven”

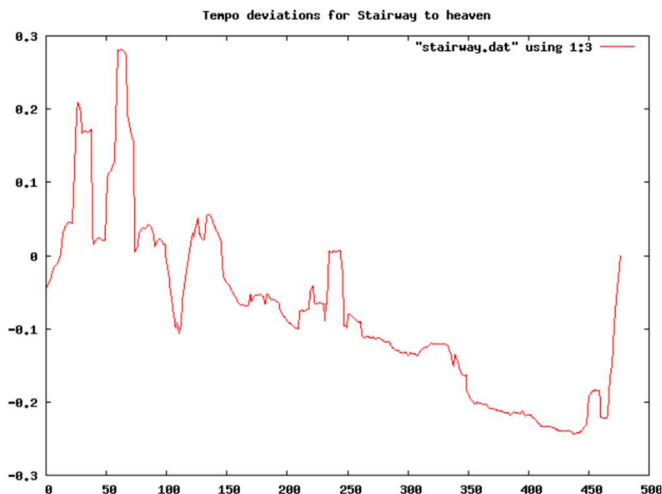


Figure 2: Tempo variations over the course of “Stairway to Heaven” by Led Zeppelin. Measures the average beat variation from the average BPM of the song. Y-axis is variation, X-axis is time.

As you can see from Figure 2, keeping consistent time is something that is virtually impossible for any given person. To fix this issue we implemented something we referred to as a

“rounding interval”. Due to our machine learning inference architecture, the input into the Recurrent Neural Network is four chords, we sampled the user input at the start of each bar line. A visual representation of the rounding interval is given in Figure 3.

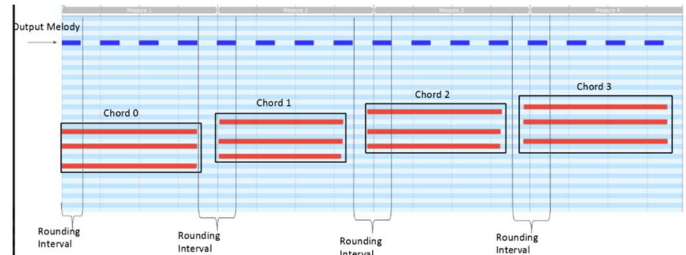


Figure 3: Visual representation of the rounding interval implemented to correct for variations in user input tempo.

To choose a proper amount of time we conducted various trials measuring each group members variation from proper timing. To do this we played a metronome at 120 BPM, in 4/4 time and simply printed out the difference in time between correct time and the time the user played a note. We found that after conducting 10 trials per team member, no variation was greater than .25s before or after the bar line. Furthermore, our rounding interval was chosen to be plus/minus .25s from the bar line.

Now that a means to account for incorrect timing from human input has been implemented, we can move on to the second step in the MIDI decode stage, which is categorizing the chords as they are being played. This step can be further broken down into two more steps which is defining what notes belong to which chord, and then how to name these chords to be used by the machine learning model. The first step in this process is defining what notes belong to the same chord. A chord is simply a sequence of three or more notes played at the same time. To determine which notes belonged to which chord we repeated a similar process to choosing a rounding interval. We simply played a series of chords and printed out the time intervals between these notes as they were played. We repeated this process several times and noticed that no notes that are part of the same chord were more than 30 MIDI ticks apart. A MIDI tick is defined to be .0010417 seconds long in 120 BPM [2], so our interval from the first note being played to each consecutive note is .031 seconds.

Now that we know which notes belong to each chord, we need to determine the relationship between these pitches. These relationships are called intervals, and the intervals of the sequence of notes determines what the chord is. For our current iteration we limited user input to be the most basic forms of simple chords, also known as triads. We limited these triads to the most basic forms of triads major and minor triads. To demonstrate this, consider Figure 4 [3].

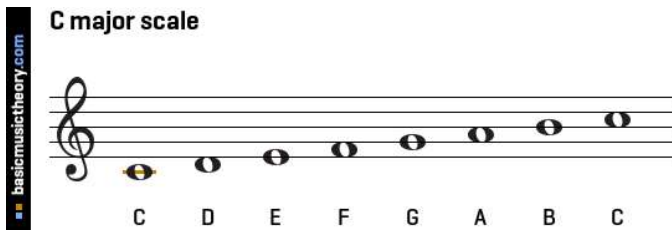


Figure 4: Notes in C Major Scale to demonstrate the structure of Major and Minor Triads

The relationship between the notes in a C Major triad are the root note (C), a major third (E), and a perfect fifth (G). Playing these notes simultaneously creates a C Major chord. You'll notice that the notes used to create the chord are simply every other note starting at the root. This is true of triads from any scale, in that you construct them using this every other note pattern. This is referred to as "stacking thirds". To differentiate the difference between a Major chord and a Minor chord lets construct a C Minor triad. The notes in a C Minor triad are the root note (C), a minor third (Eb), and a perfect fifth (G). Playing this sequence of notes simultaneously creates a C Minor chord. You'll notice that this chord does not follow a "stacking third" pattern in the C Major scale in Figure 4, that is because C Minor is not in the key of C Major. It is however in the key of C Minor, seen in Figure 5 [3].

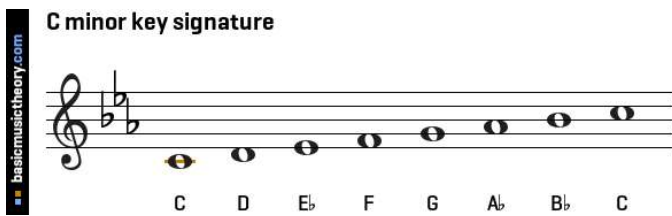


Figure 5: Notes in C Minor to demonstrate the difference between major and minor triads.

You'll also notice that the only difference between a major triad and a minor triad is the second note in the sequence (the third). Determining whether a chord is major or minor depends entirely on whether this third is a minor third or a major third, in the cases of minor chords and major chords respectively. To determine this interval, we simply use the MIDI note value, that is given in Figure 11. In music Eb and D# are enharmonic, meaning they are the same thing, what you call them depends entirely on the context. The MIDI note value for a middle C is 60, the note value for an E note is 64, and the note value for a G note is 67. These notes correspond to a C Major chord. The C is the root, the E is the third, and the G is the fifth. The distance between the root (C) and the third (E) is 4 MIDI note values, therefore a distance of four represents a Major Third. Similarly, for a C Minor chord, the root (C) is note value 60, the third (Eb) is 63, and the fifth (G) is still 67. The distance between the root and the third in a C minor chord is 3 MIDI note values, therefore a difference of 3 represents a Minor third. Once we can determine what the chord is, we created a mapping for the chords. There are 12 musical notes, each of which can

be constructed into a major and a minor chord, meaning in our limitations we have 24 possible chord inputs. We simply assigned each of these chords a value between -1 and 23 (a rest is indicated by a -1), and then filled in the appropriate value with the appropriate number to be used by the Recurrent Neural Network for melody generation.

C. Block 2 - Recurrent Neural Network

After the MIDI notes are successfully decoded, the information is then passed to the second block, the Recurrent Neural Network. A Recurrent Neural Network takes in input one at a time, preserving the sequential ordering of the data. In this case, a sequence of four chords will be passed in as input to the Recurrent Neural Network. Each value is a summary of a chord mapped to a specific integer value ranged -1 to 23 as detailed in the explanation of Block 1. Once the input has been passed in, the Neural Network performs inference, also known as forward propagation. This generates a sequence of 16 notes that correspond to the four chords that were given as input. The Neural Network generates this sequence of notes one at a time and passes the information forward. This means that every note generated has context of the notes that were previously generated before it. The raw values of the generated sequence of notes are then passed to the third block, MIDI Encode, for necessary preprocessing.

The use of the Recurrent Neural Network was a specific design choice. A Recurrent Neural Network is a very special type of Neural Network that deals with temporal data. This is data that has a sequential ordering as well non-independent instances over time. For this specific problem, every note value is directly dependent on the values of the notes around it. A Recurrent Neural Network is able to capture this information in a way that non-temporal Neural Networks are not be able to.

The second design choice corresponding to this block was the use of LSTM (Long Short Term Memory) layers inside the Recurrent Neural Network. LSTM layers are a special kind of hidden layer within a Recurrent Neural Network that allows the network to remember long-term time dependencies. The first design of this block did not include LSTM layers. After testing the network showed that it lost its context around the 10th note it generated. After which the network generated the exact same value for the rest of the 6 notes. This is known in the Machine Learning field as the Vanishing Gradient Problem. To combat this, the hidden layers of the Recurrent Neural Network were replaced with LSTM layers. This allowed the Neural Network to keep track of long-term dependencies and keep its context all the way up to the 16th note.

The current design of the Neural Network includes 3 LSTM layers with 256 nodes at each layer, which feeds into 3 Dense layers with 512 nodes at each of those layers. There are currently 3 Neural Networks each trained on a different genre of music. The current genres of music are Pop, Blues, and Classical.

Given that this is a Machine Learning problem, there are a multitude of ways in which the Neural Network's performance can be evaluated. The current metric of evaluation we are using is Mean Squared Error. Because this is a supervised learning

problem, for every input in the dataset there is also a corresponding output that is considered ideal. During the training process, the Neural Network generates an output for a given input. This output is then compared to the corresponding ideal output and the Mean Squared Error is calculated. This is the Loss Function that is defined for each of the networks.

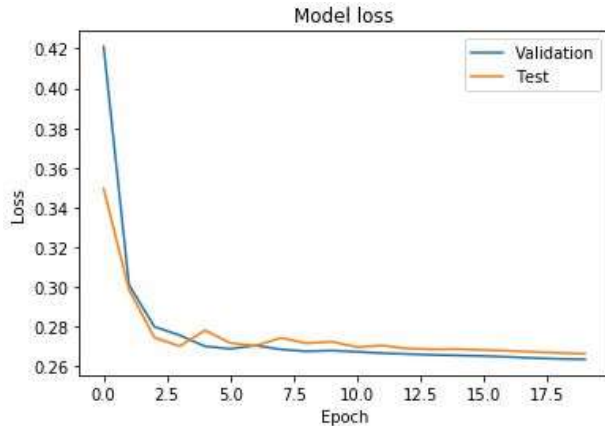


Figure 6: Training and Validation loss of Pop model

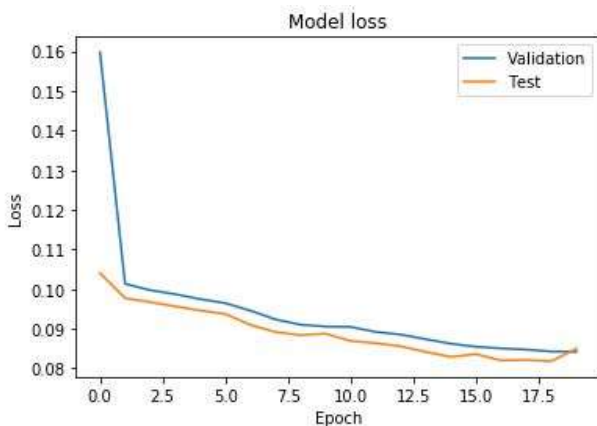


Figure 7: Training and Validation loss of Blues model

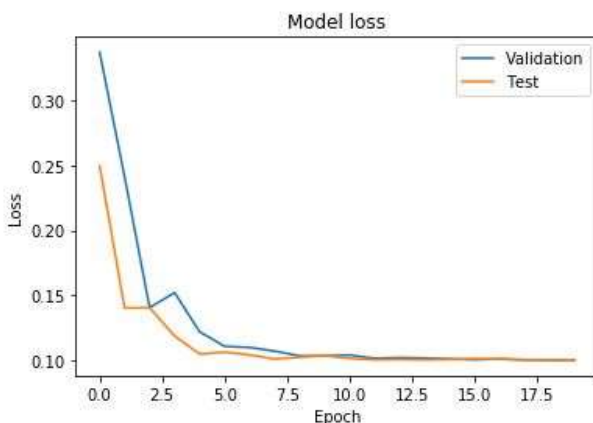


Figure 8: Training and Validation loss of Classical model

The Loss Function represents how inaccurate each Neural Network is. Over time, the Loss Function shows the rate at which each Neural Network is learning from its mistakes. It shows how each model is becoming less inaccurate after every

iteration. The three models shown were stopped after 20 training iterations.

At the start of the training process the dataset is divided into three chunks, Training, Testing, And Validation. Training takes 80% of the dataset while Testing and Validation take 10% each. The data inside the Training set is what the Neural Network learns from, and at every iteration the Neural Network uses the data inside the Testing set to calculate the Testing loss. This loss value indicates how close the Neural Network is to reaching an ideal output and is something that the Neural Network works to minimize. The data inside the Validation set is used to calculate the Validation loss. This loss metric is meant to give the most accurate representation of a Neural Network model's performance. This is because the Validation data is never shown to the Neural Network during the training or testing process. The Validation loss is silently calculated in the background at every iteration and gives us a metric of seeing how well the model performs on data that it has never seen before. This is a useful thing to measure because Neural Networks often become overfit to that data that it has been trained on. Becoming overfit means that the Neural Network performs very well on training data but poorly on new data. This generally means that the variance of this model is too high. If the Testing loss goes down while the Validation loss goes up, this is a pretty good indication that the model is becoming overfit to the data that it has been given.

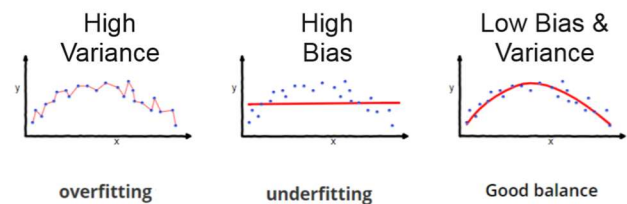


Figure 9: Relationship between overfitting, underfitting, bias and variance

As of right now, the current evaluation metric of Mean Squared Error works fine, but in the future, we have other evaluation metrics we would like to implement. One method we would like to implement is a single blind trial in which each participant listens to a compilation of melodies and tries to classify each melody as either human-composed or AI-generated. The purpose of this trial would be to get a true evaluation metric of whether a generated output is “good”. If a generated melody can fool a participant into thinking that it is human-composed, then that output would be evaluated positively and vice versa. This evaluation metric would be a lot more representative of a model's performance rather than the Mean Squared Error between a generated output and the ideal output.

A.

D. Block 3 - MIDI Encode

The last significant block in our design is the MIDI encode block. In this block we take the output from block 2, the machine learning algorithm, and turn it back into a listenable

format. This requires post processing of the data. The output of the machine learning algorithm are numbers in a continuous three-dimensional space. Thus, the numbers will almost always be floats and need to be rounded to the nearest integer value which represents the MIDI note value that the machine learning outputted. We also need to overlay the chords that the user inputted on top of the melody that is generated by the machine learning algorithm. To accomplish this, we used a function in the python library we are using to append the tracks on top of each other. We then allow the user to choose whether they want to hear just the melody that was generated, or the melody that was generated as well as the chords they played as input. Prior to making the final track that will be played the optimal octave for the melody must be selected. We decided to simply choose one octave higher than the average octave of the input chord values for now, but in the future we plan on allowing the user to select what octave to play the melody in so they can find what sounds best to them. Below is a visual representation of octaves [4].

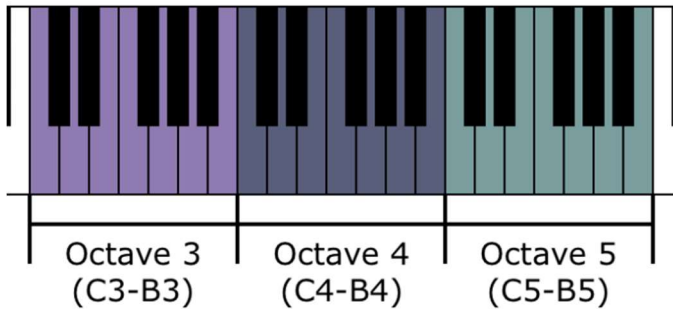


Figure 10: Visual representation of musical octaves.

Each octave has twelve notes in it that correspond to twelve sequential MIDI values. This allows us to easily scale the octave by first normalizing the data to the lowest octave, and then multiplying each note value by twelve times the number of the octave you wish to use. Below is a visual representation of the MIDI note values that demonstrates the idea of octaves well [5].

Octave Number	Notes											
	C	C#	D	D#	E	F	F#	G	G#	A	A#	B
0	0	1	2	3	4	5	6	7	8	9	10	11
1	12	13	14	15	16	17	18	19	20	21	22	23
2	24	25	26	27	28	29	30	31	32	33	34	35
3	36	37	38	39	40	41	42	43	44	45	46	47
4	48	49	50	51	52	53	54	55	56	57	58	59
5	60	61	62	63	64	65	66	67	68	69	70	71
6	72	73	74	75	76	77	78	79	80	81	82	83
7	84	85	86	87	88	89	90	91	92	93	94	95
8	96	97	98	99	100	101	102	103	104	105	106	107
9	108	109	110	111	112	113	114	115	116	117	118	119
10	120	121	122	123	124	125	126	127				

Figure 11: Visual representation of the musical octaves in MIDI

Besides allowing the user to scale the octave. In the future we will also allow them to choose the tempo at which their generated melody and input chords is played back.

III. PROJECT MANAGEMENT

Below is a list of our MDR goals.

- Live recording, melody generation, and playback ✓
- More than one available genre model to choose from ✓
- Input will be in a specific key ✓
- Playback with correct timing ✓

Figure 12: List of MDR goals set at PDR

As you can see, we completed all the goals we set for ourselves at PDR. Our current prototype generates music in real time, technically has multiple genres that you can choose to generate a melody from, allows for input of any key, and plays back in consistent timing. As for what is left to be done, we now need to create our own system on a PCB and implement our software on said new system. Below is a Gantt chart of what we need to accomplish next semester.



Figure 13: Gantt Chart for Jan-Apr 2020

As you can see by the Gantt chart, besides creating a new system and porting our existing code on to said system, we need to improve the quality of our models drastically. The amount of time we will need to dedicate to the task of data collection reflects this. Austin will be in charge of the MIDI Decode stage (recording), Vee will be in charge of the machine learning block (melody generation), Matt will be in charge of MIDI Encode (playback), and Max will provide musical expertise for the machine learning and assist both the recording and playback stages.

IV. CONCLUSION

Currently our project is a Python program on a Raspberry Pi 3 Model B. Our project is currently limited to the major/minor triad inputs outlined in Section II.B. We plan to implement a more robust chord detection algorithm that allows for more types of chords (Inversions, 7's, 9's). We also plan to implement a more robust post-processing program in our MIDI encode stage. Currently the input and output are locked at the same BPM (120) and the octave is chosen for the user. We plan to allow for the tempo to be scaled so that the output can be played faster or slower. We also plan to allow for the user to choose the octave the generated melody is in.

After consulting with our advisor, Professor Moritz, we plan to utilize a Raspberry Pi Compute Module in our larger printed circuit board. We expect this to be fairly involved, but there is extensive documentation from both Raspberry Pi and

independent sources. After this is complete, our Python program can be directly ported from the complete Raspberry Pi, to our PCB involving the Compute Module.

ACKNOWLEDGMENT

Team 21 would like to acknowledge our advisor Professor Moritz for his guidance in tackling this problem. We would also like to thank Sachin Bhat for his seemingly endless assistance during our design process.

REFERENCES

- [1] Lamere, Paul. "In Search of the Click Track." *Music Machinery*, 9 Apr. 2009, musicmachinery.com/2009/03/02/in-search-of-the-click-track/.
- [2] "MIDI Files." *Mido 1.2.9 Documentation*, mido.readthedocs.io/en/latest/MIDI_files.html.
- [3] "C Major Scale." *Basicmusictheory.com: C Major Scale*, www.basicmusictheory.com/c-major-scale.
- [4] "Octave Registers." *All About Music Theory.com*, www.allaboutmusictheory.com/piano-keyboard/octave-registers/.
- [5] "Articles." *Noterepeat.com*, www.noterepeat.com/articles/how-to/213-MIDI-basics-common-terms-explained.

APPENDIX

A. Design Alternatives

We initially planned on completing our MDR design using the De1-SoC development board. We eventually came to realize that it did not serve our purposes well and only served to overcomplicate our design process. This is mostly due to the fact that the De1-SoC has an FPGA that must be configured to gain access to a lot of the components on the board. This added a whole other element to the design process that in the end would not matter for our final design. Instead of the De1-SoC we decided to move to a raspberry pi model 3 B to complete our MDR. Not only is this board extremely easy to develop with, but it also can be used in our final PCB design as a compute module. This will make the porting of our existing software much easier, leaving only the other hardware components to be implemented (Buttons, Switches, USB interface, GPIO pins, ethernet port etc.)

B. Testing Methods

The first system specification we were able to test was number two, given in Table 1. To test if our system generated melodies in under five seconds, we simply chose 25 unique chord progression inputs and measured the time it took to generate a melody. The data for the melody generation time was then compiled into Figure 14.

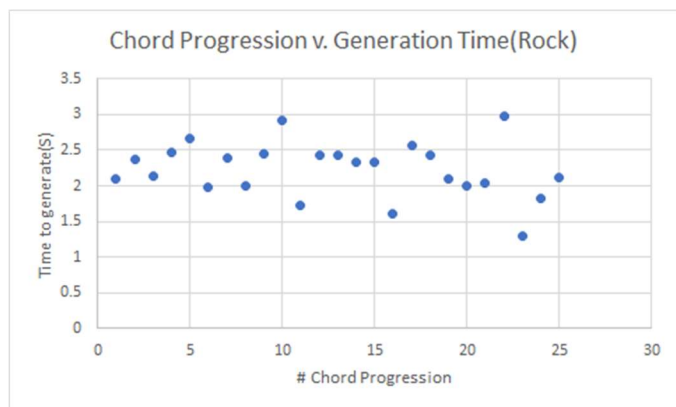


Figure 14: Chord progression inputs versus melody generation time for the rock model

A concern we emphasized at both PDR and MDR, was that our Machine Learning model would not generate melodies that sounded "good." Music is subjective, so defining a metric that determines what is good or bad is nearly impossible. Our initial idea is to simply poll people and ask them if the generated melodies sound good or not. We would have to mix in human-made melodies so that comparisons could be drawn about whether our model generated a good melody, but also a melody that sounds human-made.

C. Team Organization

Max is the project manager and he generally keeps all of us on track in terms of completing work by certain deadlines. Of course, we all have a hand in making sure things are running smoothly so we all give input from time to time and all have our individual technical work to complete.

Our team is working well together. Our individual expertises complement each other well for this project. Vee is the machine learning guru, having taken classes in the subject and conducted personal research. Austin has an interest in hardware design which has been helpful when interfacing with hardware peripherals (buttons, USB interface, GPIO pins, etc.). Matt is most interested in software which has been helpful in implementing our state machine and other code in python. Max has a good musical background which has been very important in bridging the worlds of machine learning and musical creation together.

There has not been a clear moment when a team member served as a leader. Rather we all give our feedback when pertinent and generally the whole project is directed by Max. An example of when Max helped another team member is when we came up with how our machine learning model would be constructed. He helped us decide what made sense in terms of the input and output of the model (structure of input chords and output melodies). Austin has helped Matt out extensively in creating the python code for the program, specifically in interfacing with hardware peripherals like previously mentioned. Vee is always there to discuss the machine learning model to make sure it is structured in a way that makes sense for our purposes.

We had a communication breakdown with our advisor Prof. Moritz when we were trying to use the De1-SoC for MDR. We as a team did not communicate well enough the status of how things were going with using the board which created a situation where we decided very late in the semester to switch to a Raspberry Pi. Once we established that we were not communicating effectively with our advisor we made a much more consistent effort to stay in touch with the progress of our project. This has made life a lot easier because we address problems proactively rather than waiting until there is an issue.

tune the machine learning models to be more based in music theory. This project has also helped me expand my knowledge of embedded programming when collaborating with Austin, Matt, and Vee to make sure the blocks of the block diagram interact with each other correctly. I have also gained a great deal of knowledge in team management and systems engineering. Because I am the project manager, I have been learning how to communicate with people effectively and positively so that work is accomplished in the most efficient way possible.

D. *Beyond the Classroom*

Matt: In order to write the code that implemented the MIDI encode block in our block diagram I needed to develop my python coding skills. I had only used python very briefly once before this project, so I needed to familiarize myself with python syntax and how to import libraries. Besides that, I needed to gain a basic understanding of the machine learning algorithm we utilized, just so I could give input when necessary.

Austin: Like Matt, I implemented the MIDI decode stage in Python, which I had very little experience in. Familiarizing myself with a new programming language, specifically Python, will translate very well to the professional world. Also, in the future, I would like to work in some sort of a hardware or embedded development position, so our work in the upcoming semester I feel will be very helpful in teaching me skills I may need as a professional.

Vee: I needed to do a lot of research in order to implement the Recurrent Neural Network. Before this project I had worked with other types of Neural Networks, but I had never specifically worked with Recurrent Neural Networks and temporal data. Even after I implemented the first design of the Neural Network there was still much research to be done on how to optimize the performance of the Neural Network. I consulted various research papers and the book *Machine Learning Mastery* quite a bit. I went through dozens of designs before ending up with the current design, which is still being optimized. Although the hardest skill I had to learn was how to organize and structure the data into a shape that was required for the Neural Network. I needed to find a way to mathematically represent music in three dimensions of a continuous space. It was challenging to come up with the data representation in the first place but working with the data representation after that was easy. I just needed to write a python script to extract notes information from a MIDI file and store the data in the exact three-dimensional form it needed to be for the Neural Network. I believe that the skills I developed from this data management process will greatly help me in my career in the Data Science field.

Max: Although this project is somewhat out of my interests in the field of microwaves, I have been able to apply what knowledge I have and have learned a lot. My musical knowledge has been very helpful in this music-oriented project. With this knowledge, I have been able to work with Vee to fine