Evaluating Algorithms that Learn how to Compose Music from Scratch

New long— and short—term metrics for evaluating model-generated compositions

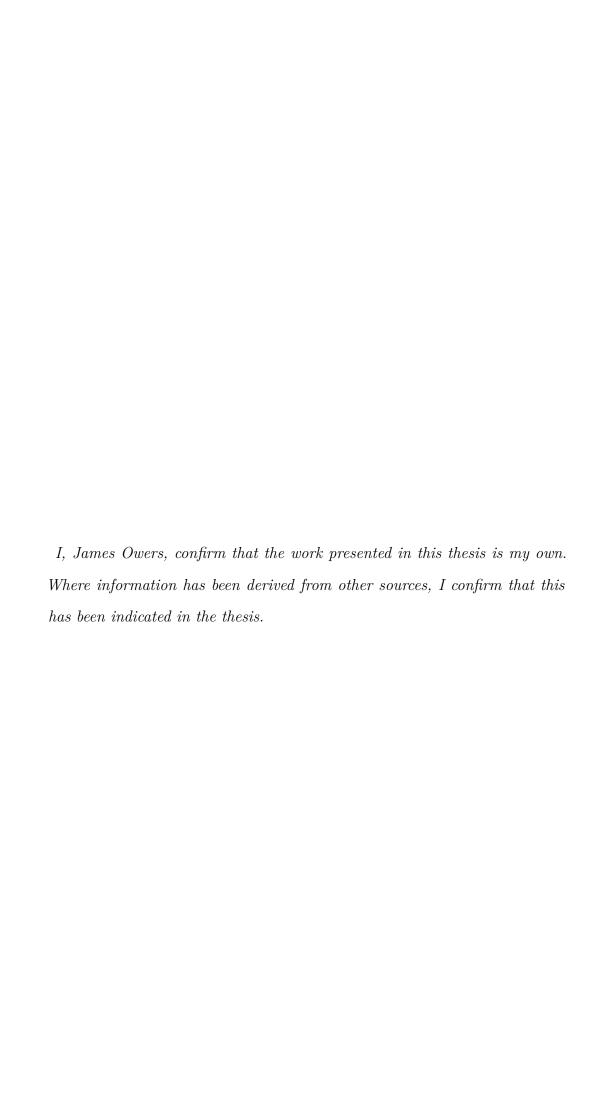
James Owers

A thesis presented for the degree of Doctor of Philosophy

> Supervised by: Amos Storkey

Mark Steedman

University of Edinburgh, UK January 2021



Abstract

Evaluating whether creative content generated by a computer is 'good,' be it music, images, or text, is unsolved and not even well defined. We identify a property of music which is not modelled well, and propose new evaluation metrics for music generation which can be used to distinguish between real and generated data, and thus be useful for automatic quantitative analysis of generation quality.

We focus on symbolic music because ...TODO... This is interesting because ...TODO... and it has implications for ...TODO...

Finally, we make recommendations for how to make progress with respect to music generation and related tasks.

Acknowledgements

Interdum et malesuada fames ac ante ipsum primis in faucibus. Aliquam congue fermentum ante, semper porta nisl consectetur ut. Duis ornare sit amet dui ac faucibus. Phasellus ullamcorper leo vitae arcu ultricies cursus. Duis tristique lacus eget metus bibendum, at dapibus ante malesuada. In dictum nulla nec porta varius. Fusce et elit eget sapien fringilla maximus in sit amet dui.

Mauris eget blandit nisi, faucibus imperdiet odio. Suspendisse blandit dolor sed tellus venenatis, venenatis fringilla turpis pretium. Donec pharetra arcu vitae euismod tincidunt. Morbi ut turpis volutpat, ultrices felis non, finibus justo. Proin convallis accumsan sem ac vulputate. Sed rhoncus ipsum eu urna placerat, sed rhoncus erat facilisis. Praesent vitae vestibulum dui. Proin interdum tellus ac velit varius, sed finibus turpis placerat.

Table of Contents

\mathbf{A}	Abstract Acknowledgements		
\mathbf{A}			
\mathbf{A}	bbre	viations	
1	Inti	roduction	1
	1.1	Target problems addressed in this thesis	1
	1.2	Historical background	1
		1.2.1 First instance of music generation	1
		1.2.2 Mozart using mechanical aids for idea generation	2
		1.2.3 Ada lovelace noting computers could generate music	2
	1.3	Modern interest and achievements	3
	1.4	What are algorithms that learn	4
	1.5	What is composing music	4
	1.6	What does it mean to compose from scratch	4
	1.7	Motivation for this work	4
	1.8	Scope of this work	4
	1.9	List of contributions in this thesis	5
2	Lite	erature Review	6

2.1	Challe	nges addı	ressed in the literature and how they are eval-	
	uated			7
2.2	A sum	mary of e	evaluation methods for creative models	7
	2.2.1	The nee	d for automated metrics	7
	2.2.2	Differen	ces between evaluating audio and symbolic	
		outputs		8
	2.2.3	The imp	possible task of satisfying all evaluation re-	
		quireme	nts with a single metric	8
	2.2.4	Evaluati	ion metrics and representations used for ex-	
		tracting	musical structures	8
2.3	Models for composing music			
	2.3.1	Models	which learn from scratch	9
	2.3.2	Models	which do not learn from scratch	9
		2.3.2.1	Heuristic models which primarily copy and	
			edit music from a database	10
		2.3.2.2	Models which incorporate expert knowledge	
			into their design	10
		2.3.2.3	Models which only work in conjunction with	
			a human composer	10
		2.3.2.4	Proprietary models: models for which ade-	
			quate details of their design are not publicly	
			available	10
2.4	Metho	ds for rep	presenting music on a computer	10
	2.4.1	Informa	tion that must be captured about a musical	
	performance			11

		2.4.2	The different representations of symbolic musical in-	
			formation	11
			2.4.2.1 Summary of differences	11
		2.4.3	Availability of data for each representation	11
		2.4.4	Availability of software for different representations .	11
		2.4.5	Evidence from the literature regarding modelling per-	
			formance differences	12
	2.5	Ethica	al considerations when designing automated methods for	
		compo	osing music	12
3	Nev	${ m w \ metr}$	ics for Evaluating Musical Generations	13
	3.1	The M	MIDI degradation toolkit	13
	3.2	A phr	ase-level metric for short-term structure	13
	3.3	A piec	ce-level metric for long-term structure	13
4	Eva	luating	g State-of-the-Art Music-generating Models	14
	4.1	Comp	arative analysis using new and existing metrics	14
	4.2	Streng	gths and shortcomings of existing models	15
	4.3	Avenu	es for improvement	15
5	A N	lew M	odel	16
6	Cor	clusio	n	18
7	Ref	erence	s	19

List of Figures

List of Tables

Abbreviations

API Application Programming Interface

 $\mathbf{J}\mathbf{SON}$ \mathbf{J} ava \mathbf{S} cript \mathbf{O} bject \mathbf{N} otation

MDTK Midi Degradation Toolkit

SOTA State of the Art

Introduction

- 1.1 Target problems addressed in this thesis
- 1.2 Historical background

...TODO... Give historical background

1.2.1 First instance of music generation

...TODO... From Section 1.2 (Briot et al. 2019)

The first music generated by computer appeared in 1957. It was a 17 seconds long melody named "The Silver Scale" by its author Newman Guttman and was generated by a software for sound synthesis named Music I, developed by Mathews at Bell Laboratories

1.2.2 Mozart using mechanical aids for idea generation

...TODO... From footnote 7 in Section 1.2 (Briot et al. 2019)

One of the first documented case of stochastic music, long before computers, is the Musikalisches Wurfelspiel (Dice Music) by Wolfgang Amadeus Mozart. It was designed for using dice to generate music by concatenating randomly selected predefined music segments composed in a given style (Austrian waltz in a given key).

1.2.3 Ada lovelace noting computers could generate music

...TODO... From (Hollings et al. 2018) Ada Lovelace, "Sketch of the Analytical Engine invented by Charles Babbage, Esq., by L. F. Menabrea," Scientific Memoirs, vol. 3, ed. Richard Taylor, 1843, pp. 666-731 (this quote on p 694)

"Note G" is the culmination of Lovelace's paper, following many pages of detailed explanation of the operation of the Engine and the cards, and of the notation of the tables. The paper shows Lovelace's obsessive attention to mathematical details - it also shows her imagination in thinking about the bigger picture.

Lovelace overseed a fundamental principle of the machine, that

the operations, defined by the cards, are separate from the data and the results. She observed that the machine might act upon things other than numbers, if those things satisfied mathematical rules.

Supposing that the fundamental relations of pitched sounds in the science of harmony and of musical composition were susceptible of such expression and adaptations, the engine might compose elaborate and scientific pieces of music of any degree of complexity or extent.

Lovelace also has the Lovelace Test of Creativity attributed to her-see (Ariza 2009).

1.3 Modern interest and achievements

...TODO...

- Imogen Heap: How AI is helping to push music creativity
- This is the AI Song Contest In the AI Song Contest teams of musicians, artists, scientists and developers take on the challenge of creating a new Eurovision-like hit with the help of artificial intelligence.
- Swooshes, Seaboards, Synths and Spawn
- David Rosen and Scott Miles on the Neuroscience of Music and Creativity
- AI Music Generation Challenge 2020 (Sturm 2020)

1.4 What are algorithms that learn

... TODO... define/introduce machine learning

1.5 What is composing music

...TODO...

1.6 What does it mean to compose from scratch

...TODO... what is the minimum information we supply as a starting point? What feedback do we give?

1.7 Motivation for this work

...TODO...

- why are we focussing on metrics and not human evaluation
- why do we care about 'from scratch'

1.8 Scope of this work

...TODO... Symbolic music only

1.9 List of contributions in this thesis

...TODO...

☐ We review the challenges addressed and models presented amounting to the current state-of-the-art with respect to algorithmic music composition

...TODO...

Literature Review

 \Box Add content and TODOs from Quip:

 $\hfill\Box$ https://quip.com/v
VjVADMamfDm/2-Literature-Review $\hfill\Box$ https://quip.com/v
0MOAQGvMH3O/Literature-Review-Org-

Notes

 \square move data and notes tables in ./tables (use YAML, allows large text blocks for notes easy to convert to table with python)

□ outline the different problems people currently try to / can solve, and how these problems relate to 'being able to compose'

 \square Review available metrics

 $\hfill\square$ Motivate the need for automated evaluation metrics

 \square Why has there been more work on audio than symbolic?

☐ Motivate the need for better evaluation for both short and long term by highlighting shortcomings for each method reviewed

$\hfill\square$ Describe state-of-the-art generative models for music composition
\Box Identify a gap with respect to modelling long term dependencies by
outlining claims and proof of them thus far - this is a specific thing we
are going show is poorly evaluated
$\hfill\Box$ Inform the reader about the multitude of different ways we can repre-
sent music and their relative strengths and weaknesses
\Box Address ethical short comings with respect to learning to compose
$\hfill\Box$ Update bibtex references to non-arxiv reference if available
2.1 Challenges addressed in the literature and how they
are evaluated
2.2 A summary of evaluation methods for creative models
 TODO how to evaluate generative models with a focus on music - how do people evaluate their success
2.2.1 The need for automated metrics
\Box Humans are expensive - show some efforts
$\hfill\square$ Humans do not agree - show some research proving this
$\hfill\square$ Humans are susceptible to change their opinion depending on context
- show some research proving this

Consistency is key when tracking performance over the long term
Attempts to unify metrics and human opinion - give WMT as an ex-
ample (Haddow 2020)

2.2.2 Differences between evaluating audio and symbolic outputs

- 2.2.3 The impossible task of satisfying all evaluation requirements with a single metric
 - □ Should tie a metric with performance for an intended task (Theis et al. 2016)
- 2.2.4 Evaluation metrics and representations used for extracting musical structures
- 2.3 Models for composing music

...TODO... State caveats about our distinctions:

- 1. Learning is essentially copying
- 2. By specifying the method of learning, we are incorporating expert knowledge
- 3. All models must work with a human composer to some extent the programmer must choose a representation for the music and is therefore a composer in some senses

Make and curate comparison table of models:
\square keep as csv
\square at min, we can use pand as to read and auto convert for insertion
here
$\hfill \Box$ is there a way to ${\tt Qreference}$ the file here and have pandoc insert?
<- do not spend time on this, cursory google!
Music Transformer (Huang et al. 2019)
MuseNet (Payne 2019)
Extend table from Chapter 7 and information from Chapter 6 in review
paper (Briot et al. 2019)
Go through https://paperswithcode.com/task/music-generation

2.3.1 Models which learn from scratch

...TODO... These are the models which our research pertains to

2.3.2 Models which do not learn from scratch

...TODO... These models are stated to highlight why they are different, have an unfair advantage in certain contexts, or explain why they are out of scope with respect to the investigation of this thesis.

- 2.3.2.1 Heuristic models which primarily copy and edit music from a database
- 2.3.2.2 Models which incorporate expert knowledge into their design

...TODO... e.g. with respect to structural hierarchy

- 2.3.2.3 Models which only work in conjunction with a human composer
- 2.3.2.4 Proprietary models: models for which adequate details of their design are not publicly available
- 2.4 Methods for representing music on a computer

...TODO... how to represent music data - (in relation to 'from scratch,' what is the minimal information supplied to the models, and is there evidence of what difference it makes (either by experiment or just by reasoning?)

- □ Note our desires with respect to our modelling challenges: we want the input to be *minimal and flexible* the model should learn as much as possible as if it were a human listener
 - \square Ideally we would work directly on sound, but this involves an

additional layer of representation.

2.4.1	Information that must be captured about a
	MUSICAL PERFORMANCE
2.4.2	THE DIFFERENT REPRESENTATIONS OF SYMBOLIC MU-
	SICAL INFORMATION
2.4.2.1	Summary of differences
	ghlight where information captured by each representation is both ferent and more/less amenable to being learned
2.4.3	Availability of data for each representation
□ Qı	nantity,
□ Qu	iality,
□ Le	gal issues
2.4.4	Availability of software for different repre-
	SENTATIONS
□ De	escribe MusPy (Dong et al. 2020) for conversion between data for-
ma	ats
□ De	escribe Music21 (Cuthbert & Ariza 2010) for conversion between

data formats

2.4.5	EVIDENCE FROM THE LITERATURE REGARDING MOD-
	ELLING PERFORMANCE DIFFERENCES
	ind any reviews (or lack thereof) of model performance differences ith respect to:
	□ evaluation metrics□ speed
2.5	Ethical considerations when designing automated

methods for composing music

New metrics for Evaluating Musical Generations

3.1 The MIDI degradation toolkit

(McLeod et al. 2020)

- 3.2 A phrase-level metric for short-term structure
- 3.3 A piece-level metric for long-term structure

Evaluating State-of-the-Art Music-generating Models

- ☐ If the main contribution is evaluating the long term structure, then ensure this is emphasized either in the title of this chapter or in the first lines
- 4.1 Comparative analysis using new and existing metrics
 - Use phrase and piece level metrics to evaluate state-of-the-art models
 - Compare and contrast, outlining the issues identified (e.g. meandering, no high-level structure)

- 4.2 Strengths and shortcomings of existing models
- 4.3 Avenues for improvement

A New Model

Potential ideas:

- An improved generative model for music
 - Training like BERT? http://jalammar.github.io/illustrated-bert/
 - Using mdtk for data augmentation in training (negative examples?), making them more robust
 - Alternative training objectives:
 - * crossentropy slow and not musically informed
 - * can we use something akin to word error rate (this has been done for text)
- Alternative ways to encode music: encoding chords and phrases in a low-rank continuous space
 - Have done some work on this with convnets and generating continuations
 - * low rank was enforced by cross-product ing two vecs

 Could investigate effect of different representations for music on performance

Conclusion

References

...TODO...

- check over using https://www.cl.cam.ac.uk/~ga384/bibfix.html
 also check with https://github.com/yuchenlin/rebiber
 Check all title casing correct (use curly braces around letters which should remain as they are). All titles should be in Title Case.
- Ariza, C., 2009. The Interrogator as Critic: The Turing Test and the Evaluation of Generative Music Systems. *Computer Music Journal*, 33(2), pp.48–70. Available at: https://www.jstor.org/stable/40301027 [Accessed January 22, 2021].
- Briot, J.-P., Hadjeres, G. & Pachet, F.-D., 2019. Deep Learning Techniques for Music Generation A Survey. arXiv:1709.01620 [cs]. Available at: http://arxiv.org/abs/1709.01620 [Accessed January 22, 2021].
- Cuthbert, M.S. & Ariza, C., 2010. music21: A Toolkit for Computer-Aided Musicology and Symbolic Music Data. In Proceedings of the 11th International Society for Music Information Retrieval Conference. Available at: https://dspace.mit.edu/handle/1721. 1/84963 [Accessed January 23, 2021].
- Dong, H.-W. et al., 2020. MusPy: A Toolkit for Symbolic Music Generation. In *Proceedings of the 21st International Society for Music Information Retrieval Conference*. Montreal, Canada. Available at: https://program.ismir2020.net/static/final_papers/187.pdf [Accessed January 26, 2021].
- Haddow, B., 2020. EMNLP 2020 Fifth Conference on Machine Translation (WMT20). 2020 Fifth Conference on Machine Translation (WMT20). Available at: http://www.statmt.org/wmt20/ [Accessed January 23, 2021].
- Hollings, C., Martin, U. & Rice, A.C., 2018. Ada Lovelace: The Making of a Computer Scientist, Oxford: Bodleian Library.
- Huang, C.-Z.A. et al., 2019. Music Transformer: Generating Music with Long-Term Structure. In *Proceedings of the 20th International Society for Music Information Retrieval Conference*. Available at: https://openreview.net/pdf?id=rJe4ShAcF7 [Accessed January 26, 2021].

- McLeod, A., Owers, J. & Yoshii, K., 2020. The MIDI Degradation Toolkit: Symbolic Music Augmentation and Correction. In *Proceedings of the 21st International Society for Music Information Retrieval Conference*. Montreal, Canada. Available at: https://program.ismir2020.net/static/final_papers/182.pdf [Accessed January 26, 2021].
- Payne, C., 2019. MuseNet. *OpenAI*. Available at: https://openai.com/blog/musenet/[Accessed January 23, 2021].
- Sturm, B., 2020. AI Music Generation Challenge 2020. The 2020 Joint Conference on AI Music Creativity. Available at: https://boblsturm.github.io/aimusic2020/ [Accessed January 22, 2021].
- Theis, L., Oord, A. van den & Bethge, M., 2016. A Note on the Evaluation of Generative Models. arXiv:1511.01844 [cs, stat]. Available at: http://arxiv.org/abs/1511.01844 [Accessed January 23, 2021].