

JAN MYCKA

# AUTOMATYCZNE GENEROWANIE MUZYKI EKSPRESYJNEJ

State-of-the-art

# Istotność ekspresji w muzyce

Oore S., Simon I., Dieleman S., Eck D., Simonyan K.: *This Time with Feeling: Learning Expressive Musical Performance*

Kompozycja



Wykonanie



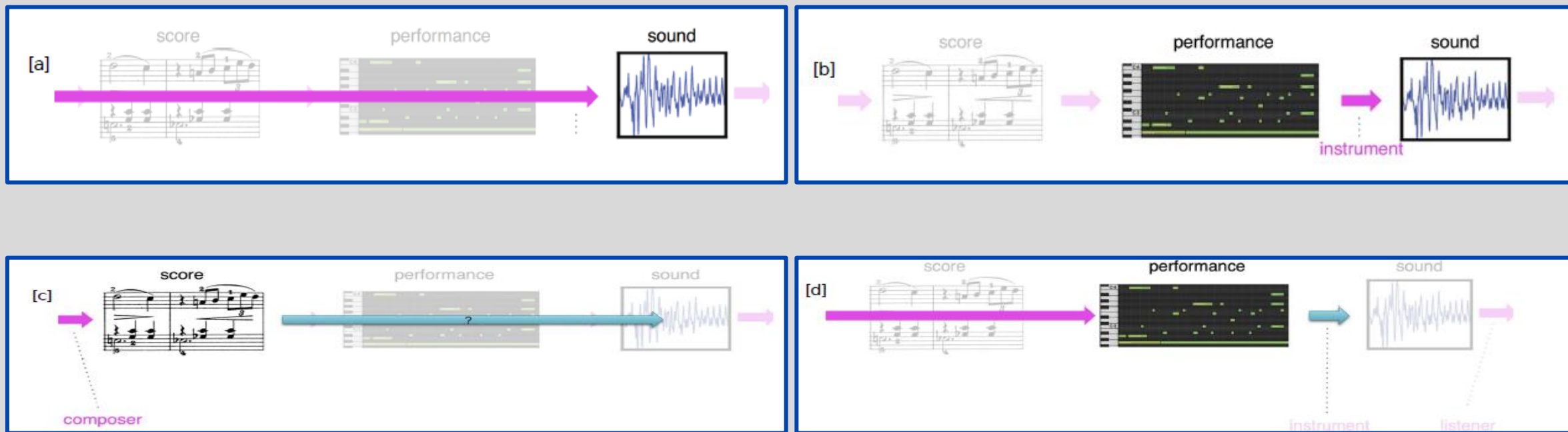
Źródło: <https://leadingtone.tumblr.com/post/4006499072/christus-der-uns-selig-macht-chorale-prelude-bwv>




Źródło: <https://www.theguardian.com/music/2015/may/11/berlin-philharmonic-imon-rattle-orchestra-vote-chief-conductor>

# Istotność ekspresji w muzyce

Oore S., Simon I., Dieleman S., Eck D., Simonyan K.: *This Time with Feeling: Learning Expressive Musical Performance*

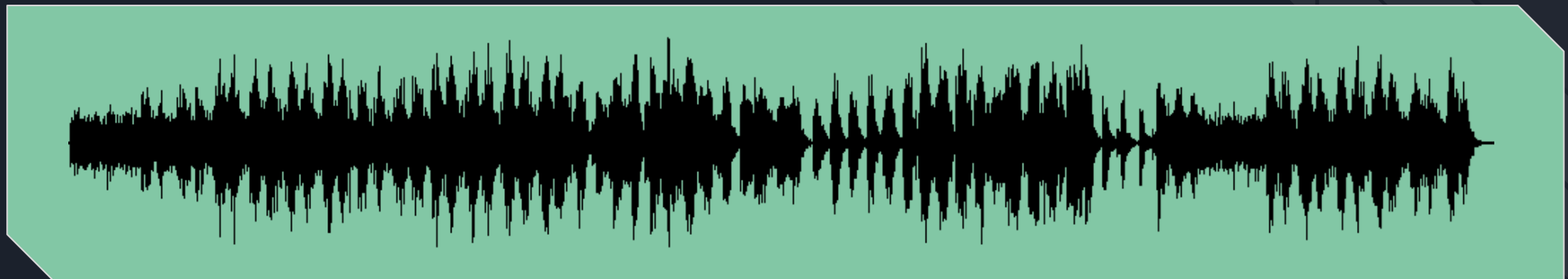




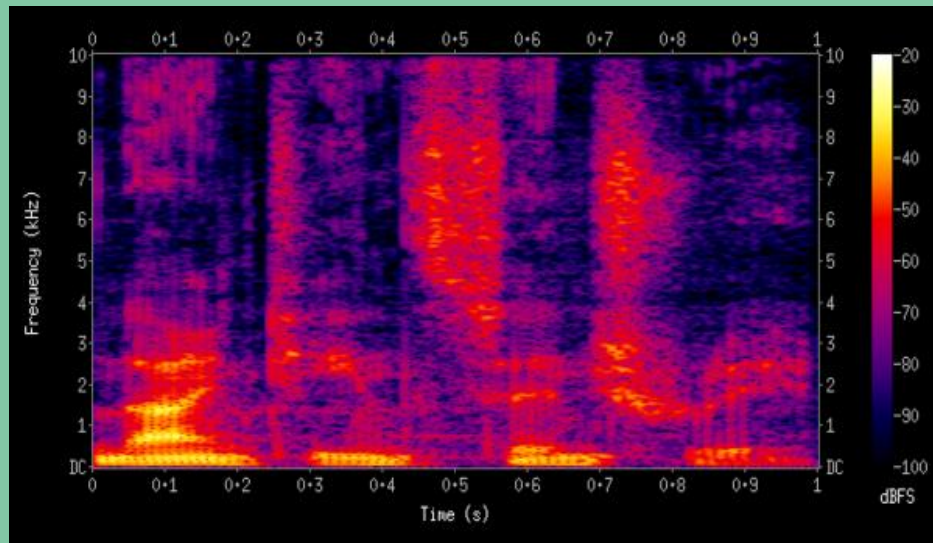
REPREZENTACJA  
EKSPRESYJNYCH  
DANYCH  
MUZYCZNYCH

# Zapis fali dźwiękowej (format *wave*)

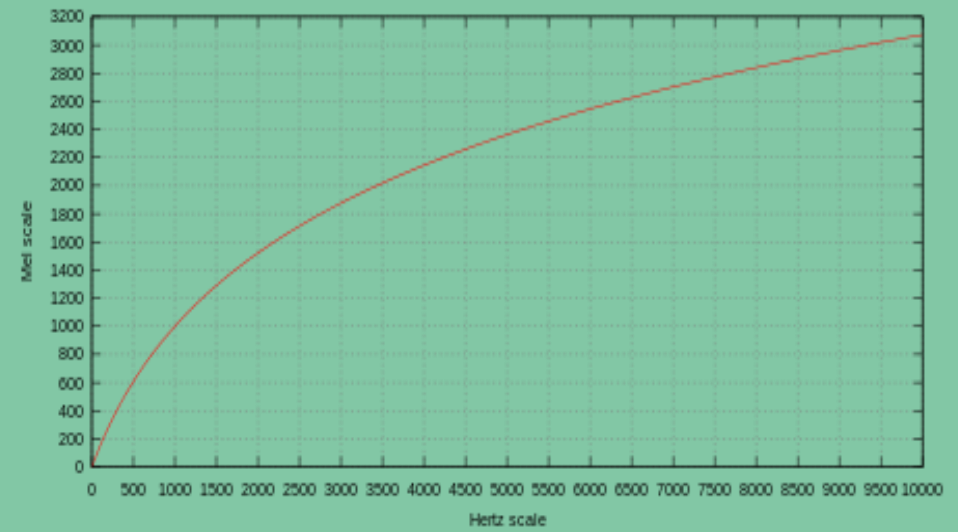
- Wszystkie możliwe informacje o wykonaniu
- Mało przystępna postać danych


































# Spektrogram



# Skala *mel*



# Standard MIDI (format *mid*)

MIDI number	Note name	Keyboard	Frequency
21	A0		27.500
22	B0		30.868
23	C1		32.703
24	D1		36.708
25	E1		38.891
26	F1		41.203
27	G1		43.654
28	A1		46.249
29	B1		48.999
30	C2		51.913
31	D2		55.000
32	E2		58.270
33	F2		61.735
34	G2		65.406
35	A2		69.296
36	B2		73.416
37	C3		77.782
38	D3		82.407
39	E3		87.307
40	F3		92.499
41	G3		97.999
42	A3		103.83
43	B3		110.00
44	C4		116.54
45	D4		123.47
46	E4		130.81
47	F4		138.59
48	G4		146.83
49	A4		155.56
50	B4		164.81
51	C5		174.61
52	D5		185.00
53	E5		196.00
54	F5		207.65
55	G5		220.00
56	A5		233.08
57	B5		246.94
58	C6		261.63

- MIDI – *Musical Instrument Digital Interface* (ang. cyfrowy interfejs instrumentów muzycznych)
- Sekwencja zdarzeń
- Opis ekspresji przez długość trwania nuty i głośność



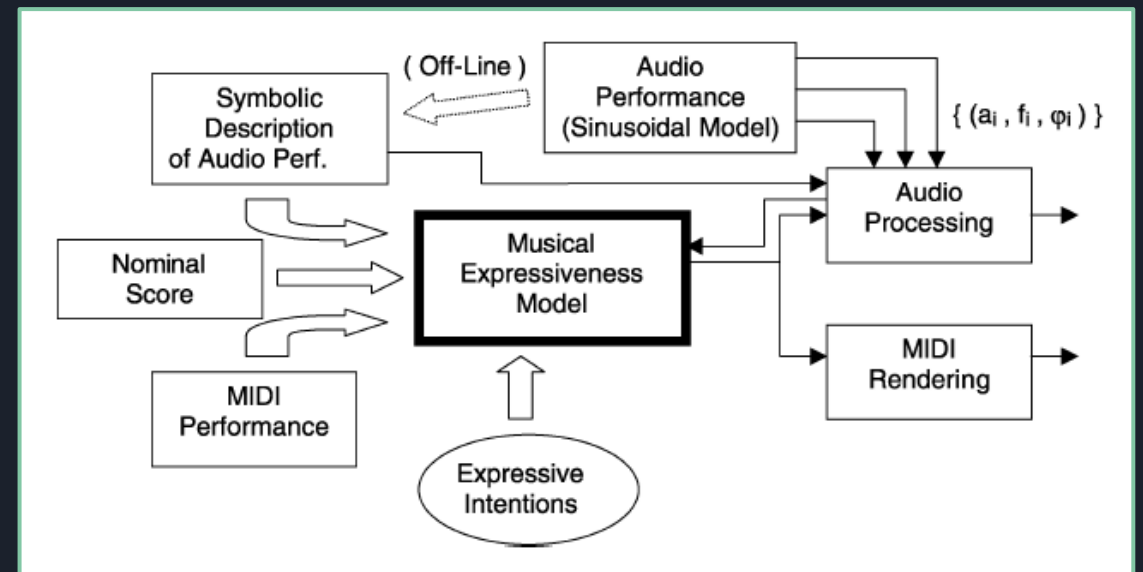
PODEJŚCIA STATE-  
OF-THE-ART



# Modyfikacja istniejącego wykonania

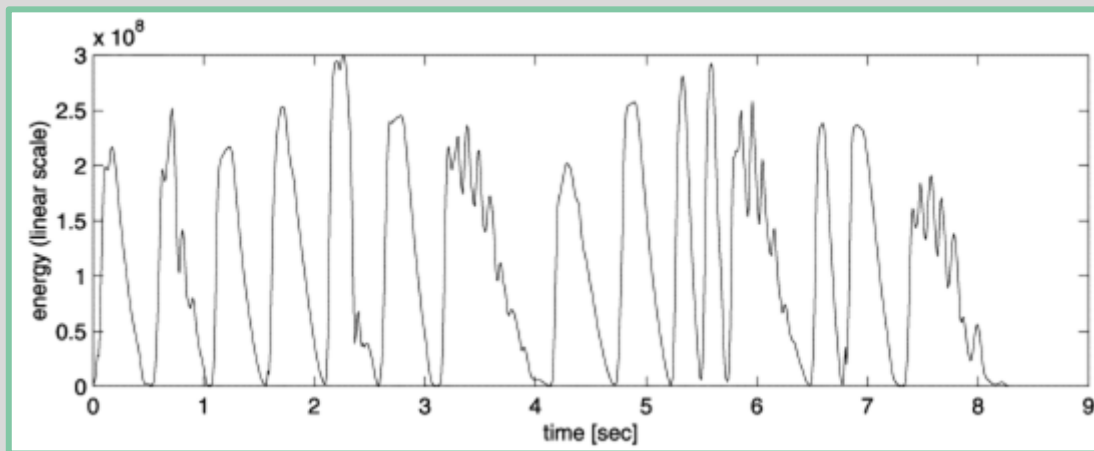
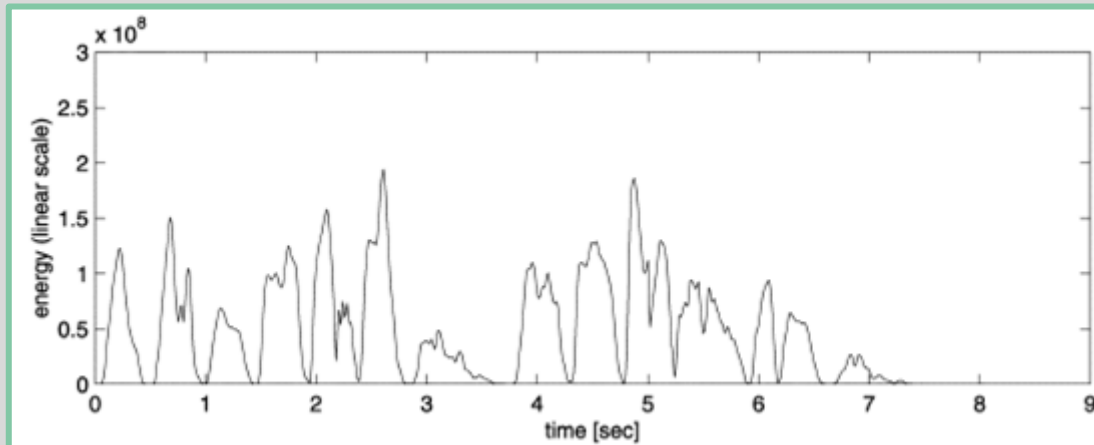
Canazza S., De Poli G., Drioli C., Roda A., Vidolin A.: *Modeling and Control of Expressiveness in Music Performance*. Proceedings of the IEEE (2004)

- Dane wejściowe – MIDI i zapis fali
- Modyfikacja ekspresyjności poprzez:
  - przesunięcie zdarzenia
  - rozciągnięcie/kompresję zdarzenia
- Możliwe intencje: *jasne, ciemne, twarde, miękkie, ciężkie, lekkie*
- Regresja liniowa



# Generowanie muzyki ekspresyjnej

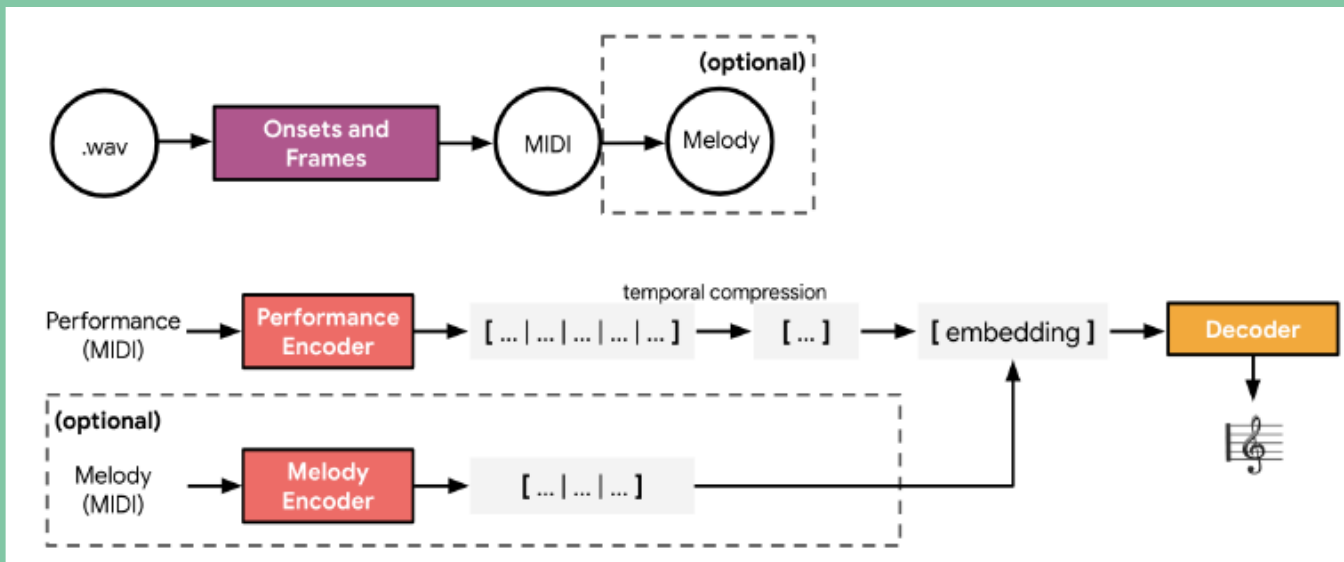
Canazza S., De Poli G., Drioli C., Roda A., Vidolin A.: *Modeling and Control of Expressiveness in Music Performance*. Proceedings of the IEEE (2004)



Energia  
wykonania  
neutralnego  
(górze) i wykonania  
ciężkiego (dół)

# Generowanie na podstawie MIDI

Choi K., Hawthorne C., Simon I., Dinculescau M., Engel J.: *Encoding musical style with Transformer Autoencoders*

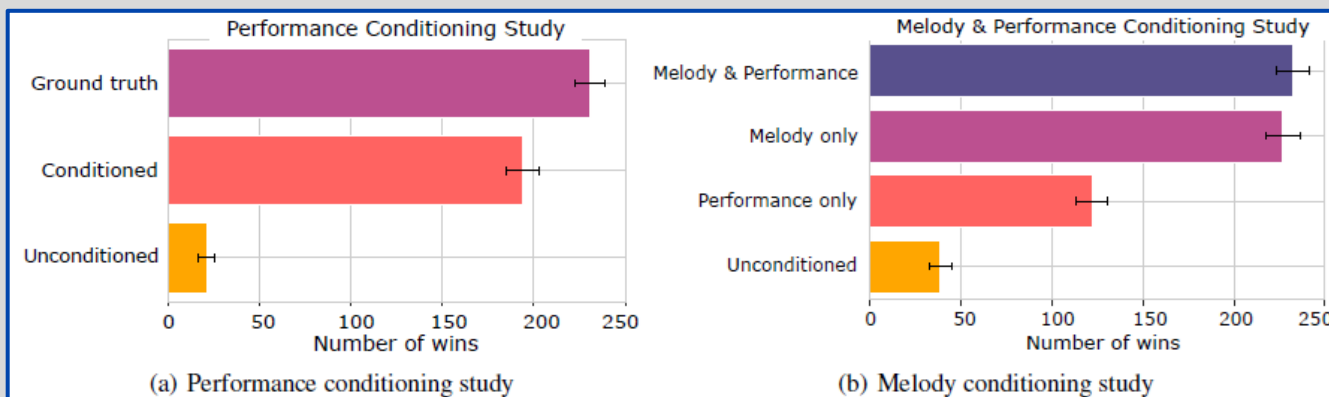


- Model *encoder-decoder*
- Opcjonalne korzystanie z reprezentacji melodii
- Wszystkie części oparte o Transformery

# Generowanie na podstawie MIDI

Choi K., Hawthorne C., Simon I., Dinculescau M., Engel J.: *Encoding musical style with Transformer Autoencoders*

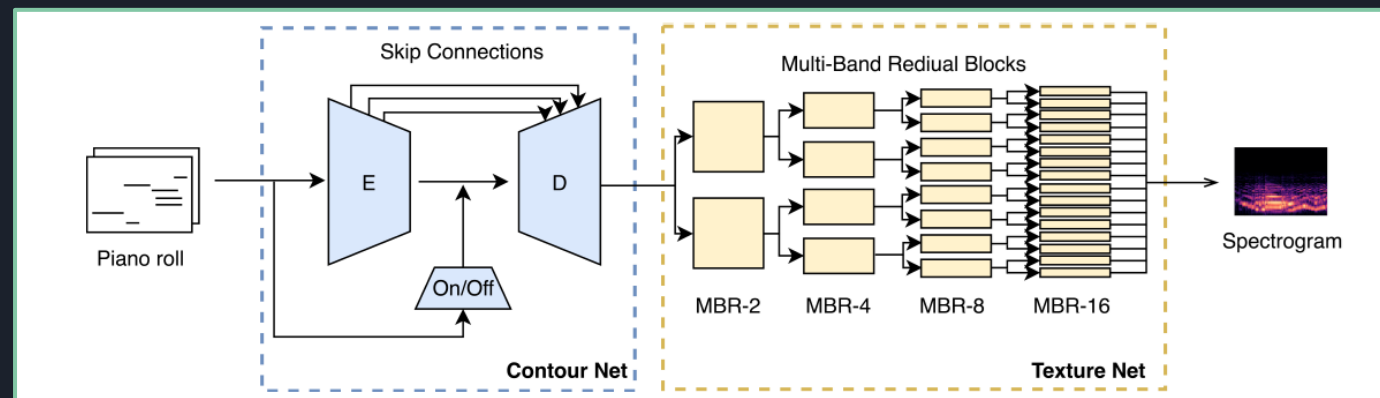
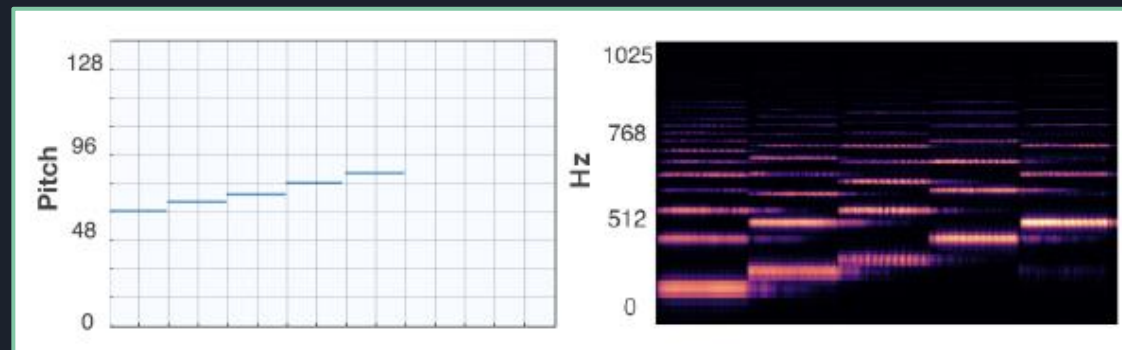
Model variation	MAESTRO	YouTube
Melody-only Transformer with rel. attention (Huang et al., 2019b)	1.786	1.302
Melody & performance autoencoder with rel. attention, sum (ours)	<b>1.706</b>	1.275
Melody & performance autoencoder with rel. attention, concat (ours)	1.713	<b>1.237</b>
Melody & performance autoencoder with rel. attention, tile (ours)	1.709	1.248



# Generowanie spektrogramu

Wang B., Yang Y.: *PerformanceNet: Score-to-Audio music generation with multi-band residual network*

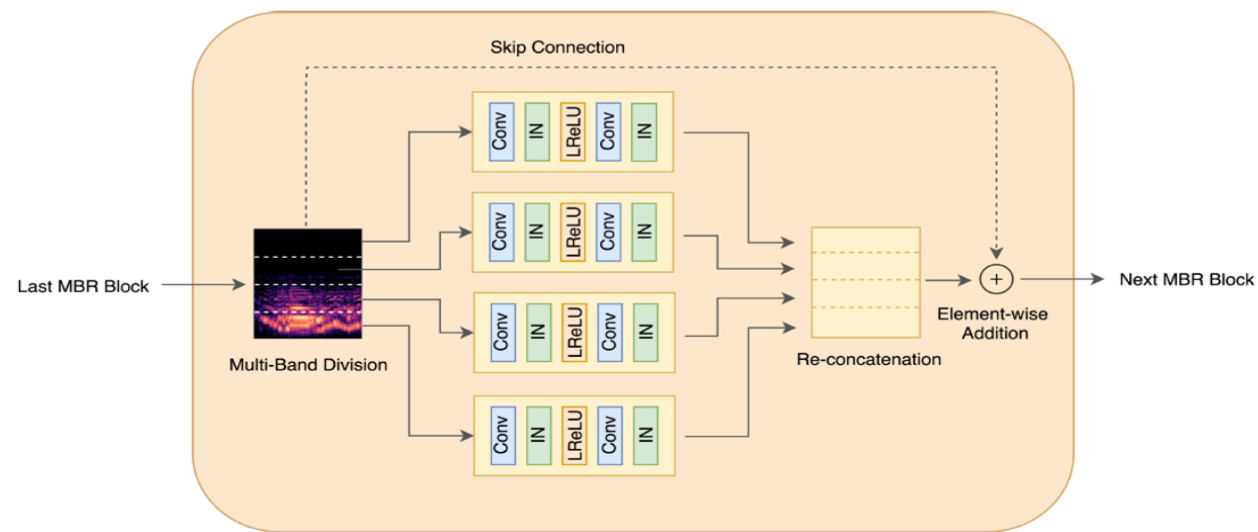
- Wstępne tworzenie spektrogramu
- Uszczegóławianie (zwiększenie rozdzielczości)



# Generowanie spektrogramu

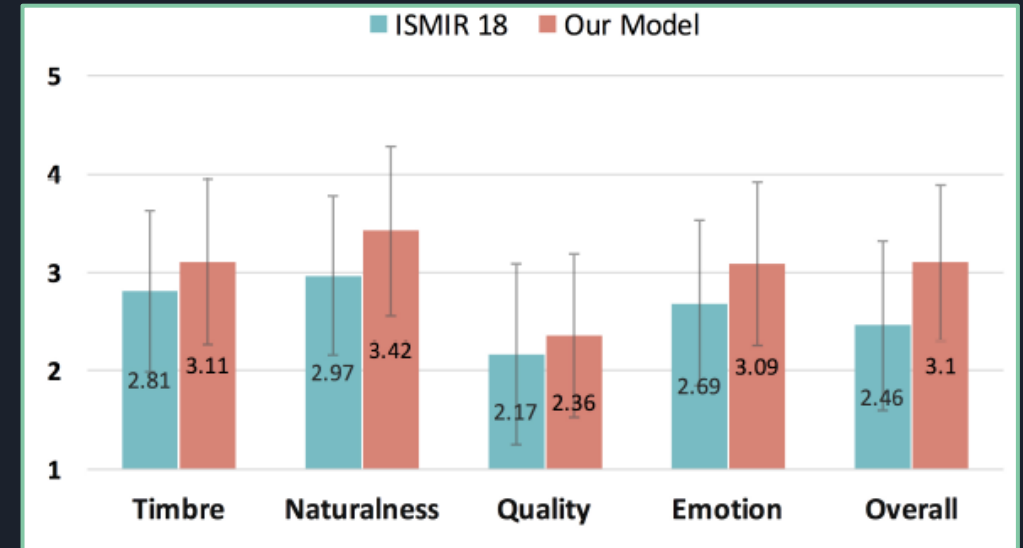
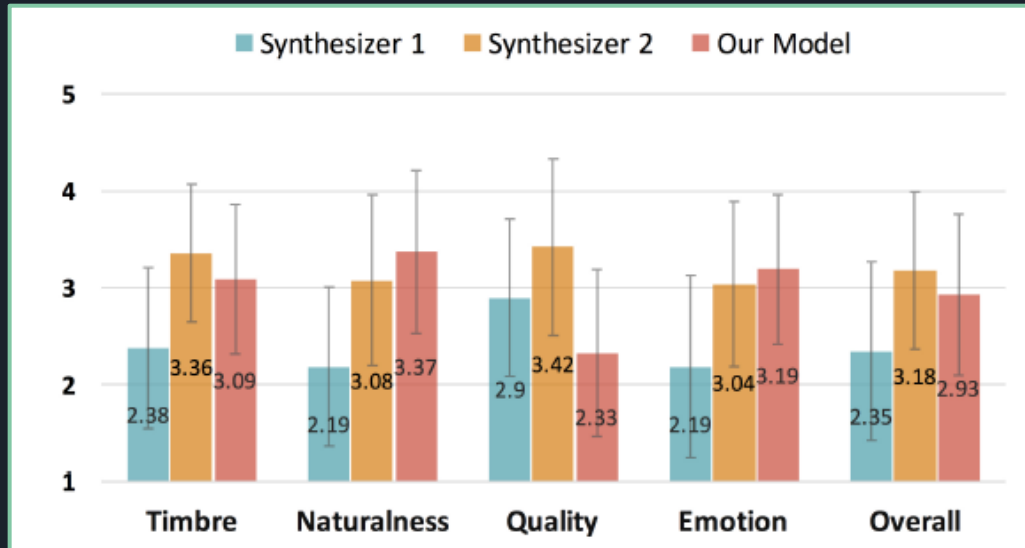
Wang B., Yang Y.: *PerformanceNet: Score-to-Audio music generation with multi-band residual network*

- Podział spektrogramu względem częstotliwości
- Uczenie różnych charakterystyk dla różnych częstotliwości



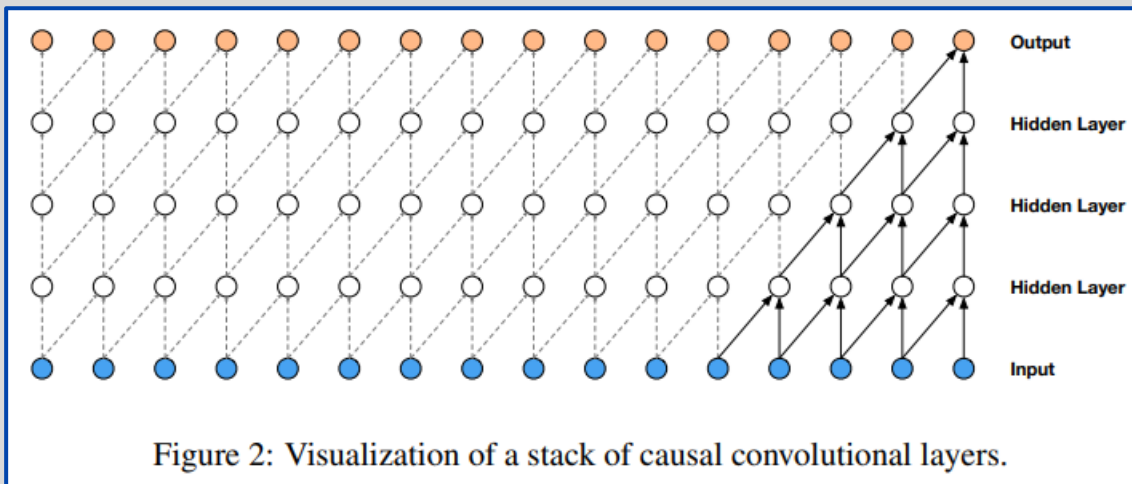
# Generowanie spektrogramu

Wang B., Yang Y.: *PerformanceNet: Score-to-Audio music generation with multi-band residual network*

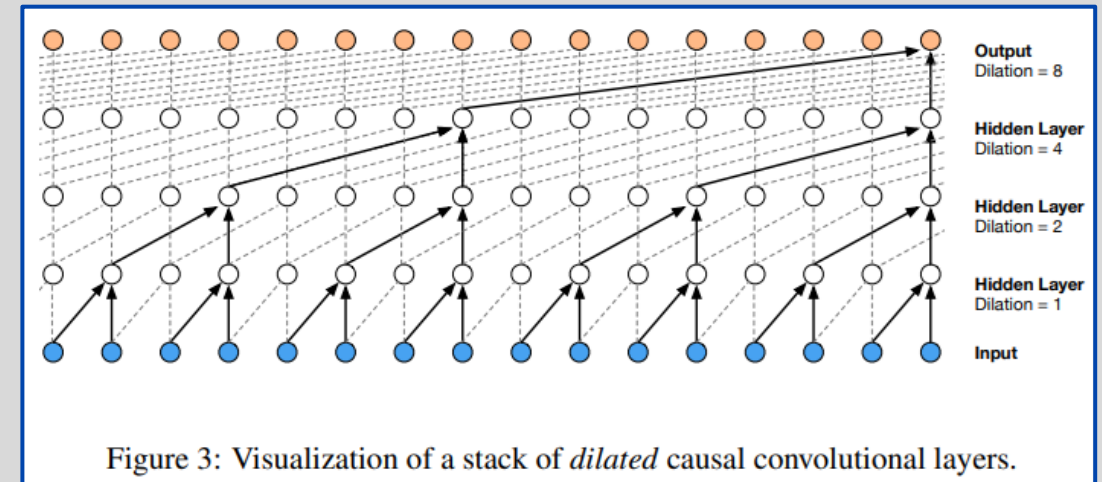


# Generowanie wave - inspiracja

van den Oord A., Dieleman S., Heiga Z., Simonyan K., Vinyals O., Graves A., Kalchbrenner N., Senior A., Kavukcuoglu K.: *WaveNet: A generative model for raw audio*



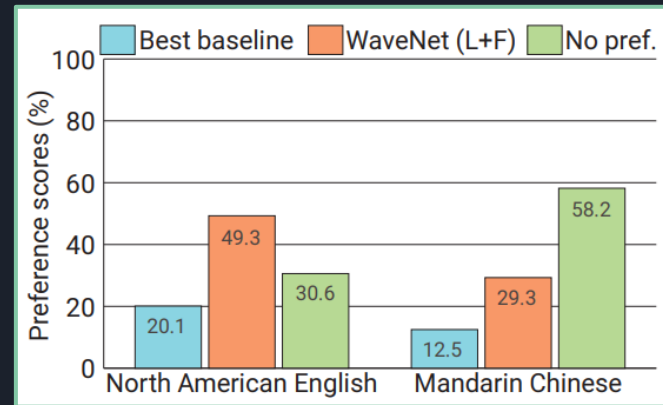
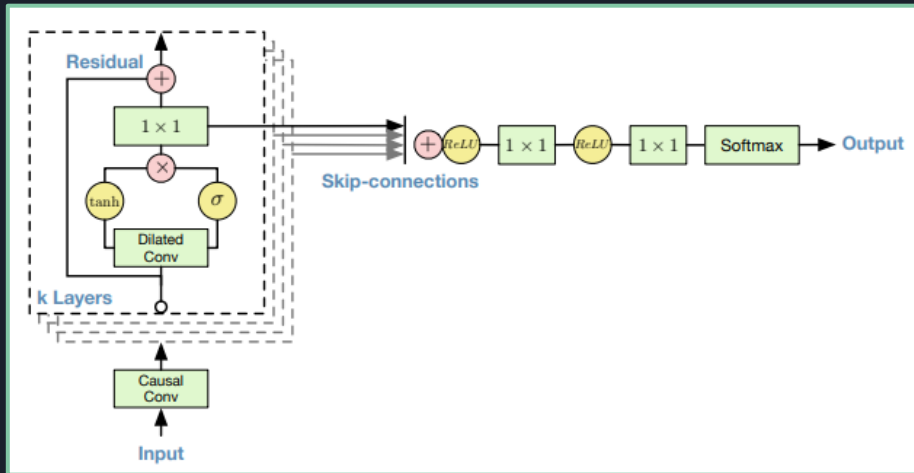
- Model text-to-speech
- Poszerzenie kontekstu





# Generowanie *wave* - inspiracja

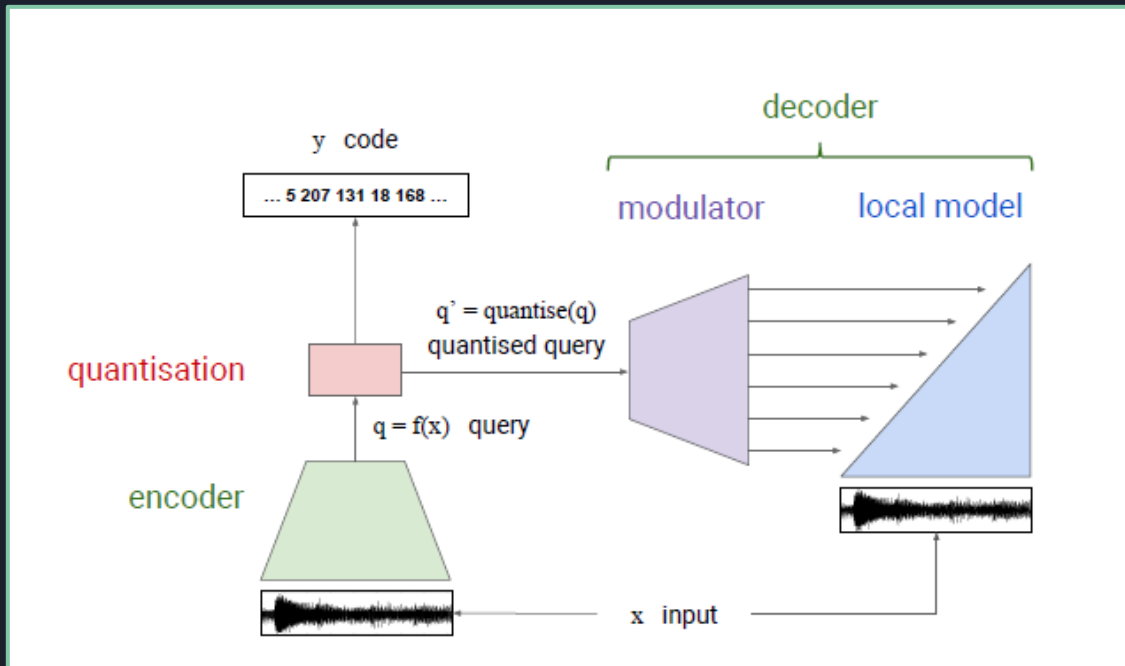
van den Oord A., Dieleman S., Heiga Z., Simonyan K., Vinyals O., Graves A., Kalchbrenner N., Senior A., Kavukcuoglu K.: *WaveNet: A generative model for raw audio*



- Wyżej oceniany niż pozostałe modele text-to-speech

# Generowanie wave

Dieleman S., van den Oord A., Simonyan K.: *The challenge of realistic music generation: modeling raw audio at scale*

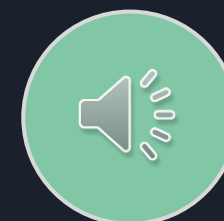
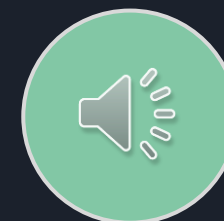


- Model *encoder-decoder*
- Wszystkie części oparte o WaveNet
- *Bottleneck* tworzony przy pomocy dyskretyzacji danych

# Generowanie wave

Dieleman S., van den Oord A., Simonyan K.: *The challenge of realistic music generation: modeling raw audio at scale*

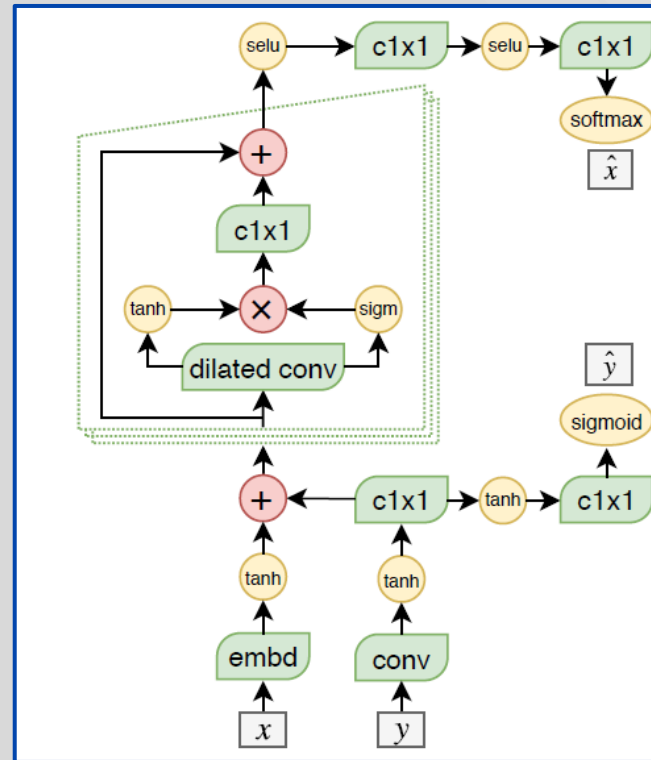
MODEL	NUM. LEVELS	RF	HUMAN EVALUATION	
			FIDELITY	MUSICALITY
Large WaveNet	1	384 ms	$3.82 \pm 0.18$	$2.43 \pm 0.14$
Very large WaveNet	1	768 ms	$3.82 \pm 0.20$	$2.89 \pm 0.17$
Thin WaveNet with large RF	1	3072 ms	$2.43 \pm 0.17$	$1.71 \pm 0.18$
hop-8 VQ-VAE + large WaveNet	2	3072 ms	$3.79 \pm 0.16$	$3.04 \pm 0.16$
hop-64 VQ-VAE + large WaveNet	2	24576 ms	$3.54 \pm 0.18$	$3.07 \pm 0.17$
VQ-VAE + PBT-VQ-VAE + large WaveNet	3	24576 ms	$3.71 \pm 0.18$	$4.04 \pm 0.14$
VQ-VAE + AMAE + large WaveNet	3	24576 ms	$3.93 \pm 0.18$	$3.46 \pm 0.15$



# Generowanie *wave*, transfer stylu

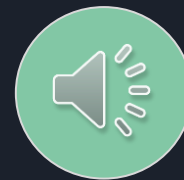
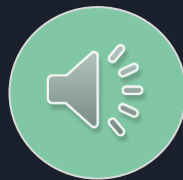
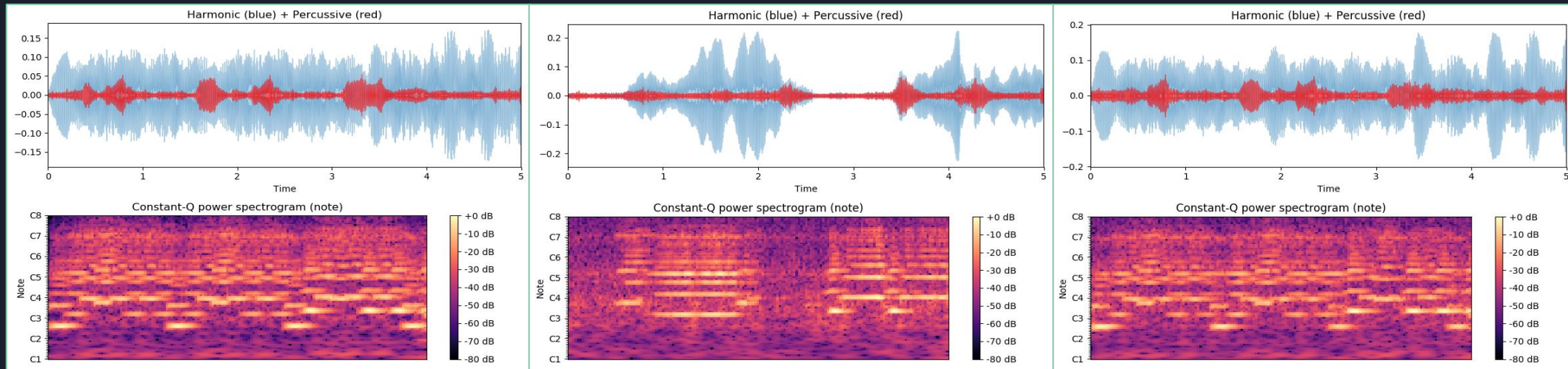
Schimbinschi F., Walder C., Erfani S. M., Bailey J.: *SynthNet: Learning to synthesize music end-to-end*

- Model oparty o WaveNet
- Korzysta zarówno z MIDI jak i *waveform*



# Generowanie wave, transfer stylu

Schimbinschi F., Walder C., Erfani S. M., Bailey J.: *SynthNet: Learning to synthesize music end-to-end*



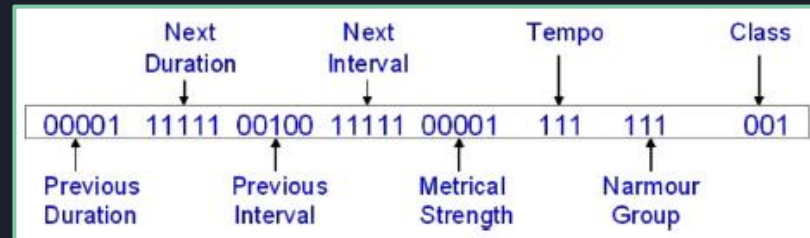
Experiment	Piano	Glockenspiel	Guitar	Cello	Trumpet	Flute	Square	All
WaveNet L26	2.22±0.25	2.48±0.23	2.18±0.25	2.37±0.28	2.18±0.29	2.37±0.22	2.30±0.09	2.30±0.10
DeepVoice L26	2.55±0.32	1.85±0.23	2.30±0.39	2.62±0.27	2.28±0.32	2.20±0.25	1.87±0.03	2.24±0.11
SynthNet L24	4.75±0.14	4.45±0.17	4.30±0.19	4.50±0.15	4.25±0.18	4.15±0.21	4.10±0.16	<b>4.36±0.07</b>

# Tworzenie reguł na podstawie wave

Ramirez R., Maestre E., Serra X.: *A rule-based evolutionary approach to music performance modeling*

```
GeneticSeqCovAlg(Class, Fitness, Threshold, p, r, m, Examples)
  Pos = examples which belong to Class
  Neg = examples which do not belong to Class
  Learned_rules = {}
  While Pos do
    P = generate p hypotheses at random
    For each hypothesis h in P,
      Compute fitness(h)
    While max(fitness(h)) < Threshold and #generations < 400 do
      Create a new generation Pnew
      P = Pnew
      For each h in P,
        Compute fitness(h)
    NewRule = the hypothesis in P that has the highest fitness
    Rpos = members of Pos covered by NewRule
    Compute PredictedValue(Rpos)
    NumericNewRule = NewRule with Class replaced by
                      Regression(Rpos)

    Learned_rules = Learned_rules + NumericNewrule
    Pos = Pos - Rpos
  Return Learned_rules
```



- Zbiór reguł opisujący zależności między dźwiękami
- Funkcja dopasowania  $\frac{tP^\alpha}{tP+fP}$

# Tworzenie reguł na podstawie wave

Ramirez R., Maestre E., Serra X.: *A rule-based evolutionary approach to music performance modeling*

Attribute	Bit Meaning
Previous duration	<i>Bit 1</i> : much shorter than current note <i>Bit 2</i> : shorter than current note <i>Bit 3</i> : same than current note <i>Bit 4</i> : longer than current note <i>Bit 5</i> : much longer than current note
Next duration	<i>Bit 1</i> : much shorter than current note <i>Bit 2</i> : shorter than current note <i>Bit 3</i> : same than current note <i>Bit 4</i> : longer than current note <i>Bit 5</i> : much longer than current note
Previous interval	<i>Bit 1</i> : much lower than current note <i>Bit 2</i> : lower than current note <i>Bit 3</i> : same than current note <i>Bit 4</i> : higher than current note <i>Bit 5</i> : much higher than current note
Next interval	<i>Bit 1</i> : much lower than current note <i>Bit 2</i> : lower than current note <i>Bit 3</i> : same than current note <i>Bit 4</i> : higher than current note <i>Bit 5</i> : much higher than current note

Metrical strength	<i>Bit 1</i> : very weak <i>Bit 2</i> : weak <i>Bit 3</i> : medium <i>Bit 4</i> : strong <i>Bit 5</i> : very strong
Tempo	<i>Bit 1</i> : slow <i>Bit 2</i> : nominal <i>Bit 3</i> : fast
Narmour group	000 = "P" group 001 = "D" group 010 = "ID" group 011 = "IP" group 100 = "VP" group 101 = "R" group 110 = "IR" group 111 = "VR" group
Class	<i>Bit 1</i> : shorten/delay/soft/ornamentation <i>Bit 2</i> : same-duration/same-onset/medium/none <i>Bit 3</i> : lengthen/advance/loud/none

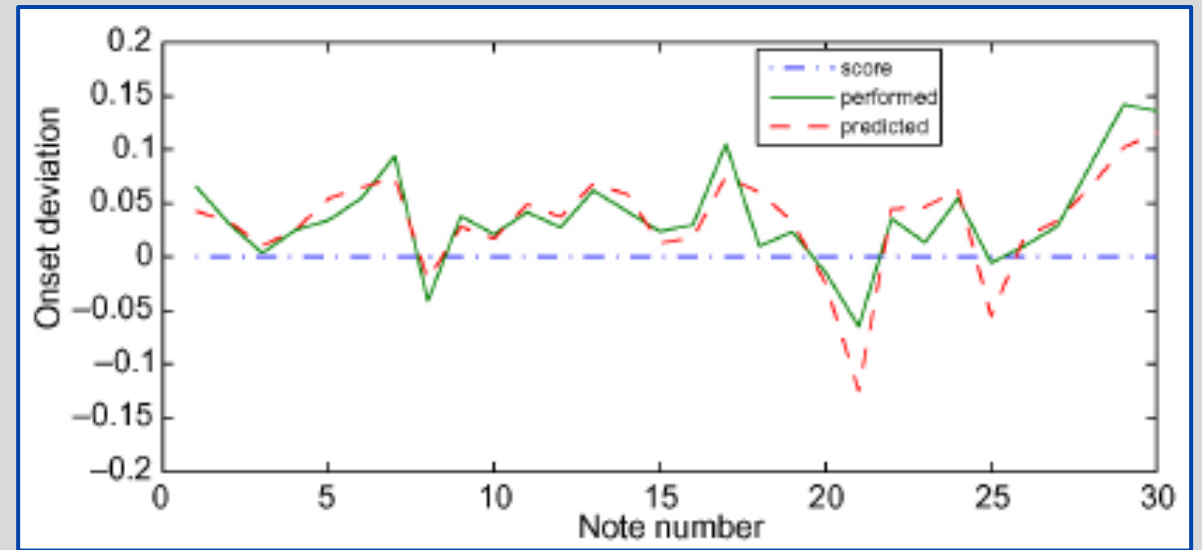
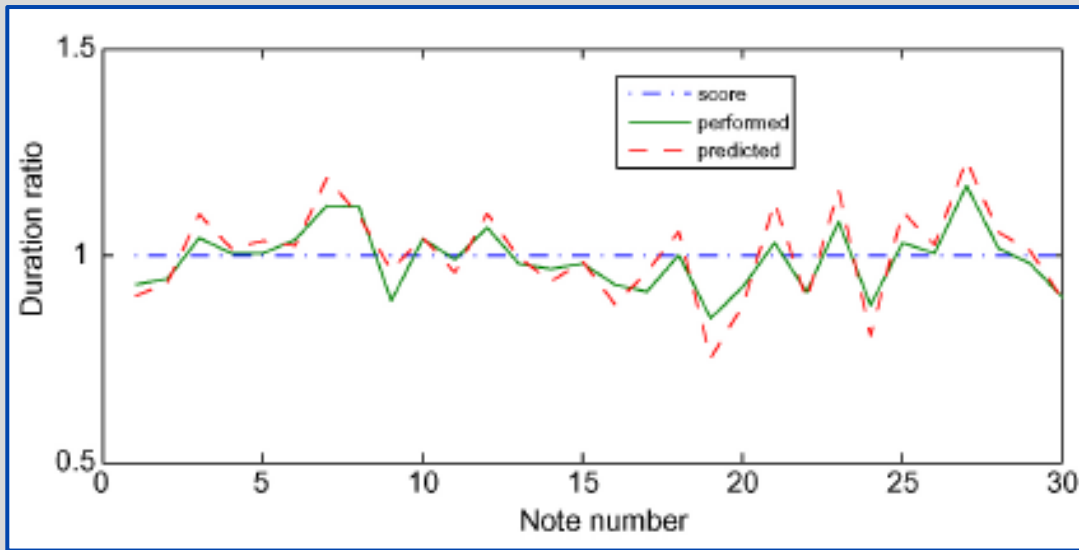
*"if the previous note duration is much longer and its pitch is the same and it is in a very strong metrical position and the current note appears in Narmour group R then lengthen the duration of the current note"*

is coded as the binary string

00001 11111 00100 11111 00001 111 110 001.

# Tworzenie reguł na podstawie *wave*

Ramirez R., Maestre E., Serra X.: *A rule-based evolutionary approach to music performance modeling*





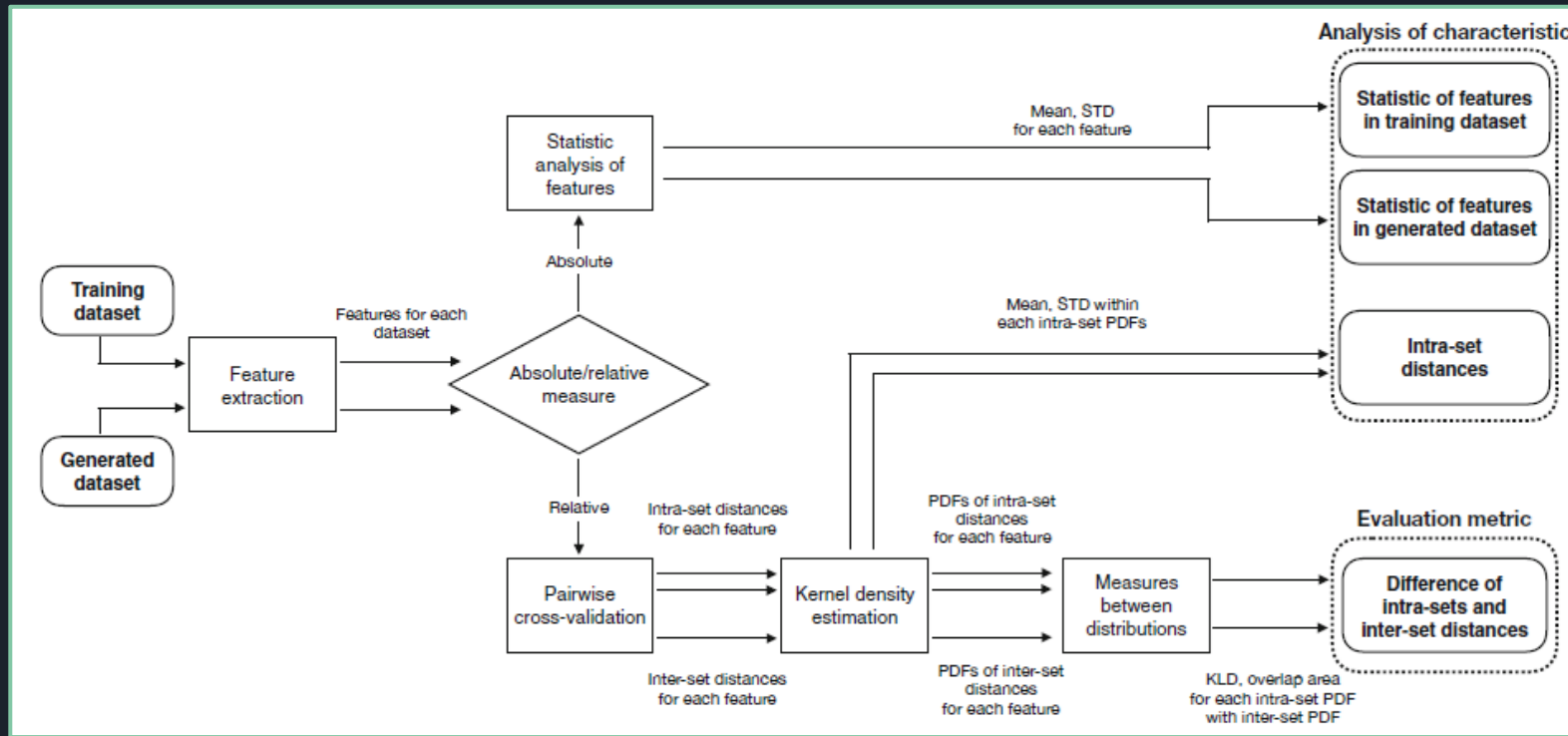


# Problemy i wyzwania

- Jakość generowanej muzyki
- Ocena generowanej muzyki

# Obiektywna ocena

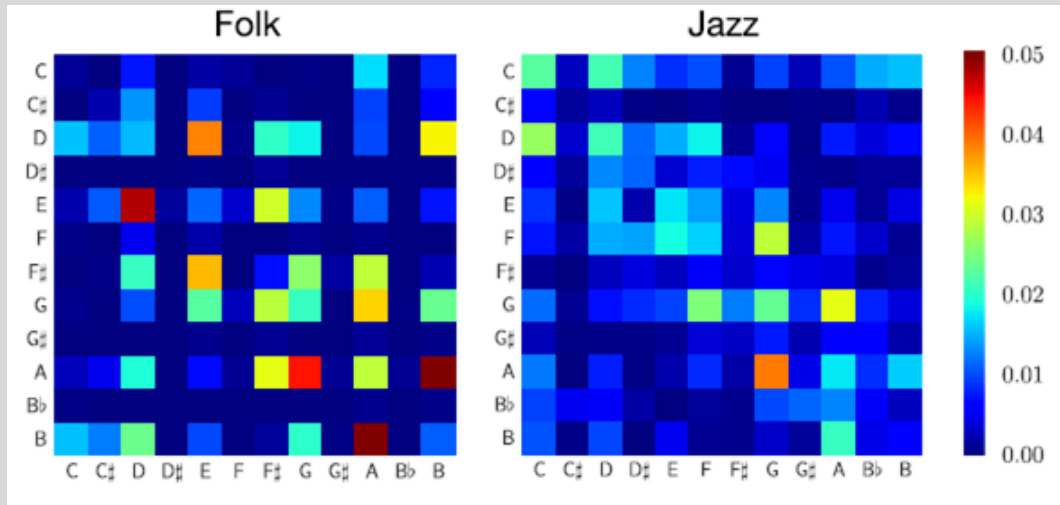
Yang L., Lerch A.: *On the evaluation of generative models in music*



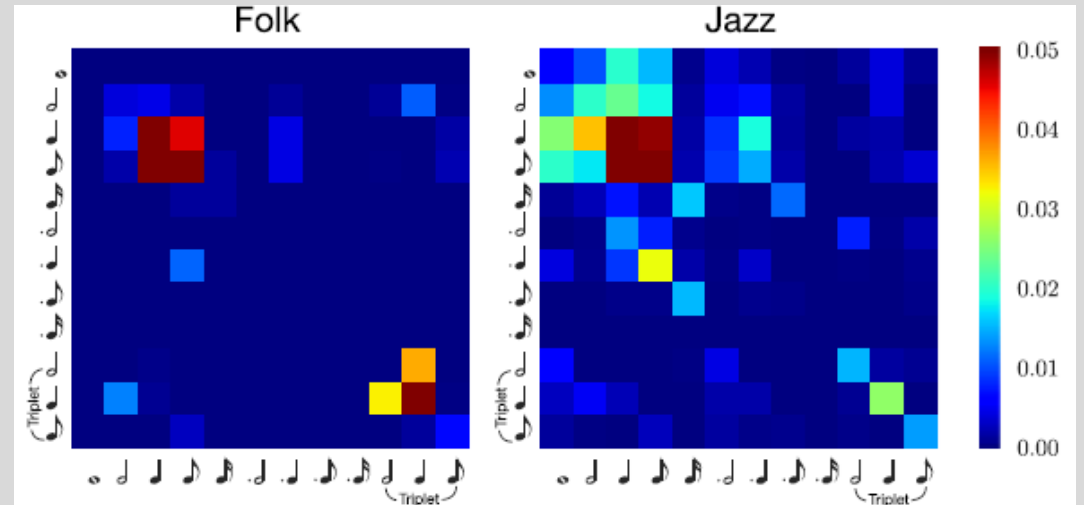
- Miara obiektywna i względna
- Oparta o właściwości statystyczne

# Obiektywna ocena - właściwości

Yang L., Lerch A.: *On the evaluation of generative models in music*



Macierz przejść  
dźwięków



Macierz przejść  
długości dźwięków

# Problem przydatności

Sturm B. L., Ben-Tal O., Monaghan U., Collins N., Herremans D., Chew E., Hadjeres G., Deruty E., Pachet F.: *Machine learning research that matters for music creation: A case study*

<b>Pieces for organ: The Glas Herry Comment &amp; X:7153</b> by folk-rnn + DeepBach (2017)	Richard Salmon <i>organ</i>
<b>Traditional Irish Sets</b> (with folk-rnn tunes in italics) <ul style="list-style-type: none"><li>• <b>Jigs</b> (The Cuil Aodha, The Dusty Windowsill, <i>The Glas Herry Comment</i>)</li><li>• <b>Slow Reels</b> (Maghera Mountain, X:2897)</li><li>• <b>Fast Reels</b> (The Rookery, X:1068, Toss The Feathers II)</li></ul>	Daren Banarsë and Musicians
<b>March to the Mainframe, Interlude, The Humours of Time Pigeon</b> by Bob L. Sturm + folk-rnn (2017)	Ensemble x.y
<b>Ed SheerAI vs XenAkIs vs Aldele</b> by Nick Collins (2017)	Ensemble x.y
<b>3 morphed pieces from "A Little Notebook for Anna Magdalena"</b> by J. S. Bach (1722) + MorpheuS (2017) <b>3 morphed pieces from "30 and 24 Pieces for Children"</b> by Kabalevsky (1937) + MorpheuS (2017)	Elaine Chew <i>piano</i>
<b>Safe Houses</b> by Úna Monaghan + folk-rnn (2017) <b>The Choice</b> by Úna Monaghan (2015) <b>The Chinwag</b> by Úna Monaghan (2015)	Úna Monaghan <i>Irish harp, concertina, electronics</i>
<b>Pieces for organ: X:633 &amp; The Drunken Pint</b> by folk-rnn + DeepBach (2017)	Richard Salmon <i>organ</i>
<b>Chicken Bits and Bits and Bobs</b> by Bob L. Sturm + folk-rnn (2017)	Ensemble x.y
<b>Bastard Tunes</b> by Oded Ben-Tal + folk-rnn (2017)	Ensemble x.y



DZIĘKUJĘ ZA  
UWAGĘ