

# Emotion-Controlled Symbolic Music Generation via Prompt Tuning with Structured Constraints

Anonymous ACL submission

## Abstract

Controlling emotional expressions while creating symbolic music is still considered an huge problem. Most of the current methods have generated music of limited lengths or do not provide precise control over the generated piece’s emotion. In the following text, we propose a hybrid approach that utilizes prompt tuning with a structured music generation process. To do this, we have trained prompt embeddings using a pre-trained SkyTNT model, keeping all base parameters frozen. We then use these learned prompt embeddings to generate music while using various constraints, increasing the music length from approximately 32 seconds to 148 seconds per piece. Our experiments, conducted using three different music pieces (sad, happy, and calm), have shown significant differentiation between each of them. Happy music had a note density of 12.0 notes per second, sad music had a note density of 5.21 notes per second, and calm music had a note density of 7.81 notes per second. We propose this work as a parameter-efficient prompt tuning application, which is currently gaining popularity in the natural language processing community, and discuss its benefits, as well as its limitations.

## 1 Introduction

In recent years, transformer models have become much better at producing symbolic music. They have become able to produce longer and more coherent music compared to previous RNN models. However, controlling what is produced by the models, especially the emotional qualities, is extremely difficult. When it comes to text generation, a prompt can be given to the model to generate text in a certain style or about a certain topic. However, in the case of music, the emotional qualities are based on the combined effect of tempo, harmony, melody, and rhythm.

The most common ways to perform controlled music generation are to train a new model from

scratch with conditioning or to fine-tune an entire pre-trained model. Both are costly, and the latter risks catastrophic forgetting of the musical structure that the model was pre-trained on. Additionally, most systems are limited to producing short music, less than 30 seconds long.

We use a different approach. Our method utilizes prompt tuning (Lester et al., 2021), a parameter-efficient technique developed within NLP, to learn emotion-specific representations on labeled MIDI data. Specifically, we fine-tune 12 prompt tokens for each emotion type (sad, happy, calm) on top of SkyTNT (Radford et al., 2019), a GPT-2 variant pre-trained on MIDI data. This only required 36,864 trainable parameters, which is only a small part of the millions of parameters in SkyTNT. The learned prompt is then used to inform a structured generation layer, which used constraints on scales, chord progressions, and rhythmic patterns to generate multi-track music pieces that can last anywhere from 30 seconds to 2.5 minutes.

This work is at the boundary of NLP methodology and music generation. Prompt tuning is originally a text-based language model technique. We believe the applicability of prompt tuning to the music world is an extension worth researching and using. Hence, we hope it is of interest to the NLP community, despite the output domain of music rather than text.

In our evaluation, there is clear differentiation in the emotions in the generated output: happy music has a note density of 12.0 notes per second compared to 5.21 for sad and 7.81 for calm, as well as note durations and scales that differ in ways that are consistent with what research on how people perceive emotions in music would suggest (Juslin and Sloboda, 2010).

080	<b>2 Related Work</b>			
081	<b>2.1 Symbolic Music Generation</b>			
082	Many transformer-based architectures have been			
083	successfully used in symbolic music generation.			
084	Music Transformer (Huang et al., 2018) was the			
085	first to use relative position representation into			
086	music sequences. SkyTNT, which we have used,			
087	is a GPT-2-based model (Radford et al., 2019),			
088	pre-trained on a large dataset of MIDI files and			
089	available online as a community-supported model			
090	through Hugging Face. It was also available in			
091	pre-trained form for easy use.			
092	<b>2.2 Controllable Generation</b>			
093	Previous work in controlled music generation has			
094	explored various methods to control the output,			
095	such as selecting a genre and assigning labels based			
096	on emotions given by the music. These methods			
097	usually required a change in the model architec-			
098	ture or tuning of the entire parameter space, which			
099	may possibly hurt the model’s music generation			
100	abilities.			
101	<b>2.3 Prompt Tuning</b>			
102	The technique of prompt tuning was first proposed			
103	by Lester et al. (2021) as a method to fine-tune			
104	frozen language models by training only a few			
105	learnable input vectors. This technique is highly			
106	effective on a variety of text processing tasks while			
107	using many fewer parameters than fine-tuning. We			
108	extend this approach to symbolic music by learning			
109	a few tokens specific to each emotion instead of			
110	fine-tuning the entire model.			
111	<b>2.4 Music and Emotion</b>			
112	Research on how people emotionally respond to			
113	music has shown that there are significant correla-			
114	tions between musical characteristics and the emo-			
115	tions felt by the people who listen to the music			
116	(Juslin and Sloboda, 2010). For instance, the tempo			
117	of the music has been shown to have one of the			
118	highest correlations with the emotions felt by the			
119	listeners. Faster tempos are associated with posi-			
120	tive or high-energy emotions, while slower tempos			
121	are associated with sad or relaxing emotions. Major			
122	keys are associated with happiness, and minor keys			
123	are associated with sad emotions. Our generation			
124	has followed this research.			
	<b>3 Method</b>			125
	<b>3.1 Base Model: SkyTNT</b>			126
	Our base model is SkyTNT, a GPT-2-style trans-			127
	former with 12 layers, 16 attention heads, and a			128
	hidden dimension of 1024. It was pre-trained on			129
	MIDI data and uses a vocabulary of 3,239 tokens			130
	representing musical events (note onsets, offsets,			131
	time shifts, etc.). The model is available via Hug-			132
	ging Face (skytnt/midi-model). We freeze all of			133
	its parameters during our training.			134
	<b>3.2 Prompt Tuning for Emotion Learning</b>			135
	For each emotion $e \in \{\text{sad, happy, calm}\}$ , we in-			136
	troduce a learnable prompt matrix $\mathbf{P}_e \in \mathbb{R}^{n \times d}$ ,			137
	where $n = 12$ tokens and $d = 1024$ is the hid-			138
	den dimension. During training, the prompt is			139
	prepended to the input token sequence:			140
	$\mathbf{h}_0 = [\mathbf{P}_e; \mathbf{E}(\mathbf{x})] \quad (1)$			141
	where $\mathbf{E}(\mathbf{x})$ is the token embedding of input se-			142
	quence $\mathbf{x}$ .			143
	We minimize cross-entropy loss over emotion-			144
	labeled MIDI sequences while keeping all base			145
	model parameters $\theta$ frozen:			146
	$\mathcal{L} = - \sum_{i=1}^T \log p(x_i   x_{<i}, \mathbf{P}_e; \theta) \quad (2)$			147
	The total number of trainable parameters is $3 \times$			148
	$12 \times 1024 = 36,864$ , compared to the millions of			149
	parameters in the full model.			150
	<b>3.3 Prompt-Guided Structured Generation</b>			151
	Then, after the prompts are trained, we use the			152
	statistics they have learned on to inform a process			153
	of structured generation. The idea is to combine			154
	the prompts’ learnings with the explicit constraints			155
	of music.			156
	We define base parameters for each emotion			157
	based on what is known about the ways in which			158
	listeners map musical features onto emotions:			159
	• <b>Scales:</b> Natural minor for sad, major for			160
	happy, pentatonic for calm.			161
	• <b>Chord progressions:</b> i-iv-VI-III (sad), I-IV-			162
	V-I (happy), I-ii-iii-I (calm).			163
	• <b>Base tempo:</b> 65 BPM (sad), 120 BPM			164
	(happy), 85 BPM (calm).			165
	• <b>Rhythm:</b> Longer note values for sad, shorter			166
	for happy, mixed for calm.			167

The learned prompt embeddings influence generation through tempo adjustment. We compute the mean activation  $\mu_p$  of the prompt embedding and scale the base tempo:

$$\text{tempo}_{\text{final}} = \text{tempo}_{\text{base}} \times (1 + \mu_p \times \alpha) \quad (3)$$

where  $\alpha = 0.3$ . The prompt variance also modulates rhythmic complexity: higher variance produces more varied rhythm patterns. The structured generator then produces 32 and 64 bar pieces with melody, chord, and bass tracks. Because tempo varies by emotion, piece durations range from roughly 32 seconds with happy at 120 BPM to roughly 148 seconds of sad at 64 BPM.

### 3.4 Multi-Track Generation

The generator has three tracks: melody, which uses step motion 70% of the time, which is motion to an adjacent scale degree, 30% of the time it has leaps, and sometimes it has octave transpositions; chords change every two measures, three-note clusters in an octave below the melody; bass has alternating chord roots and fifths, two octaves below the melody.

## 4 Experiments

### 4.1 Dataset

We use 58 jazz MIDI files that we manually labeled by emotion: 12 sad (e.g., “Georgia,” “Days of Wine and Roses”), 14 happy (e.g., “Caravan,” “Take the A Train”), and 32 calm (e.g., “Autumn Leaves,” “Stella by Starlight”). Labels were assigned based on the established emotional associations of each piece in the jazz literature. The dataset is admittedly small, which we discuss in the limitations.

### 4.2 Training Details

We train the emotion prompts for 60 epochs with AdamW (learning rate 0.0008, weight decay 0.01) and gradient clipping at 0.5. Input sequences are truncated to 256 tokens. On a single GPU, training takes several hours given the small dataset size.

### 4.3 Generation

For each of the three emotions, we create 64-bar pieces. Although we also created 32-bar versions during development, we only report results for the 64-bar versions, as one of our objectives is extended generation. Because of the different base tempos for each of the emotions, the resulting lengths are also different: for happy, it is about

Metric	Sad	Happy	Calm
Note Density (notes per second)	5.21	12.00	7.81
Avg Duration (seconds)	0.81	0.40	0.64
Pitch Range (semitones)	46	46	45
Unique Pitches	26	24	21
Polyphony	4.38	4.93	5.00
Pitch Entropy	4.53	4.40	4.19
Avg Velocity	63.5	64.1	64.4
Velocity Standard Deviation	11.6	11.7	11.6

Table 1: Evaluation metrics for the 64-bar generated pieces for each emotion.

64 seconds, for calm, it is about 90 seconds, and for sad, it is about 148 seconds. The learned prompt embedding is then used to adjust the base parameters before the structured generator is applied.

### 4.4 Evaluation Metrics

We try to assess the generated pieces with respect to standard symbolic music metrics, as suggested by recent surveys on evaluating music generation (Ji et al., 2023). We started out trying to use the mgeval toolkit (Dong et al., 2020), but had some compatibility problems with its Python 2 codebase, so we ended up calculating all metrics with the pretty\_midi library (Raffel and Ellis, 2014) instead.

The metrics we use are: note density (notes per second), pitch range (semitones between highest and lowest pitch), polyphony (average number of notes played at the same time), pitch entropy (a measure of how spread out the pitch decisions are), average note duration, and velocity statistics (mean and standard deviation).

## 5 Results

Table 1 shows metrics for each emotion from the 64-bar generated pieces.

### 5.1 Emotional Differentiation

The first thing to notice is the note density. For happy music, it is 12.0 notes per second, sad music has 5.21 notes per second, while calm music has a note density of 7.81 notes per second. This is what we would expect given their difference in terms of tempo and rhythm, as a faster tempo and smaller note values mean more notes are packed into a given second. However, note duration shows a very different result, as sad music has an average note duration of 0.81 seconds, while happy music has an average note duration of 0.40 seconds, with calm music sitting in between at 0.64 seconds.

Pitch range is very similar between the three emotions, at 45-46 semitones, which is what we would expect given that all three emotions sit in a similar range. However, note that the number of unique pitches does differ between each emotion, as sad music uses 26 pitches, while happy music uses 24 pitches, with calm music sitting at 21 pitches. Sad and happy music used seven-note scales while calm music used a five-note pentatonic scale. Pitch entropy follows the same pattern.

Polyphony is around 4.4 to about 5.0 for all three emotions, supporting that the chord voicings are creating a similar level of chord level for all three emotions. The differences between the emotions are in timing and melody, but not in the number of notes being played at the same time.

### 5.2 Consistency Across Lengths

We have also created 32-bar, which is the medium-sized output, versions of all pieces during development. The metrics for the 32-bar and 64-bar versions were similar, with standard deviation less than 0.05 for all metrics. This says that the generation does not get impacted by the issue of worse quality when the output gets longer, which is a problem that is seen in single-token generation. The 64-bar versions of the pieces are presented in the table from earlier.

## 6 Discussion

The hybrid approach is an effective approach in the sense that it creates musically coherent music with different characteristics per emotion. However, we would like to see what exactly differences are between the outputs.

The majority of the emotional differences in the output is a result of the generation parameters, for example, the scale that is used, the tempo it is set to, and the rhythms that are included. These are set by hand according to what we know about how humans perceive emotions in music; these parameters are not learned. The P-tuning aspect does play a role in this as well, as it adjusts tempo according to what it has learned about prompt statistics. In practice, however, this is a small change, usually in the range of 0-1 BPM. With a larger training set, the prompts could have a greater effect. However, with a set of 58 examples, there is only so much signal to be learned.

Of course, the system has a few things that benefit it. First, the parameter efficiency is real, where

we are only using 36K parameters compared to full fine-tuning, which uses millions. Second, the use of a structured layer means that we can generate music that is a few minutes long without it completely falling apart, which is hard to create when one is only trying to generate the next word. Finally, it's all interpretable, meaning a researcher and musician can look at the parameters and see what scale and tempo are being used, and change it based on what they need to change it to.

We view this as a decent starting point. The prompt tuning adds a learned component that is not just based on the rules, while the use of a structured layer provides the musical unity that neural generation has difficulty with on a dataset of this size. However, it remains a question whether this type of hybrid approach will be useful as datasets grow in size and generation becomes easier.

## 7 Conclusion

We introduced a hybrid method for emotion-controlled symbolic music generation with prompt tuning and a generation based off of a rule that was set. Emotion prompts are learned from a frozen SkyTNT model and used to control parameters in a rule-based music generation component. The model generates music pieces with 30 seconds to 2.5 minutes of length on multiple music tracks in MIDI format. Happy music has 2.3 times the density of sad music, while other metrics such as note lengths, pitch, and scale vary consistently amongst emotions.

The model can create music pieces of any length without a significant drop in quality. Future work includes testing the model on a larger set of emotion-labeled music, exploring the relationship between learned prompts and music generation, and most importantly, testing the model on human listeners to ensure that the metric differences correspond to actual differences in music perception.

### Limitations

We used only 58 labeled MIDI files, which are from the jazz genre. This is a small data set and may affect the prompt tuning performance. If there is a more diverse data set, prompt tuning may perform better and may result in a better output.

Furthermore, we use only three basic emotions: sad, happy, and calm. Emotions are more complex and diverse, but we used the three most basic ones. Thus, using more diverse emotions could result in

348 different outputs.

349 The influence of the P-tuning component is also  
350 small. The tempo can only change by 0-1 BPM.  
351 The majority of the influence is from the hand-  
352 designed structured parameters. It is worthy to  
353 restate this one more time to stress the importance  
354 of this.

355 On top of that, we did not use any human lis-  
356 tening study. Our evaluation is objective and thus,  
357 we use only the objective metrics from the MIDI  
358 output. Note density and note duration relate to  
359 tempo-related emotions, but we did not test whether  
360 humans can really hear the emotions. This is the  
361 most important thing we did not do, and it is the  
362 most important thing we missed.

363 Finally, the generation method that we used is  
364 less flexible but more coherent. A fully end-to-  
365 end learning method may discover more ways of  
366 expressing emotions than our predefined rules can.

## 367 References

368 Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszko-  
369 reit, Noam Shazeer, Ian Simon, Curtis Hawthorne,  
370 Andrew M. Dai, Matthew D. Hoffman, Monica  
371 Dinulescu, and Douglas Eck. Music Transformer:  
372 Generating Music with Long-Term Structure. *arXiv*  
373 *preprint arXiv:1809.04281*, 2018.

374 Brian Lester, Rami Al-Rfou, and Noah Constant. The  
375 power of scale for parameter-efficient prompt tuning.  
376 In *Proceedings of the 2021 Conference on Empirical*  
377 *Methods in Natural Language Processing (EMNLP)*,  
378 pages 3045-3059, 2021.

379 Patrik N. Juslin and John A. Sloboda. *Handbook of*  
380 *Music and Emotion: Theory, Research, Applications*.  
381 Oxford University Press, 2010.

382 Shulei Ji, Jing Luo, and Xinyu Yang. A comprehen-  
383 sive survey on deep music generation: Multi-level  
384 representations, algorithms, evaluations, and future  
385 directions. *arXiv preprint arXiv:2506.05104*, 2023.

386 Hao-Wen Dong, Ke Chen, Julian McAuley, and Tay-  
387 lor Berg-Kirkpatrick. MusPy: A toolkit for symbolic  
388 music generation. In *Proceedings of the 21st Interna-*  
389 *tional Society for Music Information Retrieval Con-*  
390 *ference (ISMIR)*, pages 142-149, 2020.

391 Colin Raffel and Daniel P. W. Ellis. Intuitive analy-  
392 sis, creation and manipulation of MIDI data with  
393 `pretty_midi`. In *Proceedings of the 15th International*  
394 *Society for Music Information Retrieval Conference*  
395 *(ISMIR)*, Late Breaking and Demo Papers, 2014.

396 Alec Radford, Jeff Wu, Rewon Child, David Luan,  
397 Dario Amodei, and Ilya Sutskever. Language models  
398 are unsupervised multitask learners. *OpenAI Blog*,  
399 1(8):9, 2019.