



TECHNISCHE UNIVERSITÄT  
ILMENAU



Fraunhofer  
IDMT

# Generation of Symbolic Music Based on MusicVAE

Jakob Lerch

Supervisors:

Prof. Dr. Shardt

Dr. Andrew McLeod

Dr. Jakob Abeßer

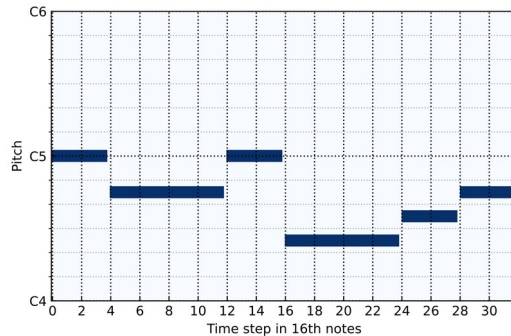
July 27, 2023



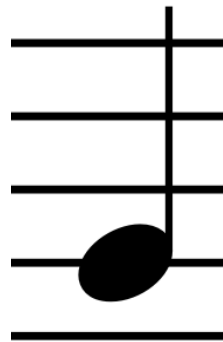
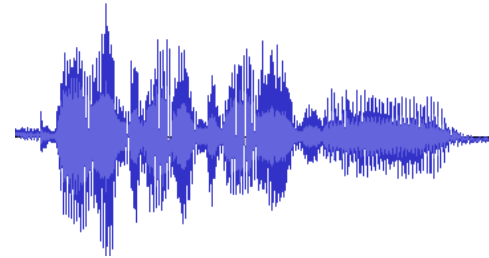
Music generation could be **supportive for composition or live performances.**

Picture retrieved July 16, 2023 from

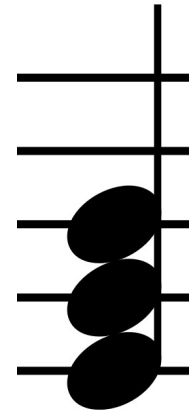
<https://pixabay.com/photos/music-producer-studio-actor-audio-4507819/>



or



or

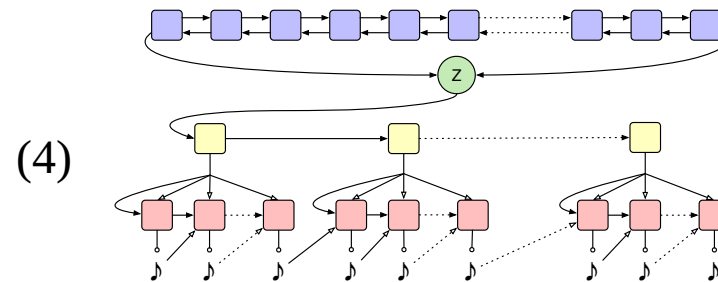
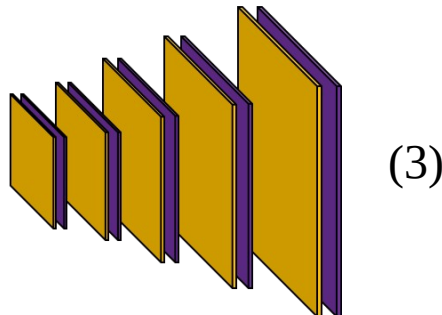
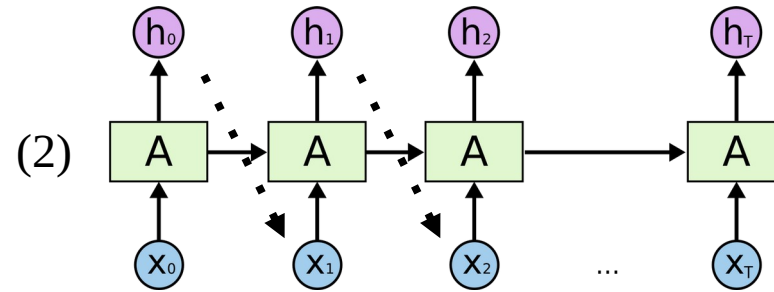
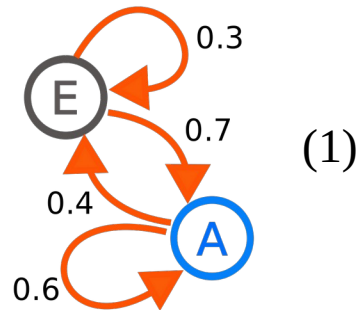


unconditioned or conditioned

- (1) review **state of the art**
- (2) **re-implement** MusicVAE
- (3) **evaluate quality** of generated excerpts



# State of the Art



(1) Picture retrieved July 17, 2023 from [https://en.wikipedia.org/wiki/Markov\\_chain#/media/File:Markovkate\\_01.svg](https://en.wikipedia.org/wiki/Markov_chain#/media/File:Markovkate_01.svg)

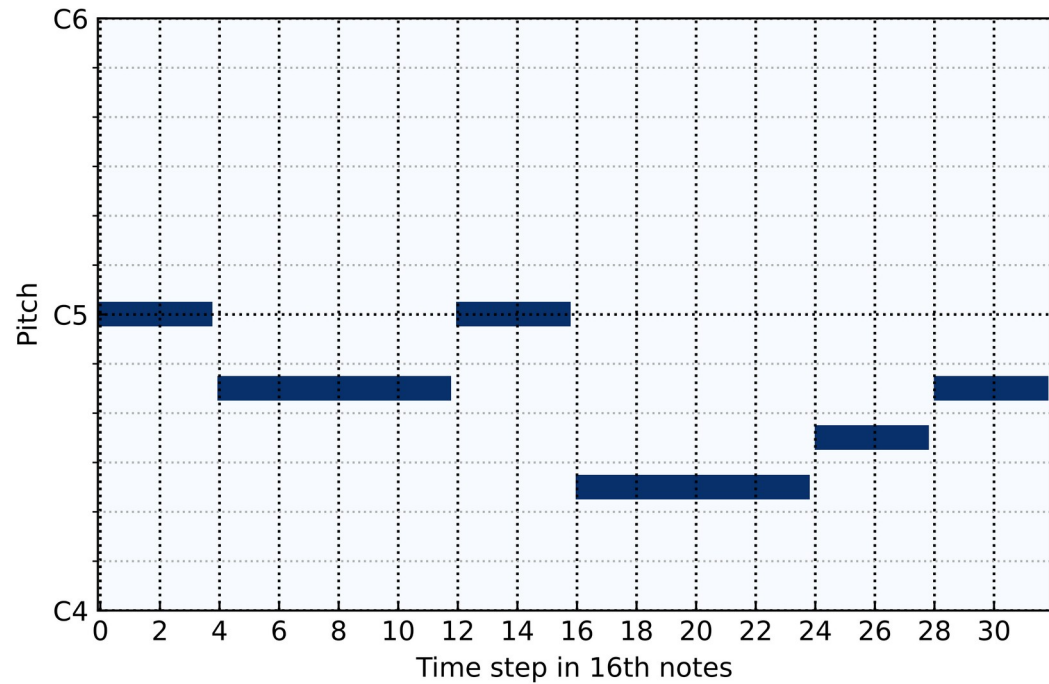
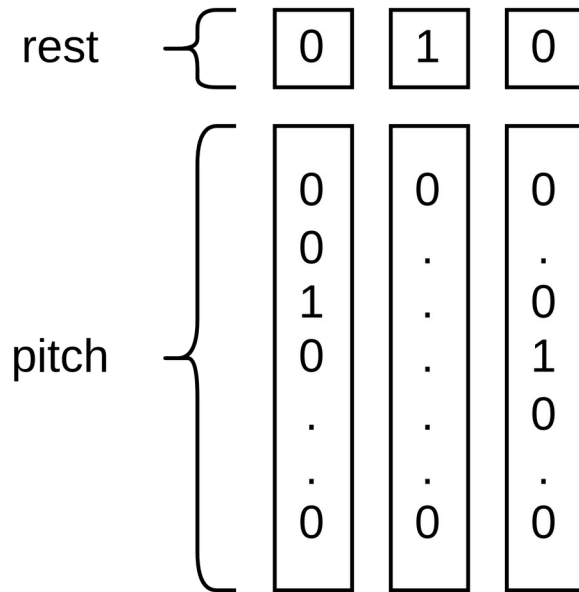
(2) Picture retrieved and changed July 17, 2023 from <https://colah.github.io/posts/2015-08-Understanding-LSTMs/img/RNN-unrolled.png>

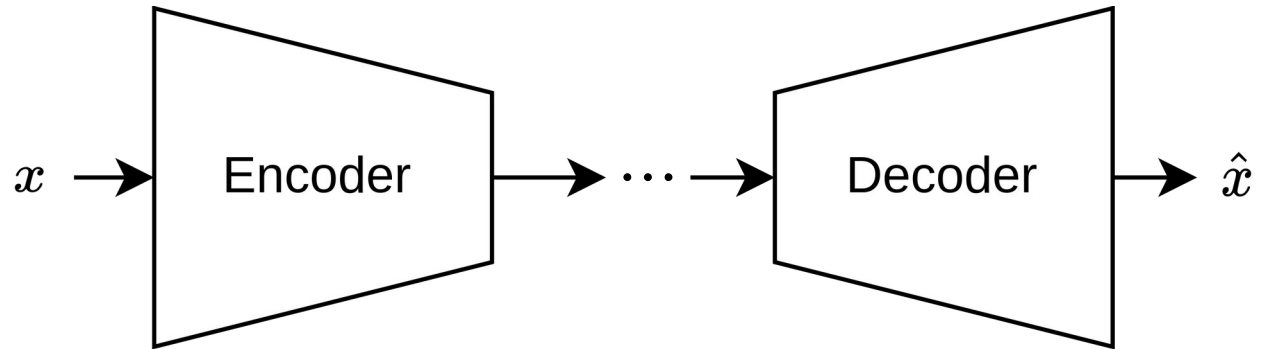
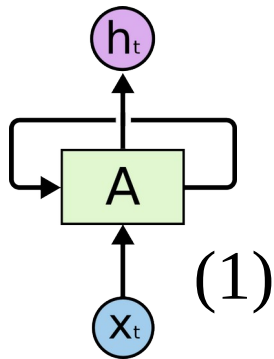
(3) Picture retrieved and changed July 17, 2023 from [https://clinicadl.readthedocs.io/en/latest/images/transfer\\_learning.png](https://clinicadl.readthedocs.io/en/latest/images/transfer_learning.png)

(4) Picture retrieved and changed July 17, 2023 from [https://magenta.tensorflow.org/assets/music\\_vae/architecture.png](https://magenta.tensorflow.org/assets/music_vae/architecture.png)

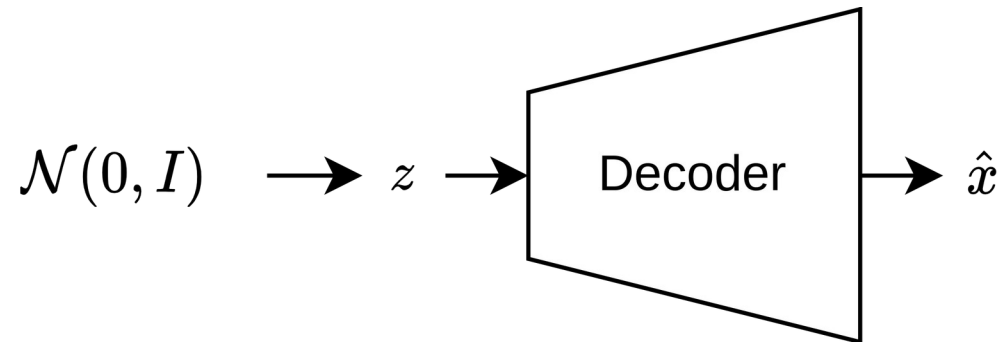
# MusicVAE

# Representation



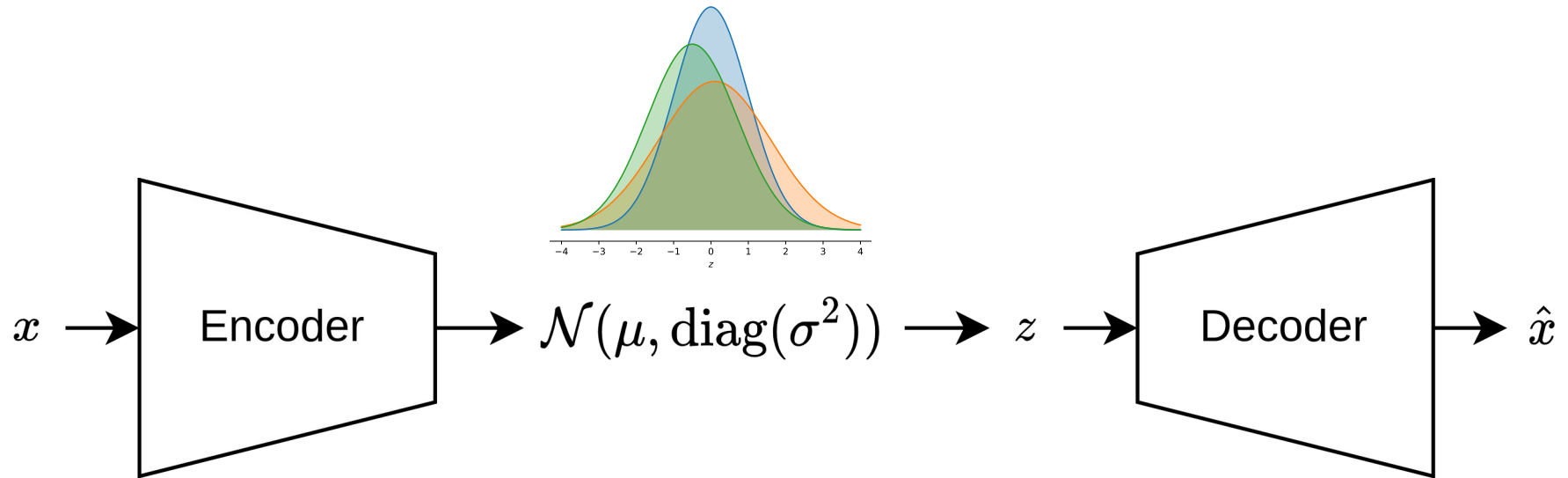


Train a model to  
generate music  
**from a random  
vector.**

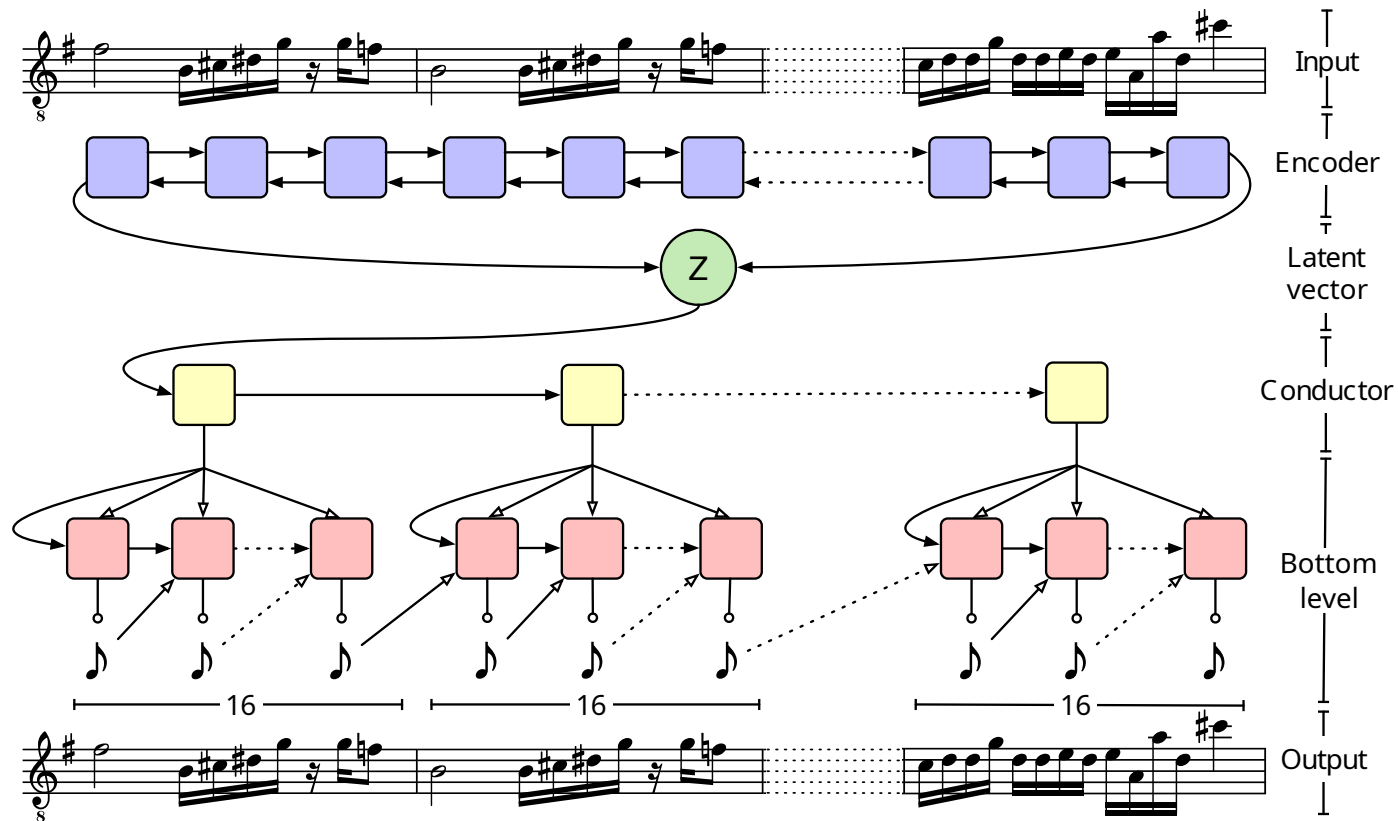


(1) Picture retrieved July 8, 2023 from

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/img/RNN-rolled.png>



$$L(x) = \text{rec. loss} + D_{KL}(\underbrace{\mathcal{N}(\mu, \text{diag}(\sigma^2))}_{\text{latent distribution}} || \underbrace{\mathcal{N}(0, I)}_{\text{prior distribution}})$$

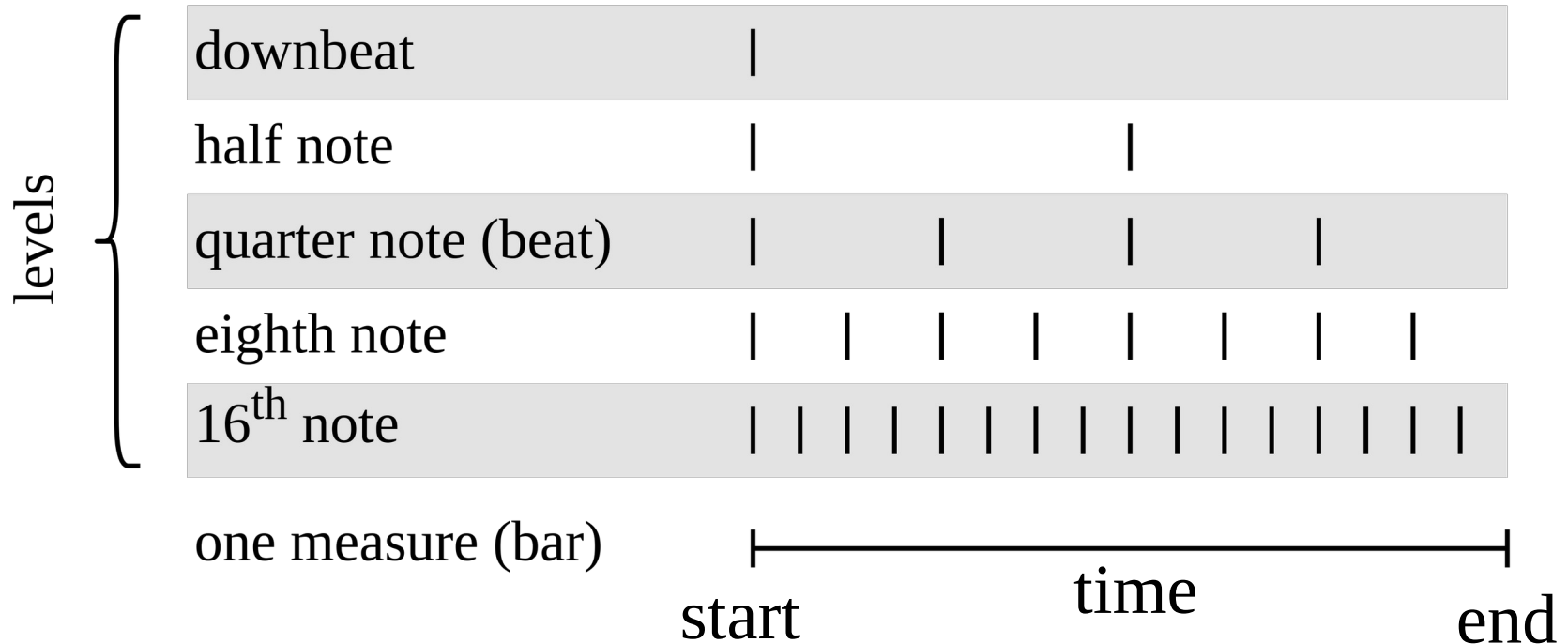


$$L(x) = \text{rec. loss} + \underline{\beta} \max[ D_{KL}, \underline{\lambda} ]$$

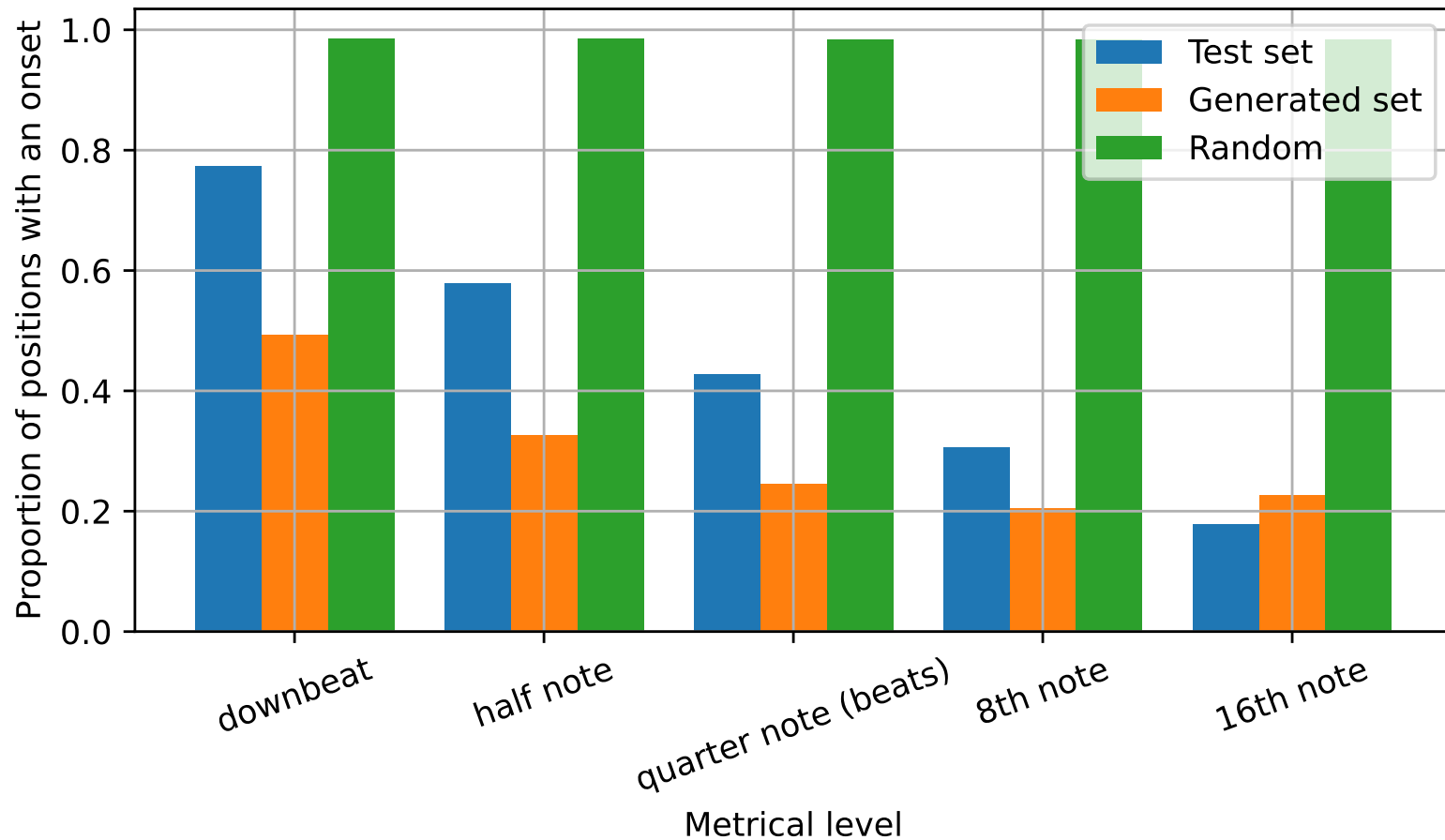
- **$\beta$  and  $\lambda$  chosen after grid search**
  - $\beta = 1; \lambda = 33.3$
- optimizer
  - Adam
  - learning rate (LR) =  $10^{-3}$
- batch size = 64
- weight decay
  - $L_2$  regularization with weight  $10^{-6}$
- LR scheduling
  - customized variant of ReduceLROnPlateau
- early stopping was used



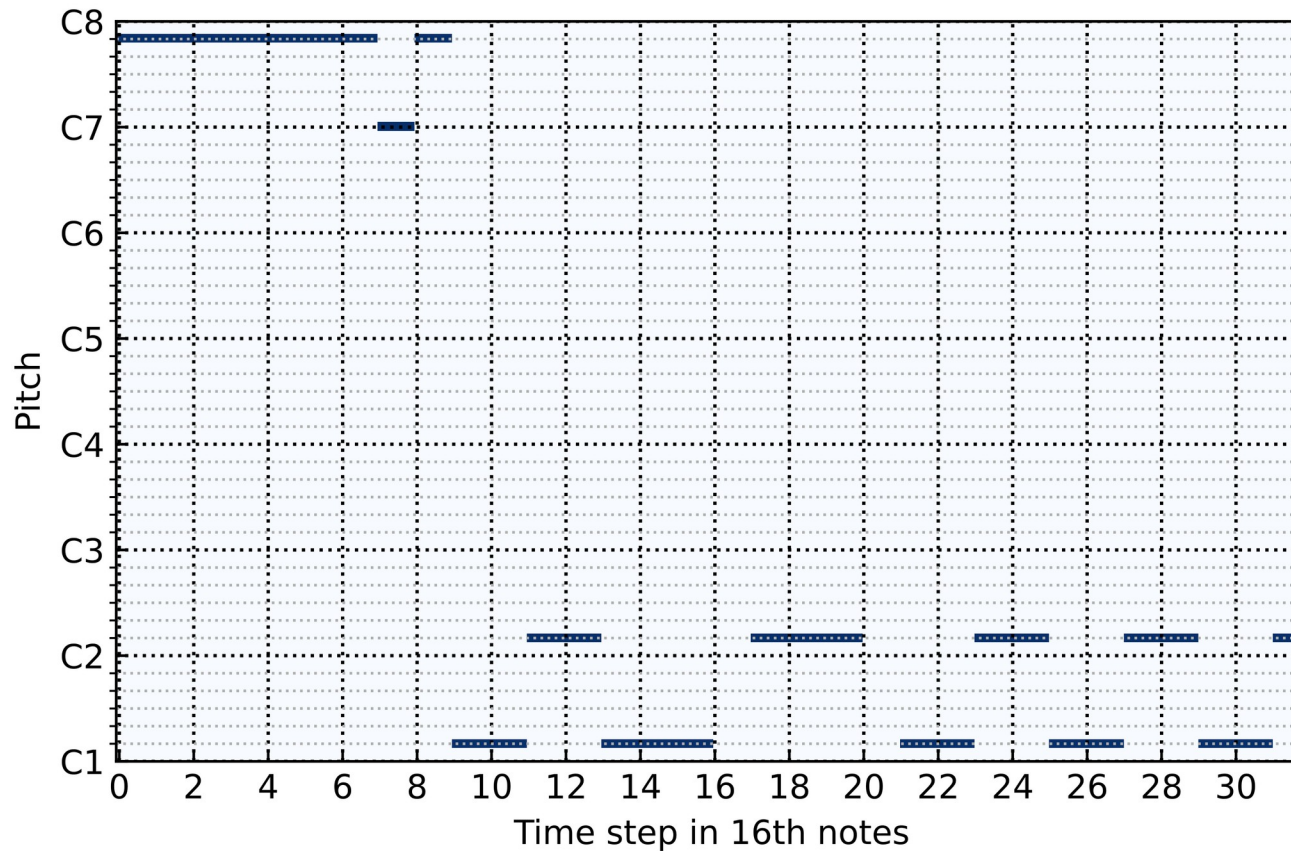
# Results: Rhythmic Features



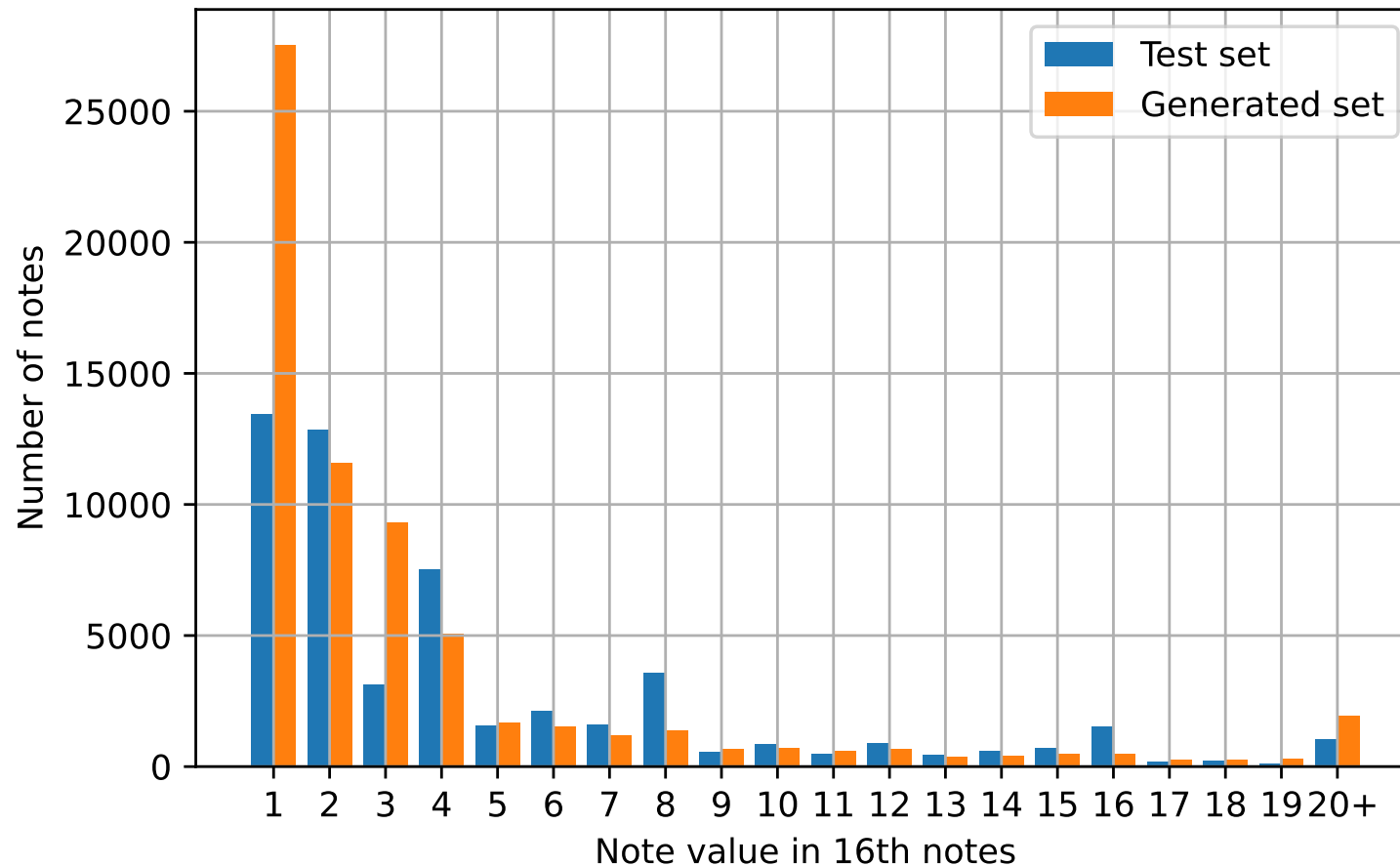
Note onsets can be assigned to a **metrical level**.



In the generated set there are **more onsets on uneven 16<sup>th</sup> notes** than on even ones.



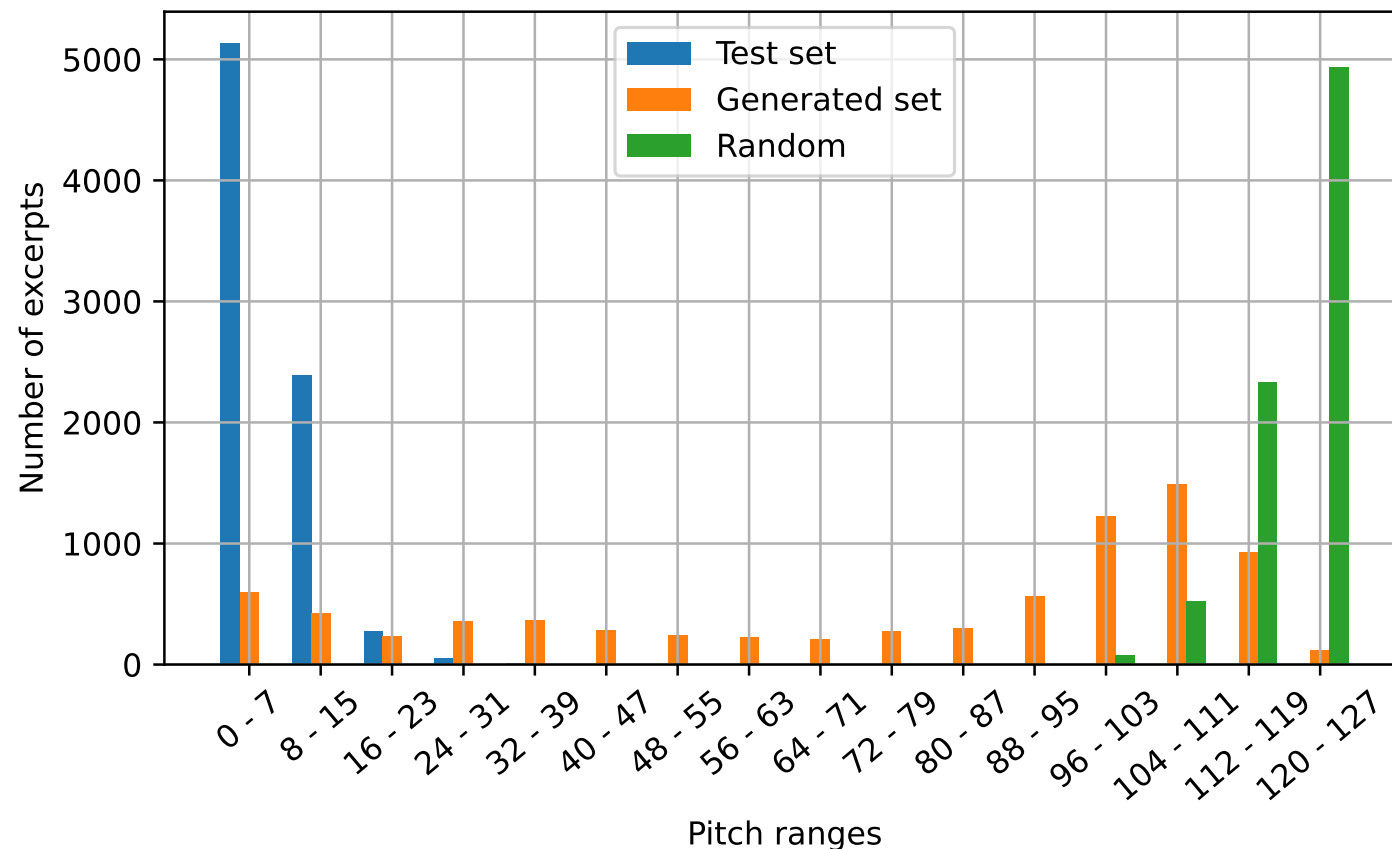
Sequence Nr. 30 (Figure 6.4)



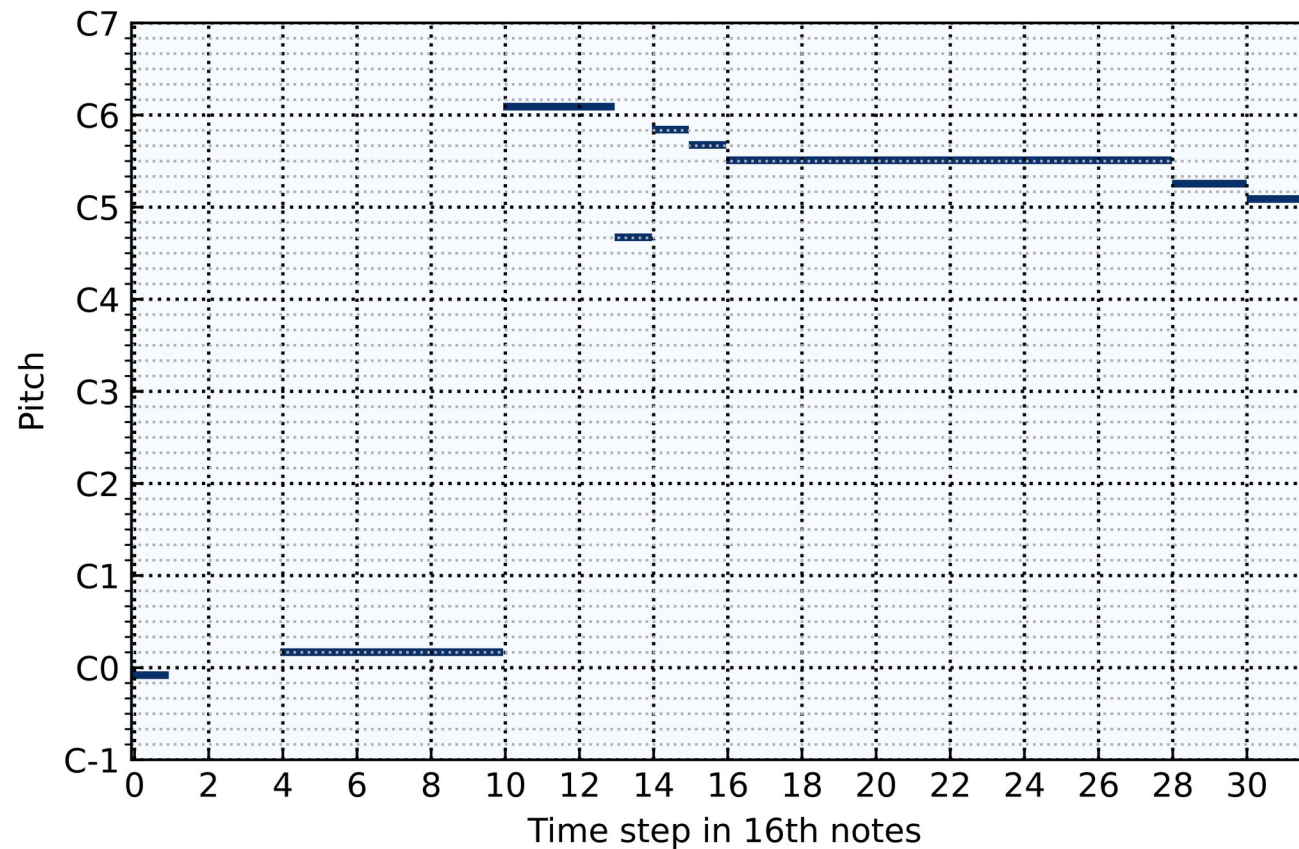
Peaks in the note length were not copied.

# Results:

# Melodic Features

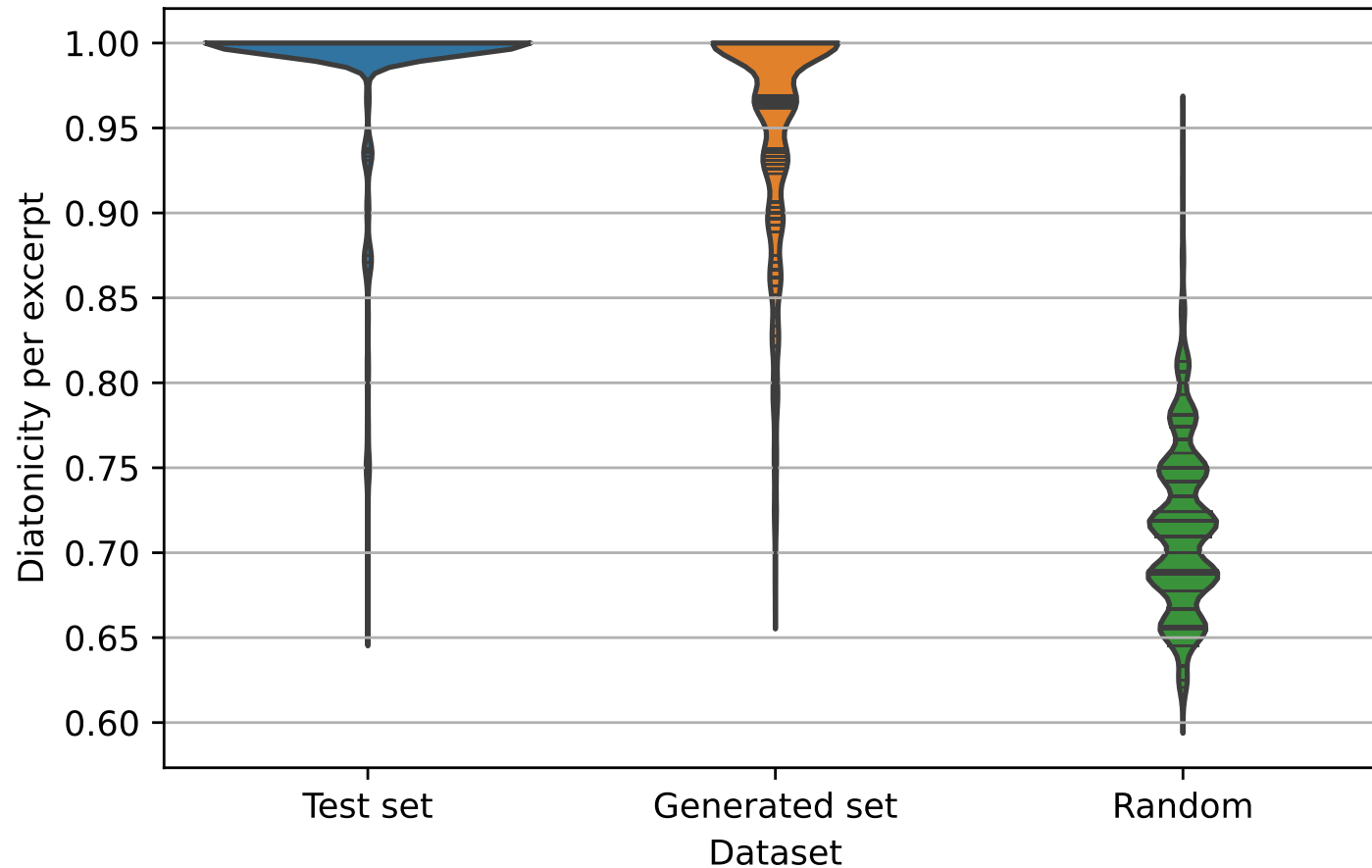


There are **high pitch jumps** in the generated excerpts.

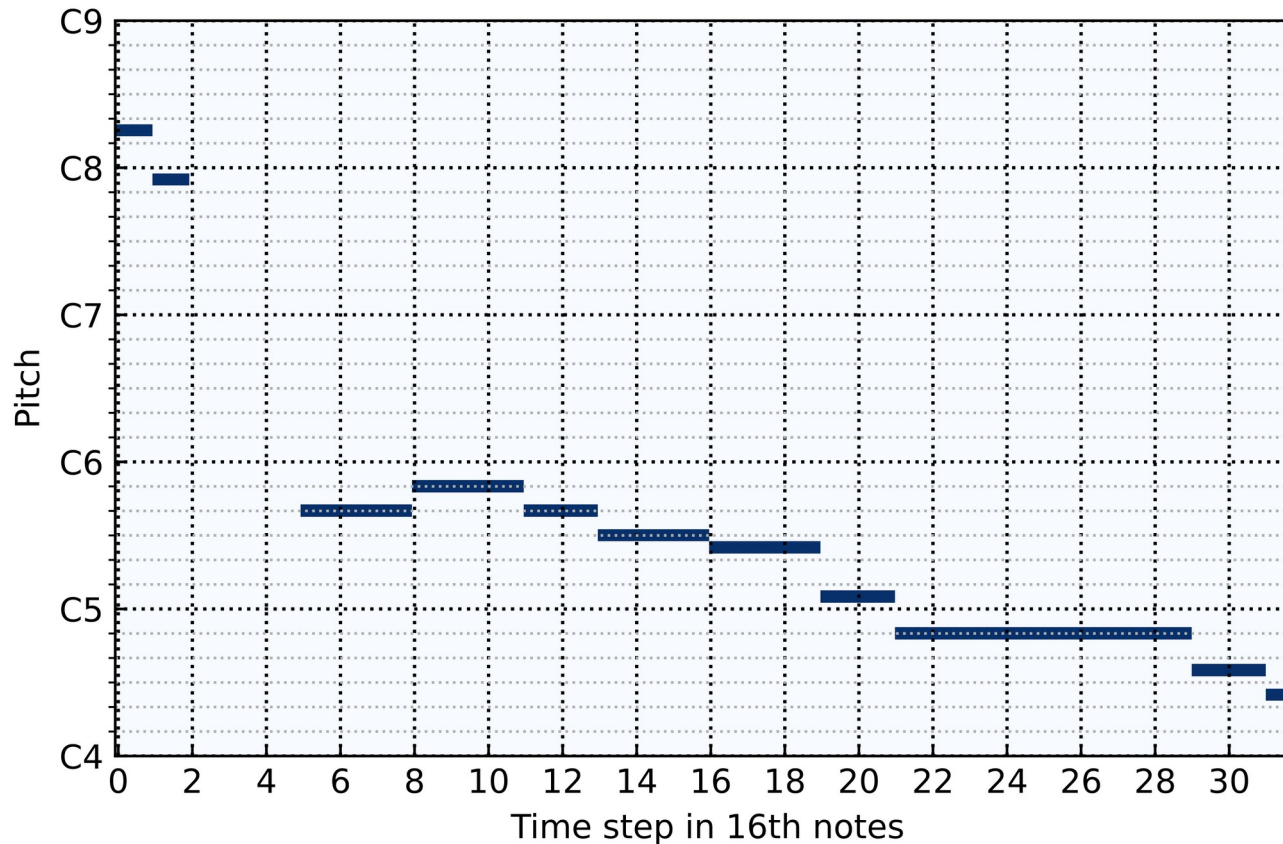


Sequence Nr. 24 (Figure 6.7)



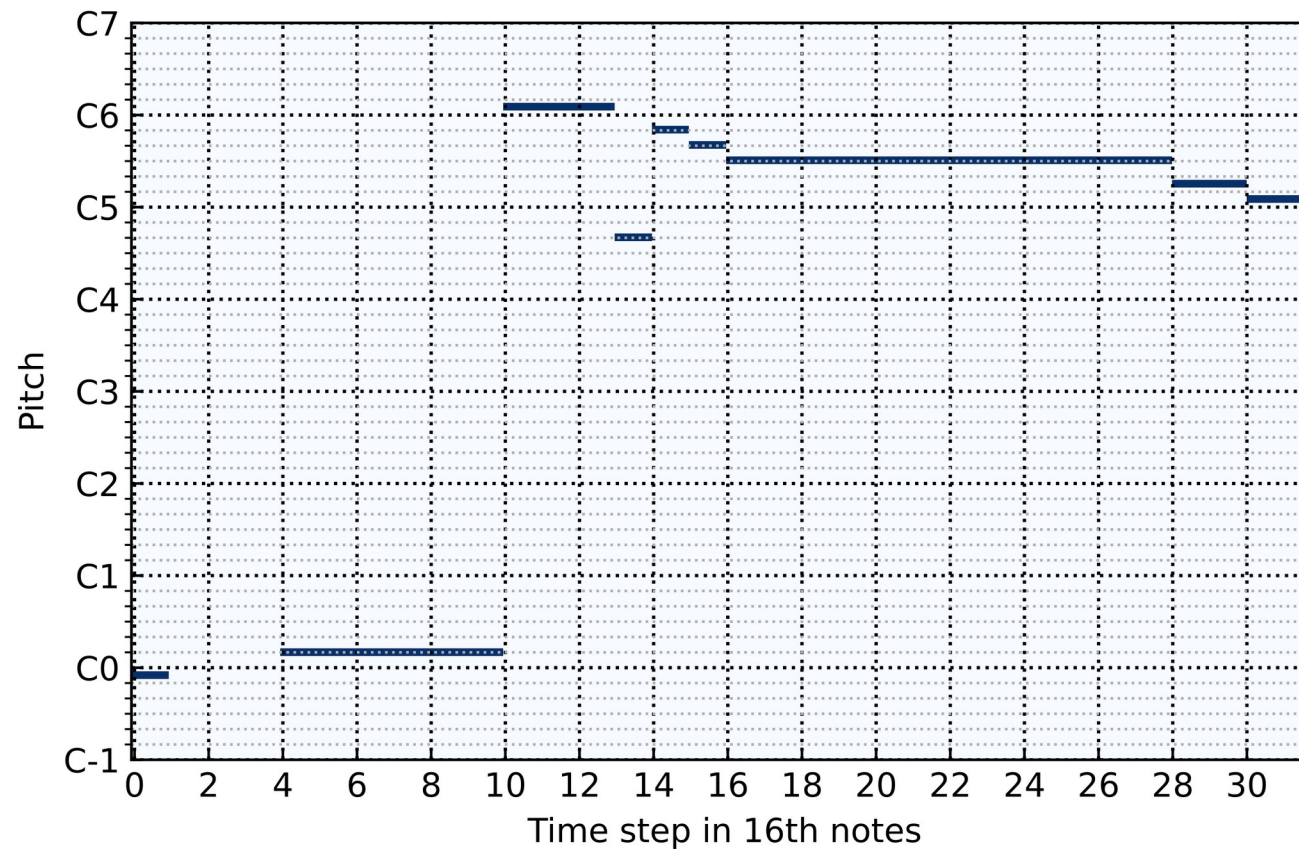


Generated excerpts are **mostly diatonic**, but there are odd notes.

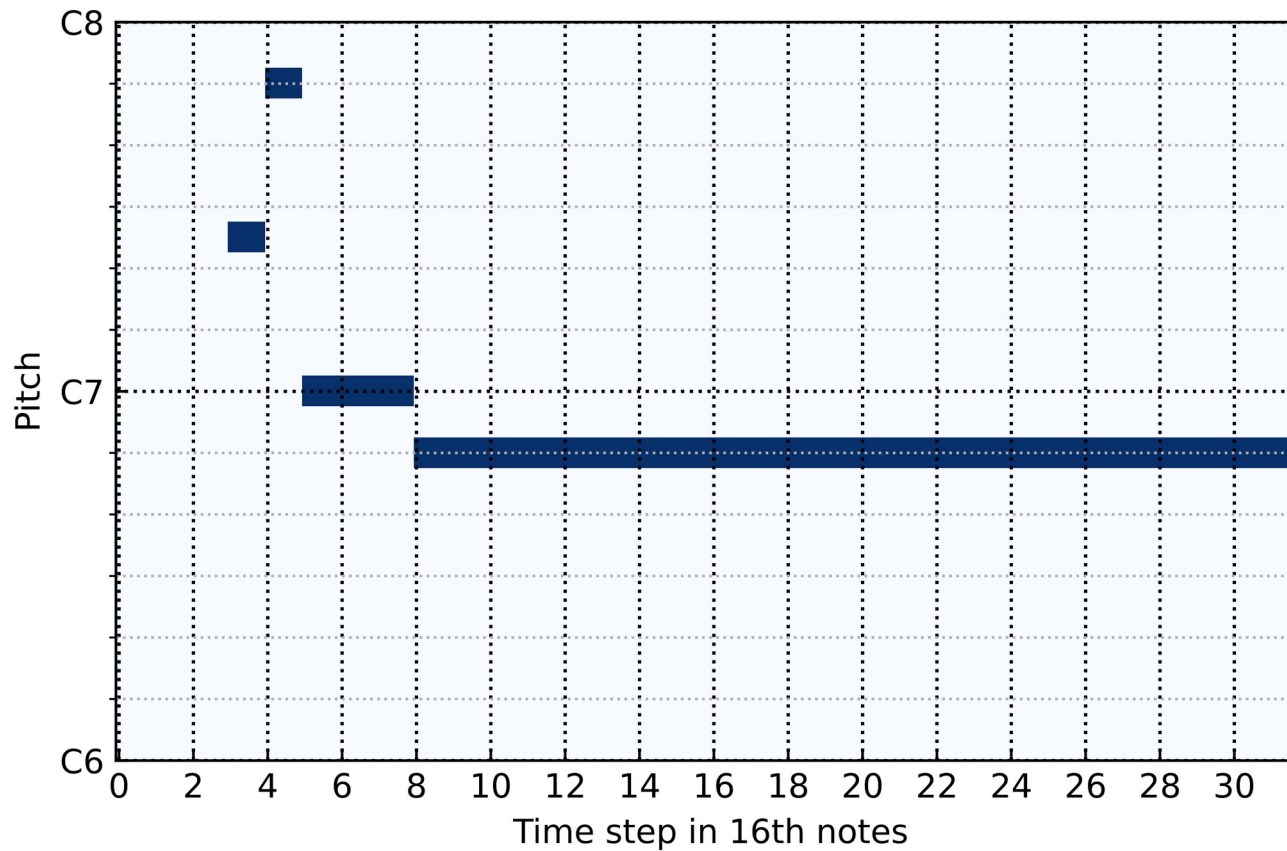


Sequence Nr. 13 (Figure 6.11)

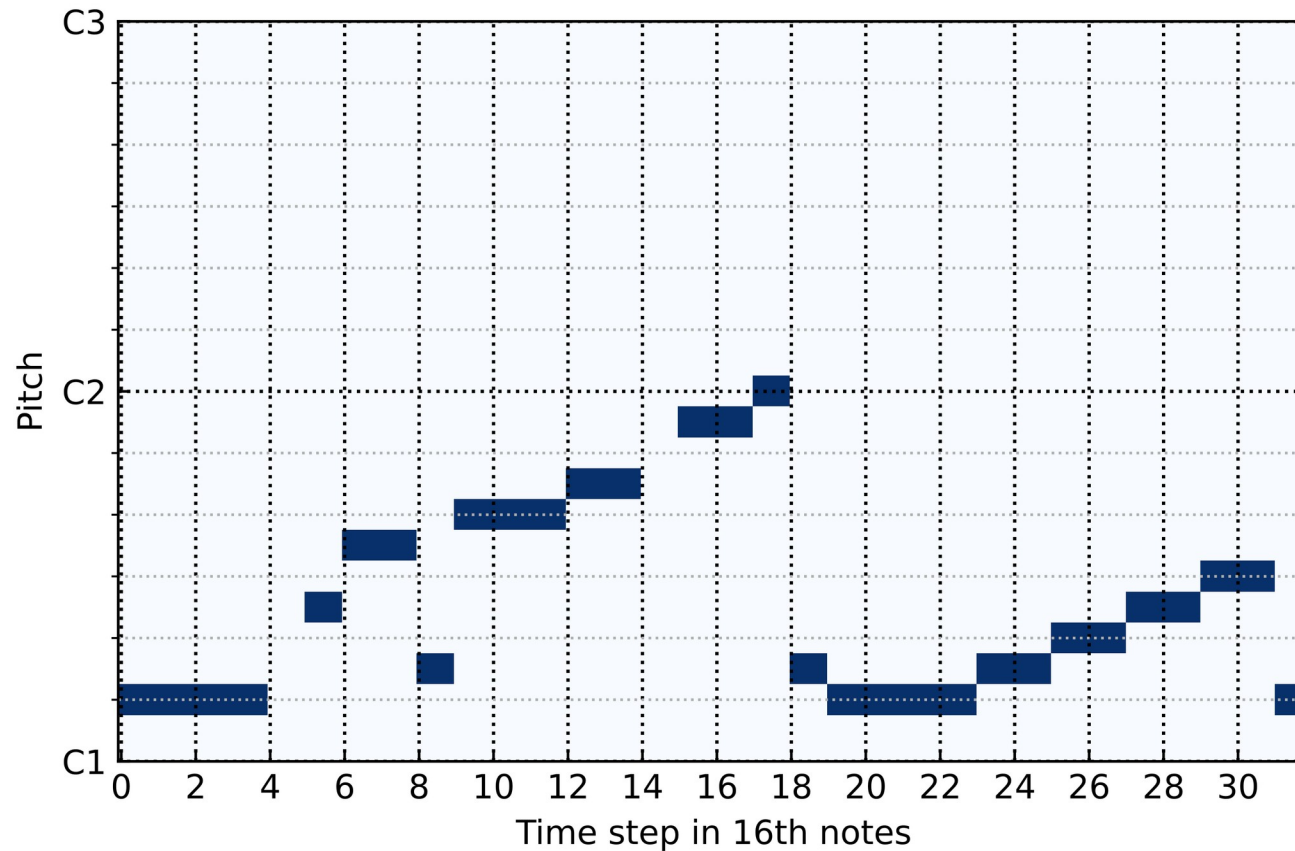
# Results: Qualitative Evaluation



Sequence Nr. 24 (Figure 6.7)



Sequence Nr. 50 (Figure 6.12)



Sequence Nr. 34 (Figure 6.15)

# Conclusion

- state of the art has been reviewed
- re-implementation
  - MusicVAE's flat variant implemented
  - excerpts generated
- generated excerpts were evaluated



- more **onsets on uneven 16<sup>th</sup> notes**
- single **non-diatonic notes & pitch jumps**
- mostly musically coherent & **pleasant-sounding**

- recreate training **dataset**
- adjust training **procedure**
- **extend** model
  - hierarchical, polyphonic, conditioned



TECHNISCHE UNIVERSITÄT  
ILMENAU

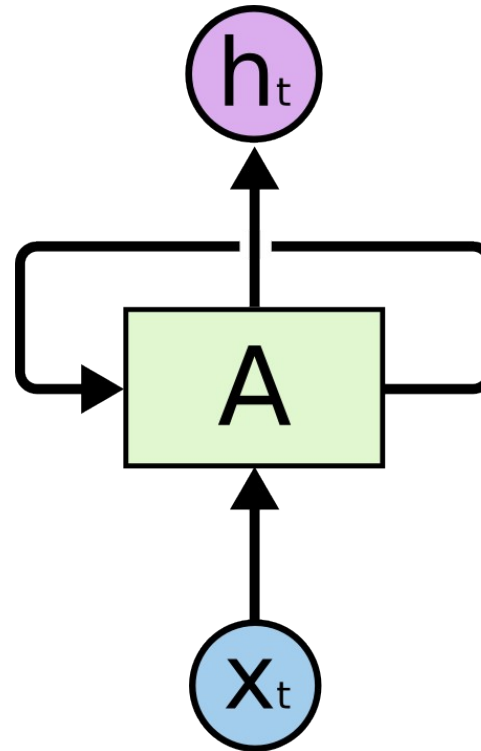


**Fraunhofer**  
IDMT

**Thank you!**

# Recurrent Neural Networks

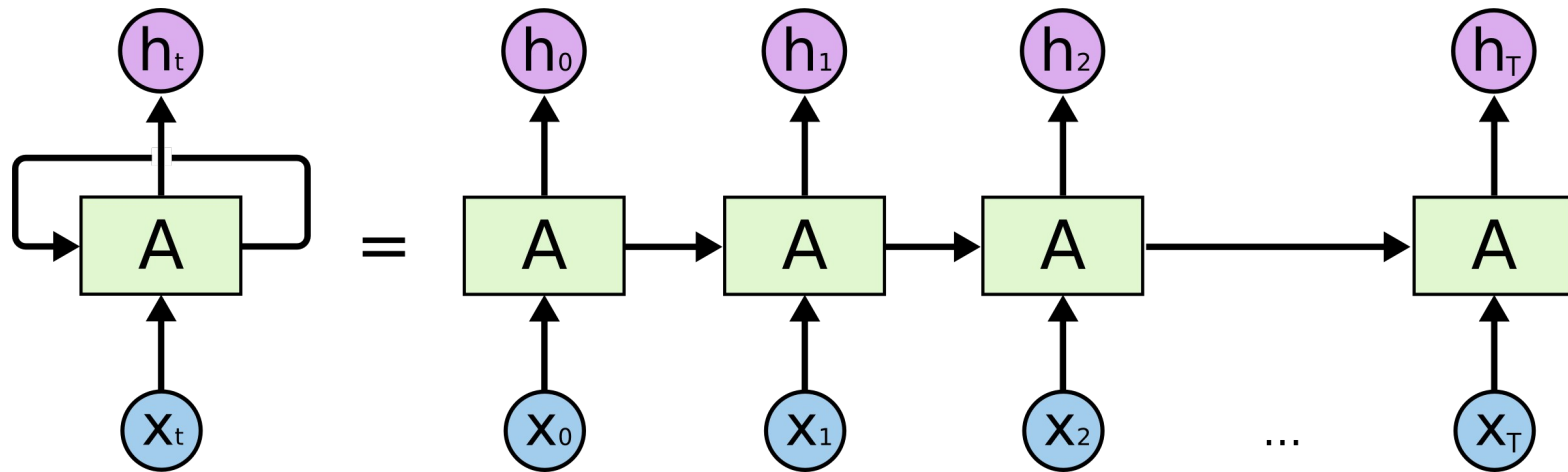
RNNs are neural networks with feedback.



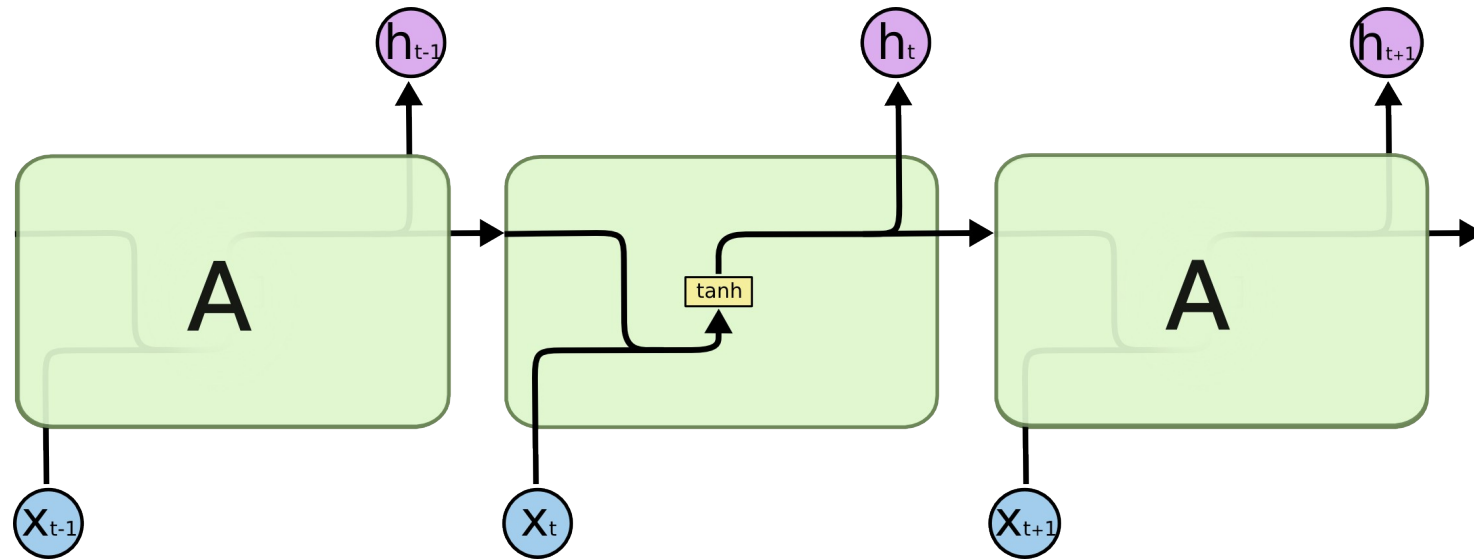
Picture retrieved July 8, 2023 from

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/img/RNN-rolled.png>

RNNs can be represented unrolled.



Picture retrieved and changed July 8, 2023 from  
<https://colah.github.io/posts/2015-08-Understanding-LSTMs/img/RNN-unrolled.png>

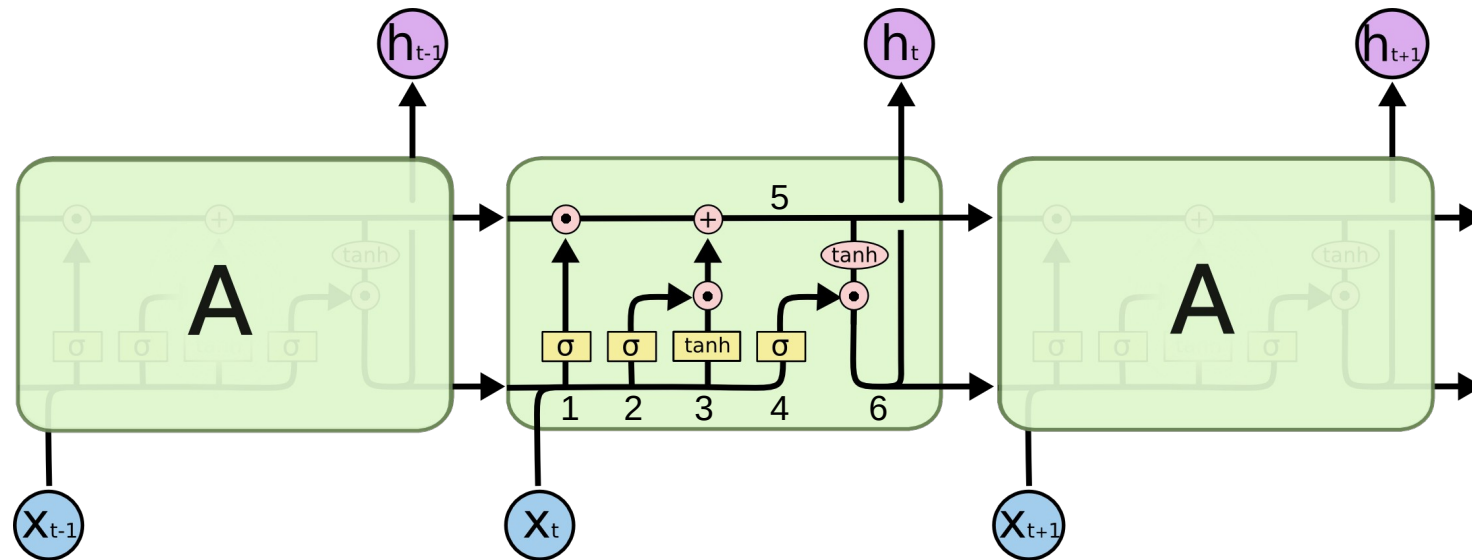


$$h_t = \tanh(W_i x_t + b_i + W_h h_{t-1} + b_h)$$

Standart RNNs have some drawbacks.

Picture retrieved July 8, 2023 from

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-SimpleRNN.png>



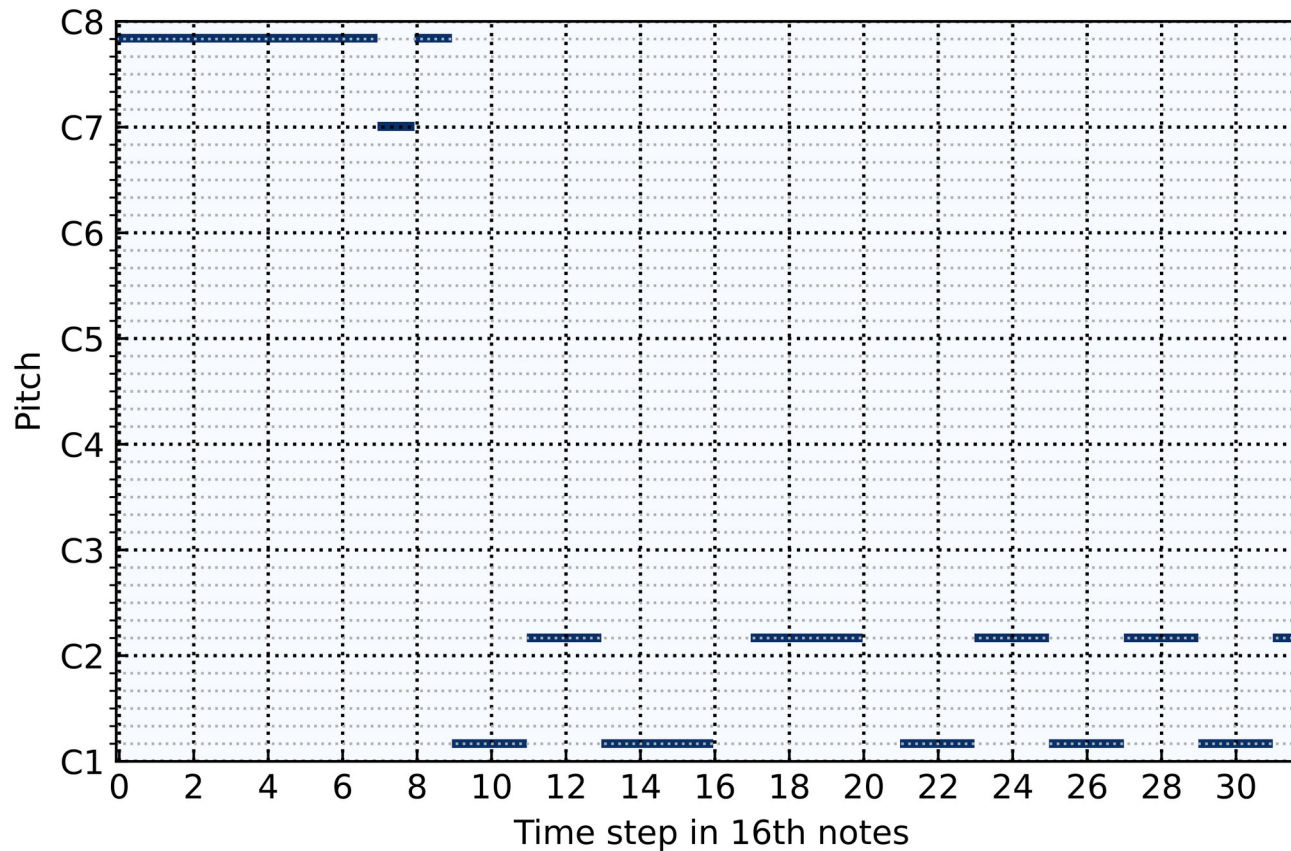
$$\begin{aligned}
 1: i_t &= \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{t-1} + b_{hi}) & \sigma(x) &= \text{sigmoid}(x) \\
 2: f_t &= \sigma(W_{if}x_t + b_{if} + W_{hf}h_{t-1} + b_{hf}) \\
 3: g_t &= \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{t-1} + b_{hg}) & 5: c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\
 4: o_t &= \sigma(W_{io}x_t + b_{io} + W_{ho}h_{t-1} + b_{ho}) & 6: h_t &= o_t \odot \tanh(c_t)
 \end{aligned}$$

Picture retrieved and changed July 8, 2023 from

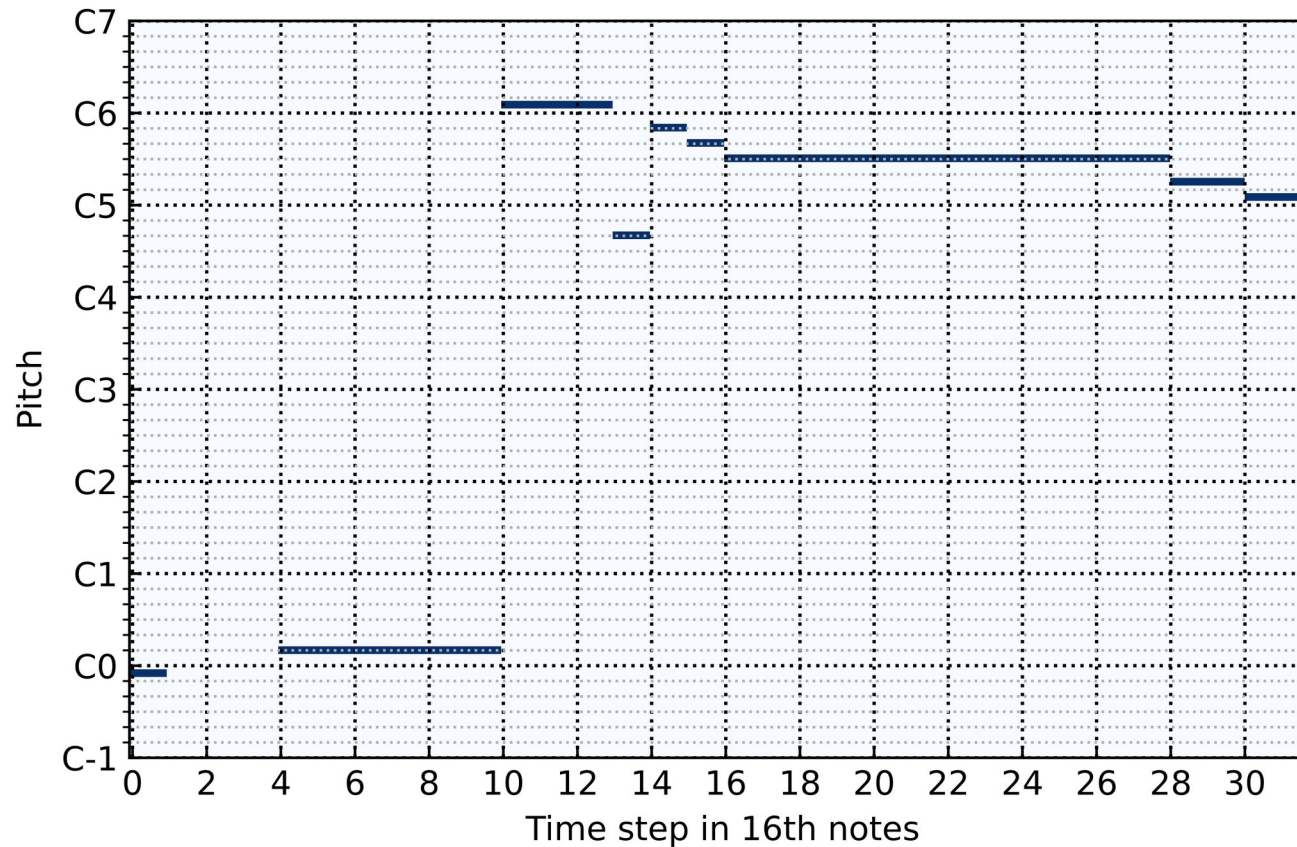
<https://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-chain.png>



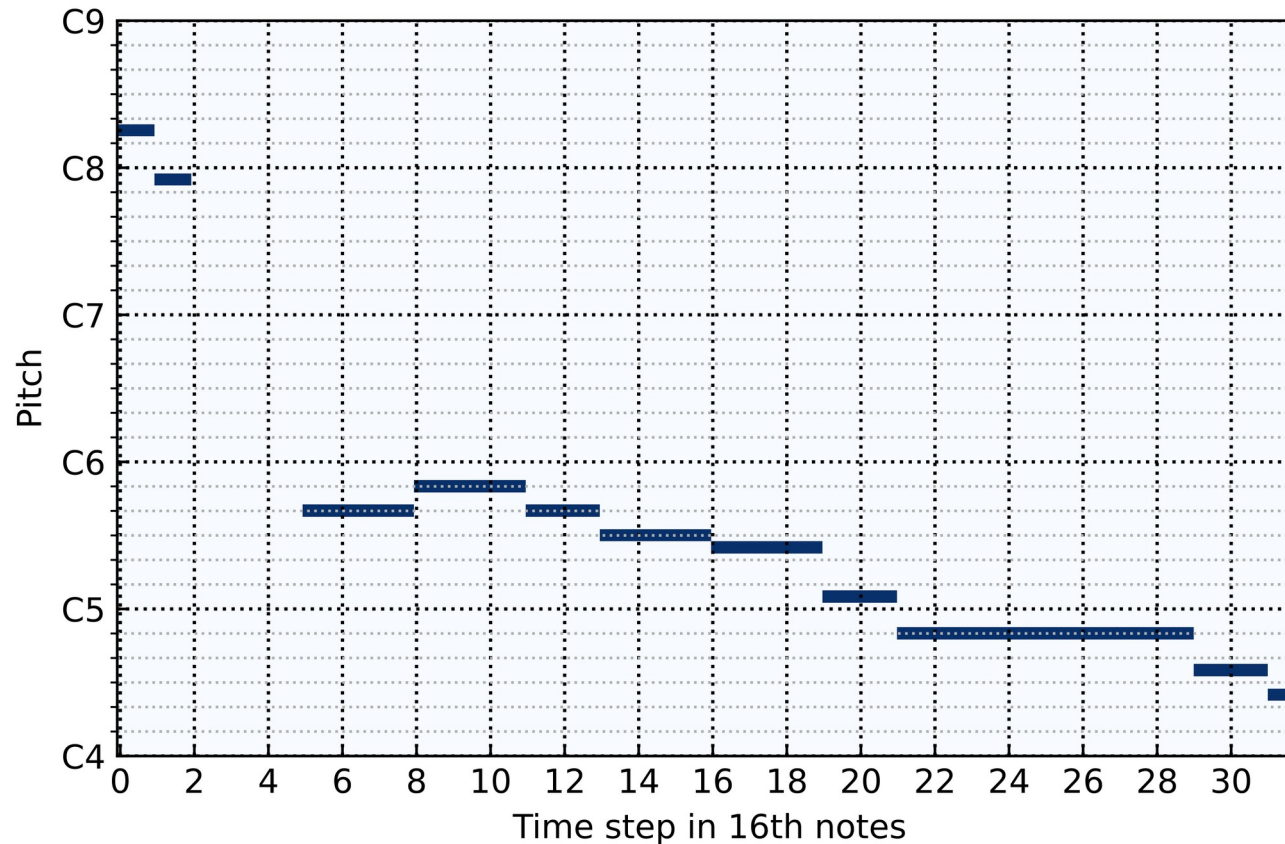
# Generated Excerpts



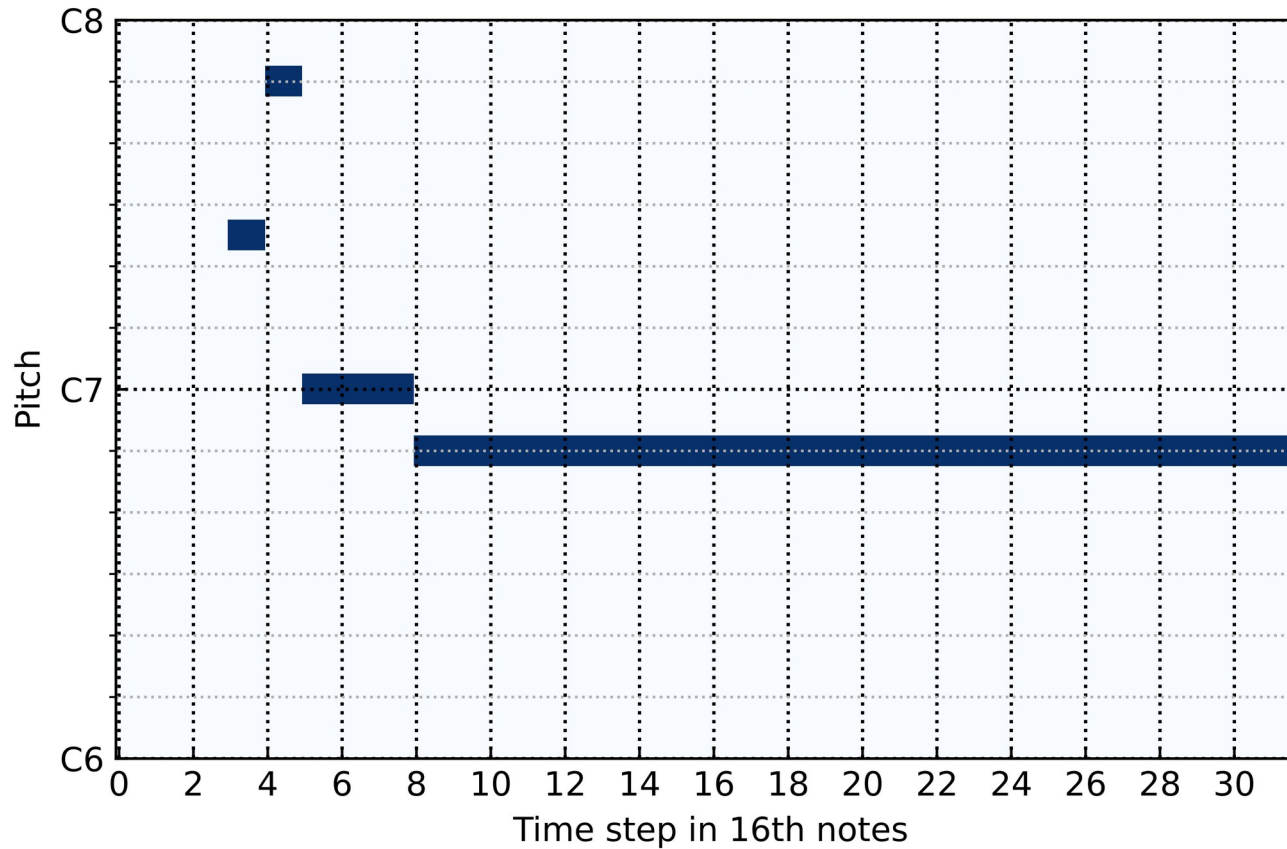
Sequence Nr. 30 (Figure 6.4)



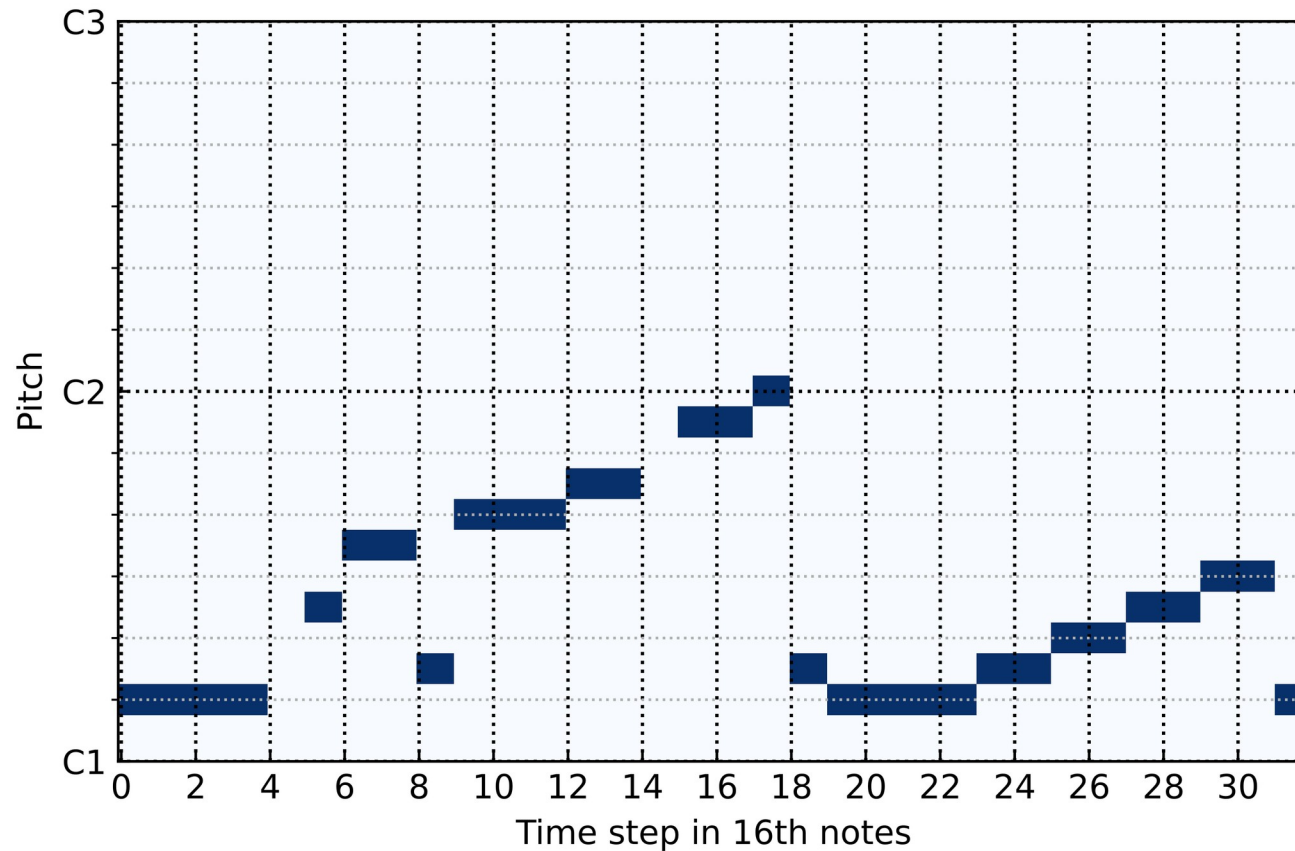
Sequence Nr. 24 (Figure 6.7)



Sequence Nr. 13 (Figure 6.11)



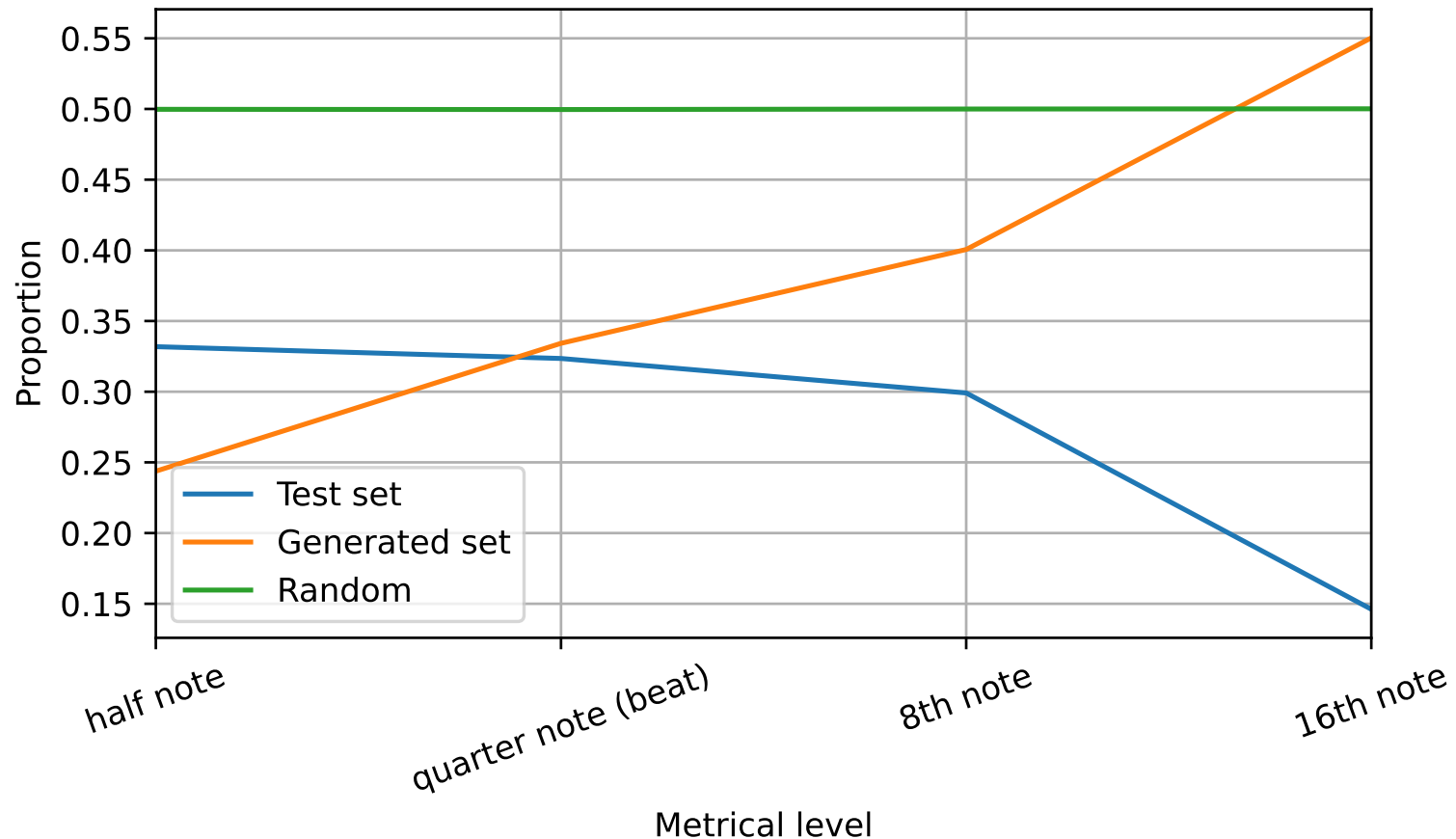
Sequence Nr. 50 (Figure 6.12)



Sequence Nr. 34 (Figure 6.15)

# Evaluation

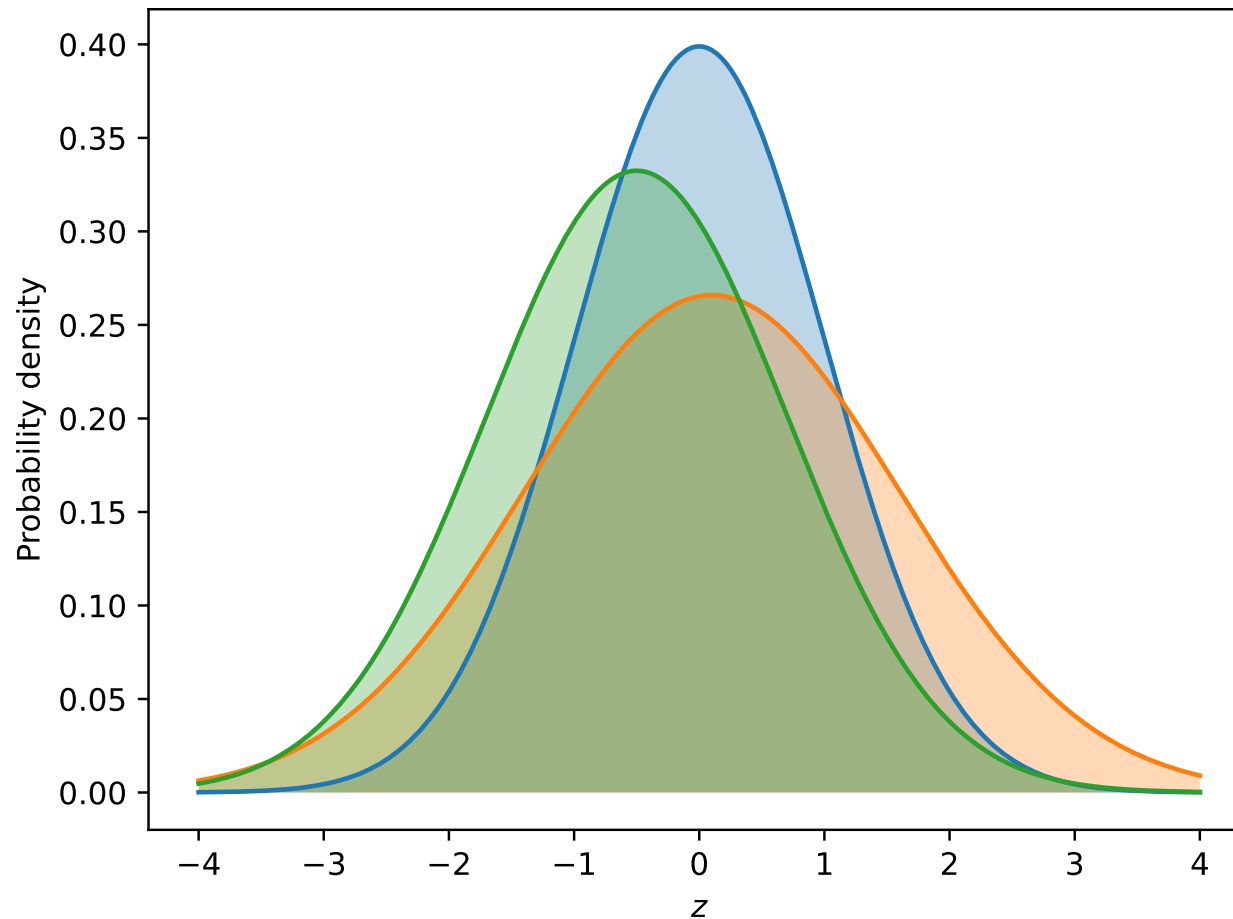




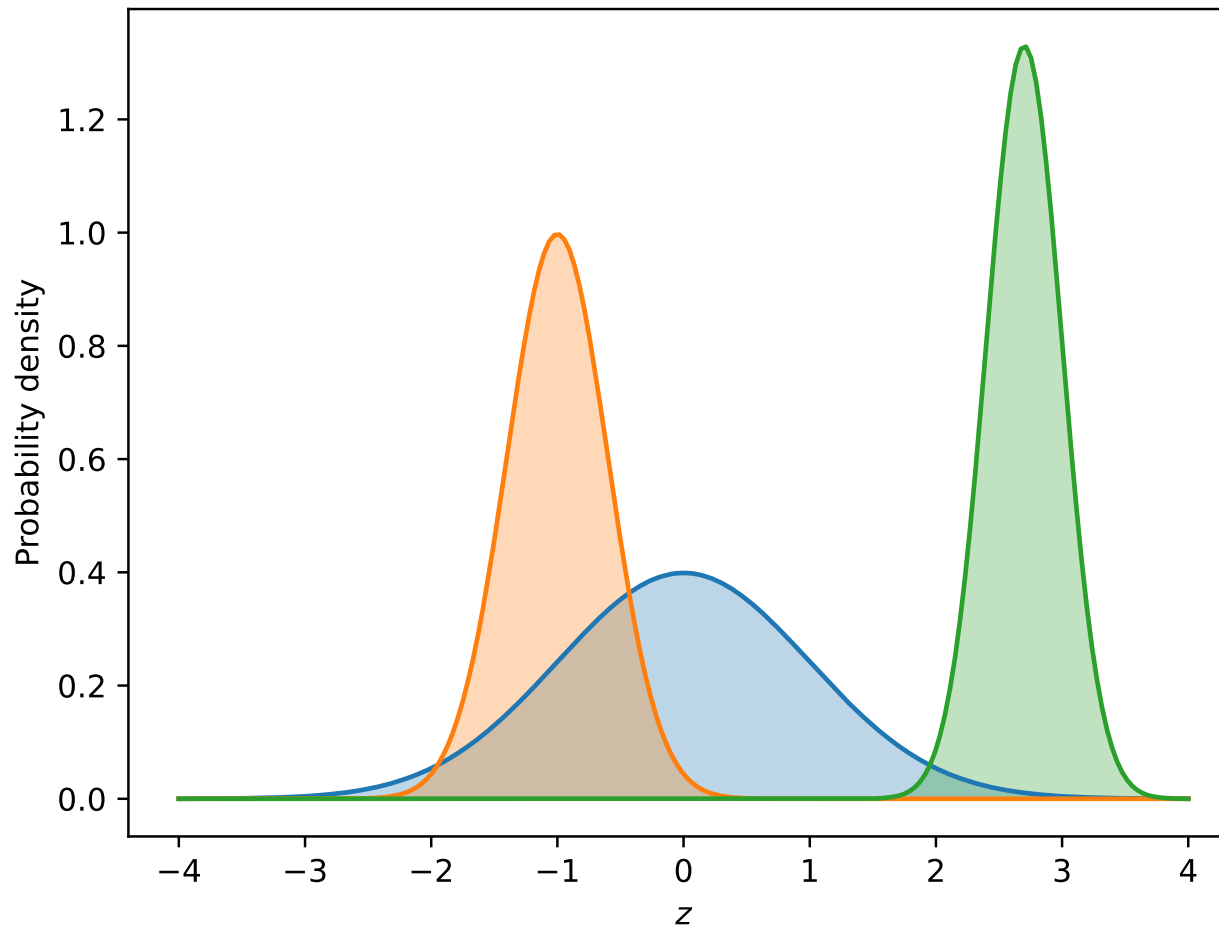
In the generated set there are **more onsets on uneven 16<sup>th</sup> notes** than on even ones.



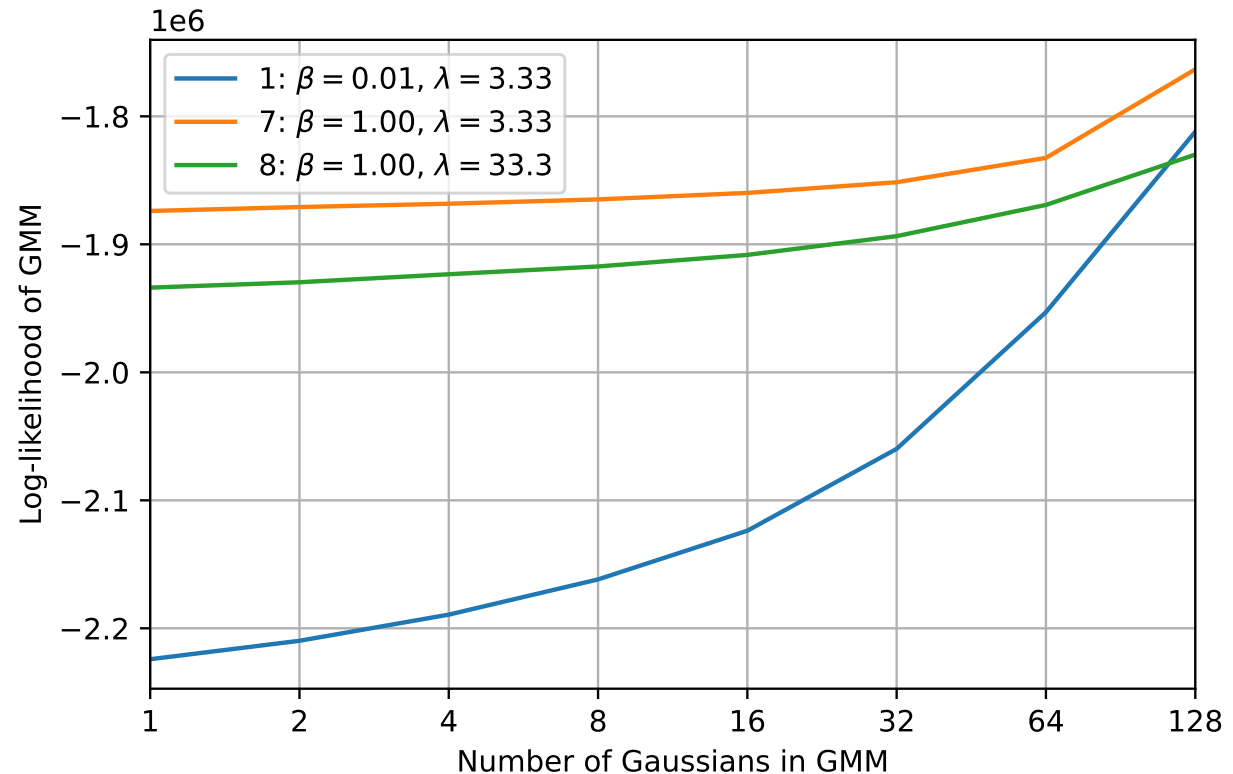
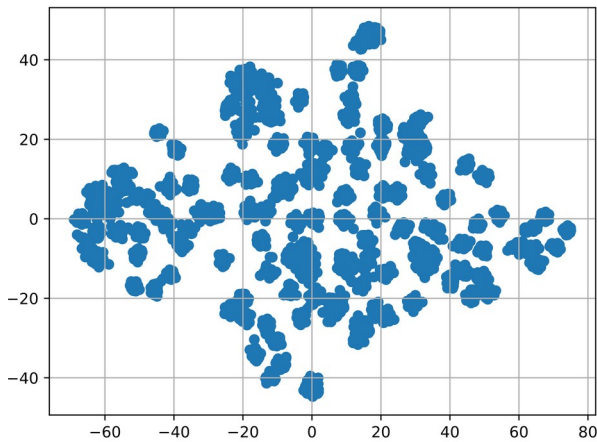
# Latent Space



well-formed



not well-formed



The structure of the latent spaces of Models 1, 7, and 8 was examined.

# Generation of Symbolic Music Based on MusicVAE

Jakob Lerch

Supervisors:

Prof. Dr. Shardt

Dr. Andrew McLeod

Dr. Jakob Abeßer

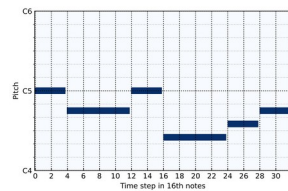
July 27, 2023



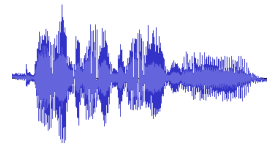
Music generation could be **supportive for composition or live performances.**

Picture retrieved July 16, 2023 from  
<https://pixabay.com/photos/music-producer-studio-actor-audio-4507819/>

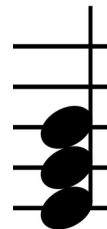
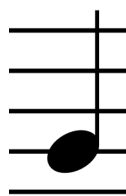
- **What is music generation?**
- **why do we care?**
  - create music for creative productions
  - support during composition or live performances



or



or



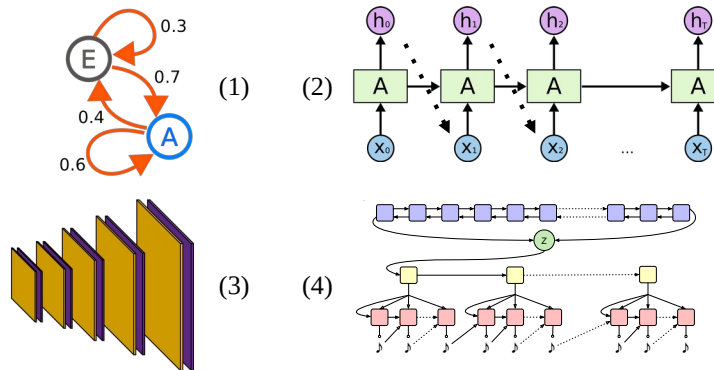
unconditioned or conditioned

- **symbolic vs. waveform**
- **monophonic vs. polyphonic**
- **conditioned vs. unconditioned**
- **Western music**
- **objectives**
  - summarize state of the art
  - re-implement MusicVAE
    - training data: two-bar monophonic sequences
  - generate excerpts and discuss their quality

- (1) review **state of the art**
- (2) **re-implement** MusicVAE
- (3) **evaluate quality** of generated excerpts



# State of the Art



(1) Picture retrieved July 17, 2023 from [https://en.wikipedia.org/wiki/Markov\\_chain#/media/File:Markovkate\\_01.svg](https://en.wikipedia.org/wiki/Markov_chain#/media/File:Markovkate_01.svg)

(2) Picture retrieved and changed July 17, 2023 from <https://colah.github.io/posts/2015-08-Understanding-LSTMs/img/RNN-unrolled.png>

(3) Picture retrieved and changed July 17, 2023 from [https://clinicadl.readthedocs.io/en/latest/images/transfer\\_learning.png](https://clinicadl.readthedocs.io/en/latest/images/transfer_learning.png)

(4) Picture retrieved and changed July 17, 2023 from [https://magenta.tensorflow.org/assets/music\\_vae/architecture.png](https://magenta.tensorflow.org/assets/music_vae/architecture.png)

## • approaches

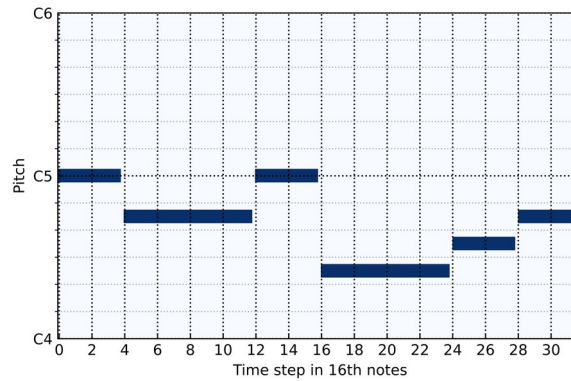
- historically: Illiac Suite, use markov models
- just RNNs, LSTMs
  - generate **autoregressively**
  - **difference to Markov models:** MM looks back limited amount of time, RNNs potentially unlimited time
- MidiNet
  - generate **pianoroll** using **CNN**
- MusicVAE
  - **hierarchical**
- TransformerVAE
  - leads to **similar reconstruction acc.**

## • evaluation procedures

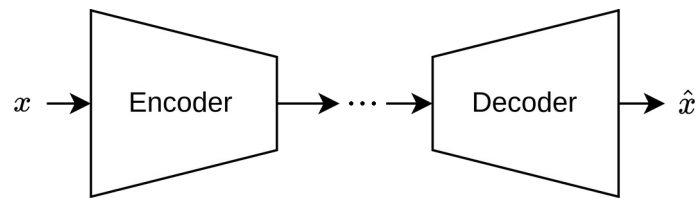
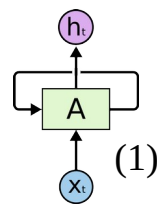
- quality?
  - listening study
  - objective measures
- originality?
  - how many notes are equal (regardless of transp.)
- deficit: only compared to other DL appr.

# MusicVAE

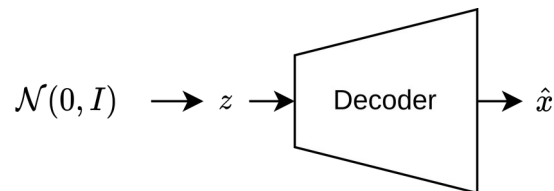
rest	{	0	1	0
pitch	{	0	0	0
		0	.	.
		1	.	0
		0	.	1
		.	.	0
		0	0	0



- **data and representation**
  - there are symbols for each pitch plus a symbol for “no pitch”

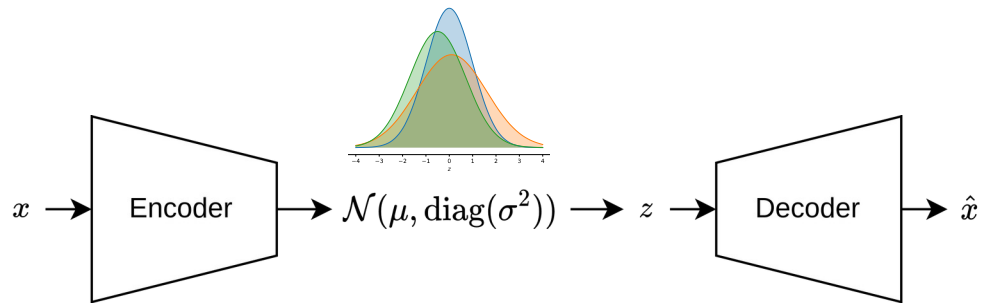


Train a model to  
generate music  
from a random  
vector.



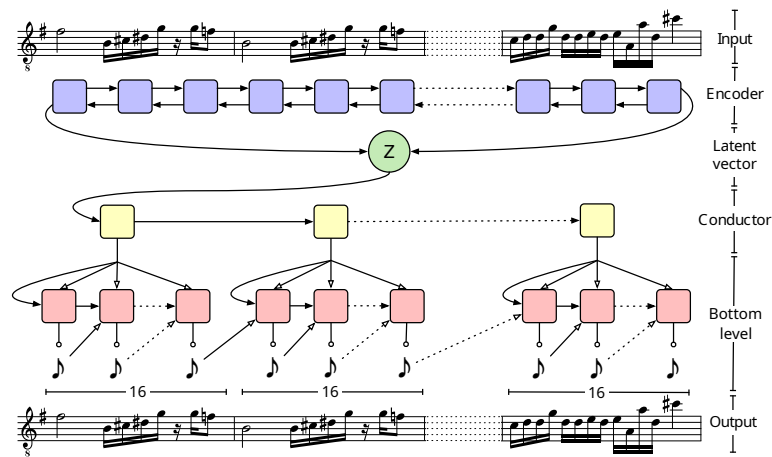
(1) Picture retrieved July 8, 2023 from  
<https://colah.github.io/posts/2015-08-Understanding-LSTMs/img/RNN-rolled.png>

- **general structure** of MusicVAE
  - music  $\rightarrow$  encoder  $\rightarrow$   $z$  with special requirements  $\rightarrow$  decoder  $\rightarrow$  music
- **music is encoded as sequence of vectors.**
  - **each vector represents one 16<sup>th</sup> note**
- first: RNNs, then VAEs



$$L(x) = \text{rec. loss} + D_{KL}(\mathcal{N}(\mu, \text{diag}(\sigma^2)) \parallel \mathcal{N}(0, I))$$

- $x \rightarrow \mu, \sigma \rightarrow \text{dist} \rightarrow z \rightarrow \hat{x}$
- **reconstruction loss**
  - cross entropy
- **KL divergence**



## Hierarchical

- **Encoder:**
  - two stacked BLSTMs
- **Decoder:**
  - hierarchical
  - conductor LSTM
    - one vector per bar
  - bottom-level LSTM
    - bar from conductor vector

## Flat

- **Encoder:**
  - simple BLSTM
- **Decoder:**
  - multi-layer LSTM

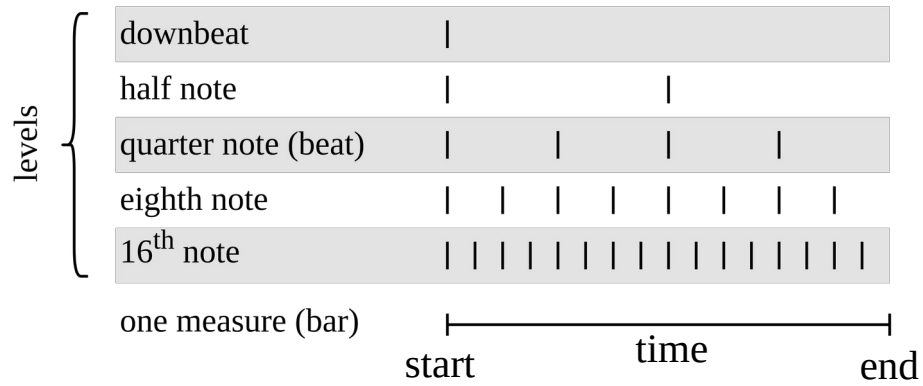
$$L(x) = \text{rec. loss} + \underline{\beta} \max[ D_{KL}, \underline{\lambda} ]$$

- **$\beta$  and  $\lambda$  chosen after grid search**
  - $\beta = 1; \lambda = 33.3$
- optimizer
  - Adam
  - learning rate (LR) =  $10^{-3}$
- batch size = 64
- weight decay
  - $L_2$  regularization with weight  $10^{-6}$
- LR scheduling
  - customized variant of ReduceLROnPlateau
- early stopping was used

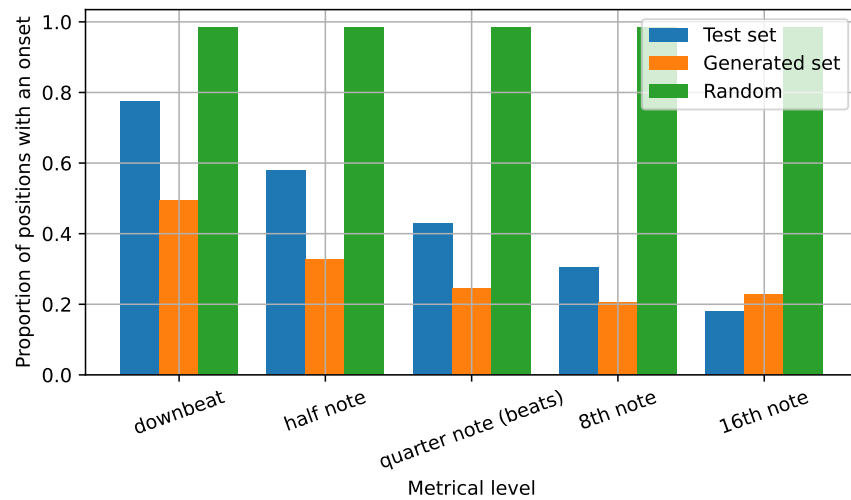
- **loss function**
  - $\beta, \lambda$  constant
- **training techniques**
  - Adam, 1e-3
  - batch size 64
  - L2 wd: 1e-6
  - customized ReduceLROnPlateau
  - early stopping



# **Results: Rhythmic Features**

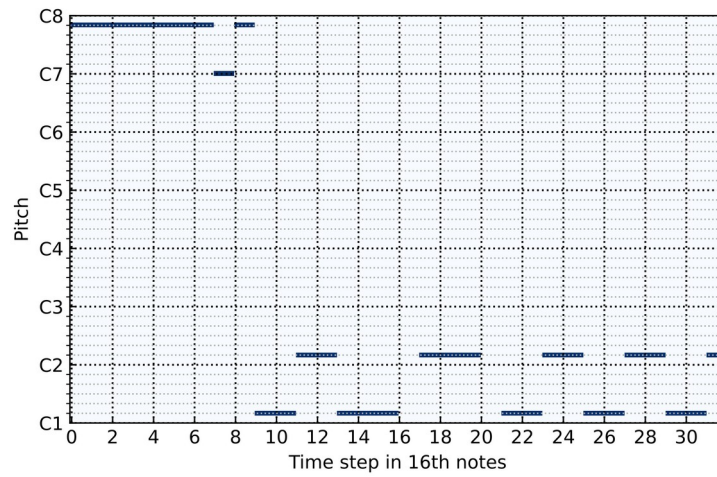


Note onsets can be assigned to a **metrical level**.

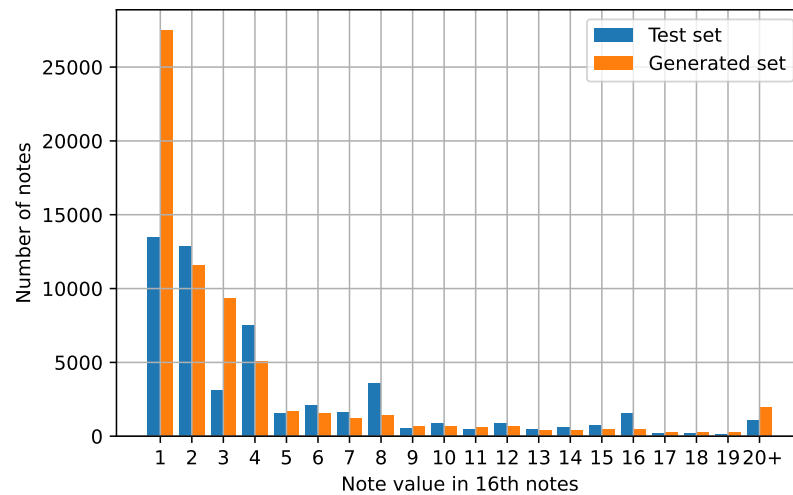


In the generated set there are **more onsets on uneven 16<sup>th</sup> notes** than on even ones.

- onset proportion
  - ...
- note length, avg note length
  - ...



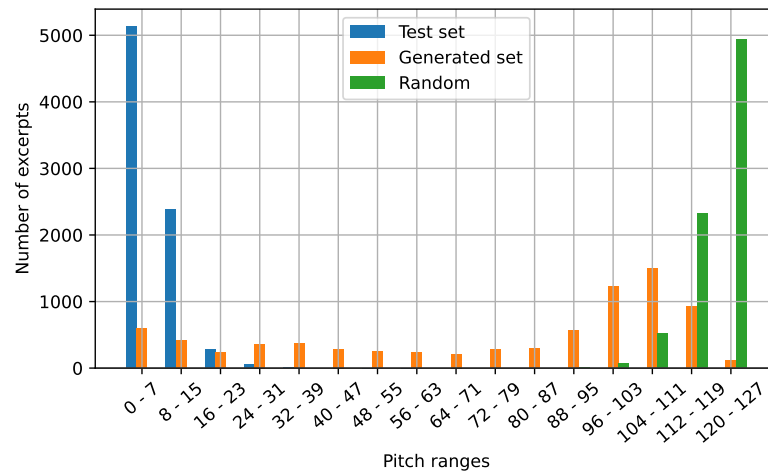
Sequence Nr. 30 (Figure 6.4)



Peaks in the note length were not copied.

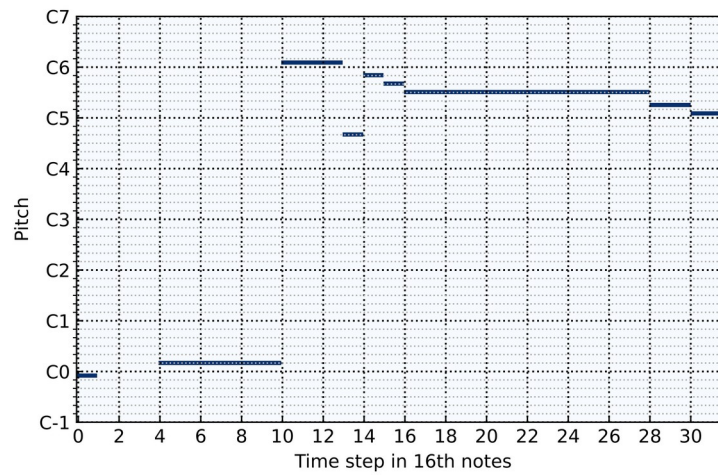
- onset proportion
  - ...
- note length, avg note length
  - ...

# **Results: Melodic Features**



There are **high pitch jumps** in the generated excerpts.

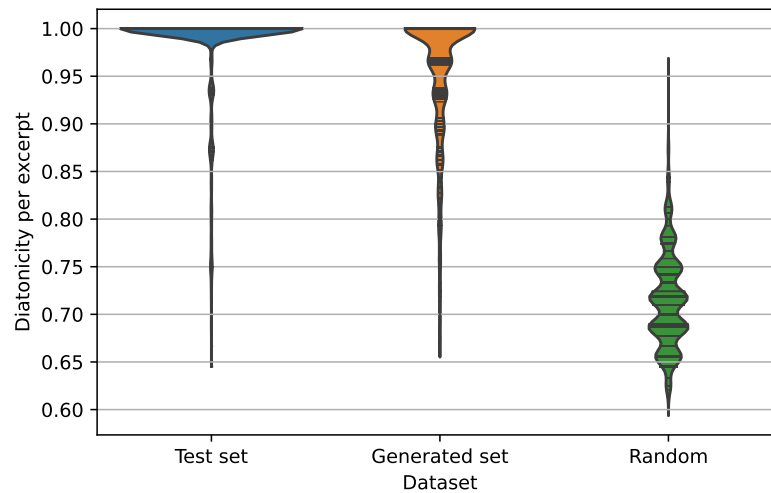
- pitch range
- ...



Sequence Nr. 24 (Figure 6.7)

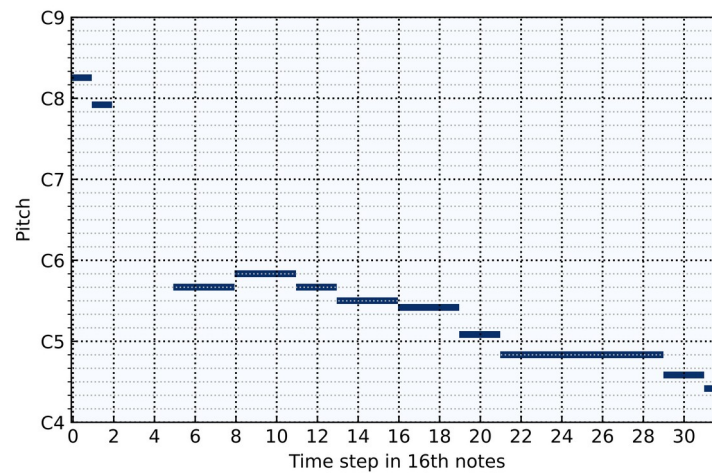


## Results: Melodic Features



Generated excerpts are **mostly diatonic**, but there are odd notes.

- diatonicity
  - how much it stays in one key

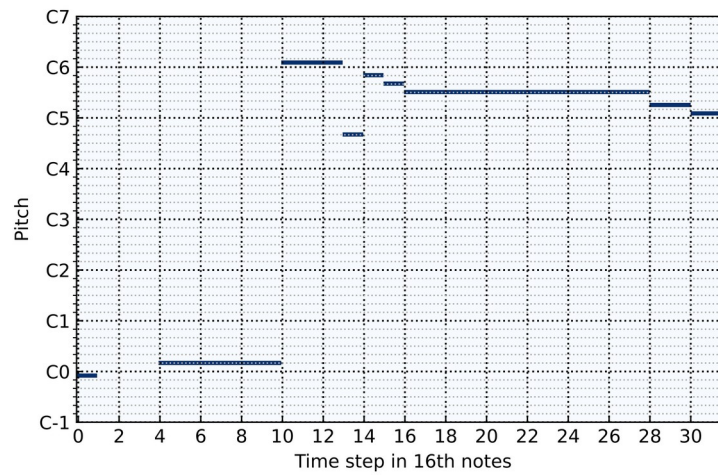


Sequence Nr. 13 (Figure 6.11)

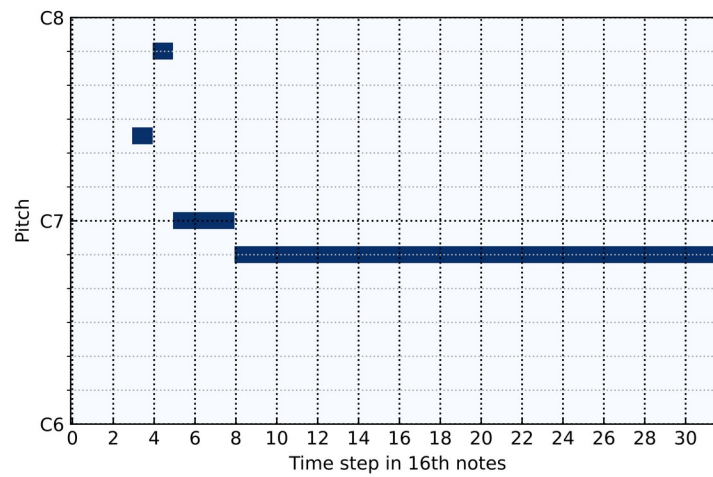
# Results: Qualitative Evaluation

23 / 30

- high pitch jumps
- rhythmic differences where not noticed
- ignoring few odd notes, still pleasant and coherent

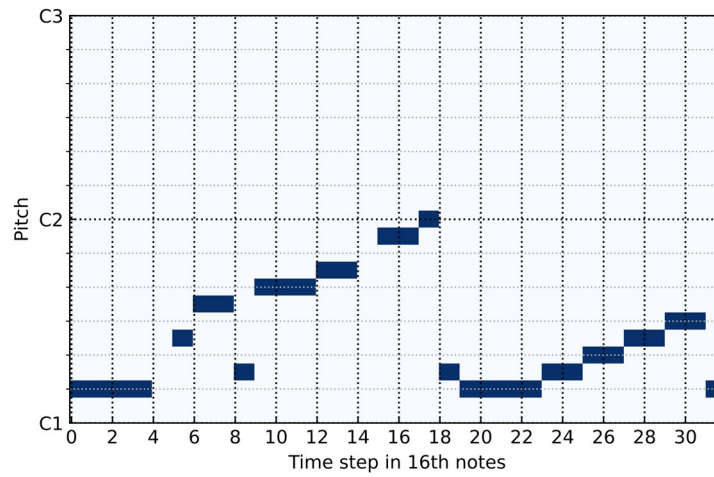


Sequence Nr. 24 (Figure 6.7)



Sequence Nr. 50 (Figure 6.12)

- chimes



Sequence Nr. 34 (Figure 6.15)

- rhythmically interesting

# Conclusion

- state of the art has been reviewed
- re-implementation
  - MusicVAE's flat variant implemented
  - excerpts generated
- generated excerpts were evaluated

- What has been done
  - MusicVAE re-implemented, excerpts generated & evaluated
  - generated excerpts similar to training data, but still differences
  - apart from these, generated melodies sounded pleasant
- Future work
  - larger dataset, analysis of it
  - regularization techniques to reduce overfitting
  - sudden jumps in the loss → investigate reason and improve training
  - improve hierarchical model
  - condition the latent space →
  - evaluate latent space in a different way



- more **onsets on uneven 16<sup>th</sup> notes**
- single **non-diatonic notes & pitch jumps**
- mostly musically coherent & **pleasant-sounding**

- What has been done
  - MusicVAE re-implemented, excerpts generated & evaluated
  - generated excerpts similar to training data, but still differences
  - apart from these, generated melodies sounded pleasant
- Future work
  - larger dataset, analysis of it
  - regularization techniques to reduce overfitting
  - sudden jumps in the loss → investigate reason and improve training
  - improve hierarchical model
  - condition the latent space →
  - evaluate latent space in a different way

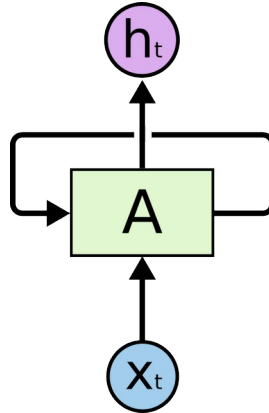
- recreate training **dataset**
- adjust training **procedure**
- **extend** model
  - hierarchical, polyphonic, conditioned

- recreate training dataset
  - larger dataset
  - analyse training dataset more thoroughly before training
  - consider alternative representations
- adjust training procedure
  - tune weight decay
  - scheduled sampling & annealing of  $\beta$
  - more thorough hyperparameter search
  - find out **reason for jumps in the loss curve** & adjust training accordingly
- extend model
  - implement hierarchical model

**Thank you!**

# Recurrent Neural Networks

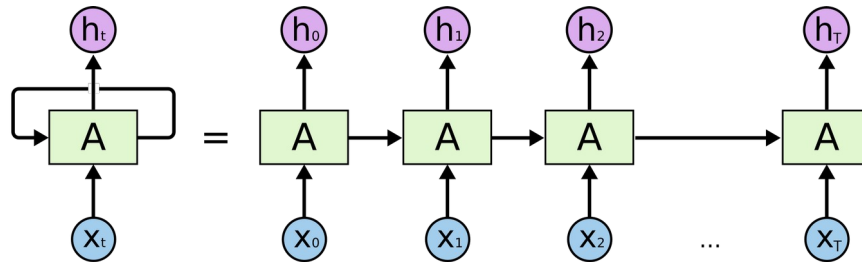
RNNs are neural networks with feedback.



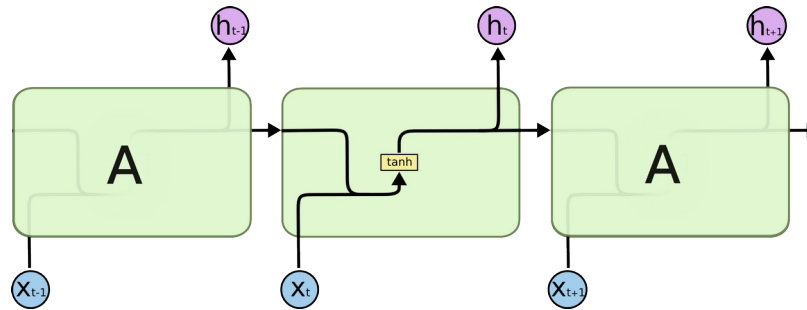
Picture retrieved July 8, 2023 from  
<https://colah.github.io/posts/2015-08-Understanding-LSTMs/img/RNN-rolled.png>

- RNNs can be used to process time-series data

RNNs can be represented unrolled.



Picture retrieved and changed July 8, 2023 from  
<https://colah.github.io/posts/2015-08-Understanding-LSTMs/img/RNN-unrolled.png>



$$h_t = \tanh(W_i x_t + b_i + W_h h_{t-1} + b_h)$$

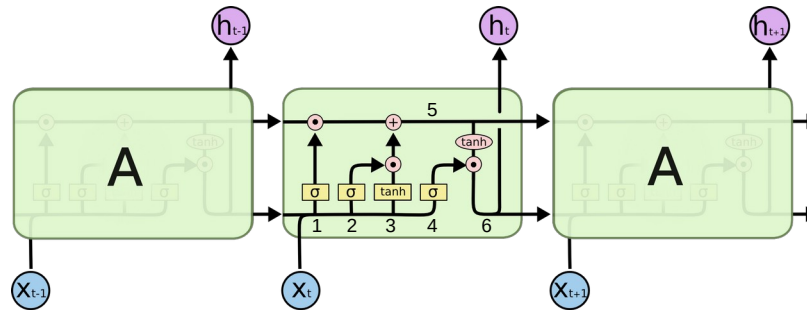
Standart RNNs have some drawbacks.

Picture retrieved July 8, 2023 from

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-SimpleRNN.png>

- **problems:**

- ...



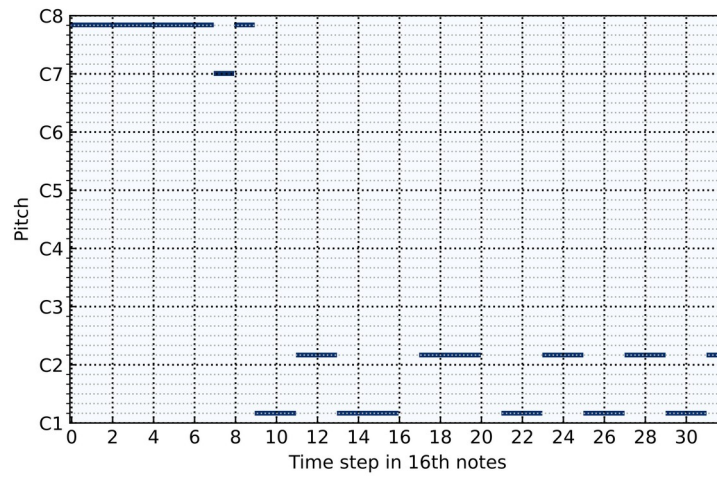
$$\begin{aligned}
 1: i_t &= \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{t-1} + b_{hi}) & \sigma(x) &= \text{sigmoid}(x) \\
 2: f_t &= \sigma(W_{if}x_t + b_{if} + W_{hf}h_{t-1} + b_{hf}) \\
 3: g_t &= \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{t-1} + b_{hg}) & 5: c_t &= f_t \odot c_{t-1} + i_t \odot g_t \\
 4: o_t &= \sigma(W_{io}x_t + b_{io} + W_{ho}h_{t-1} + b_{ho}) & 6: h_t &= o_t \odot \tanh(c_t)
 \end{aligned}$$

Picture retrieved and changed July 8, 2023 from  
<https://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-chain.png>

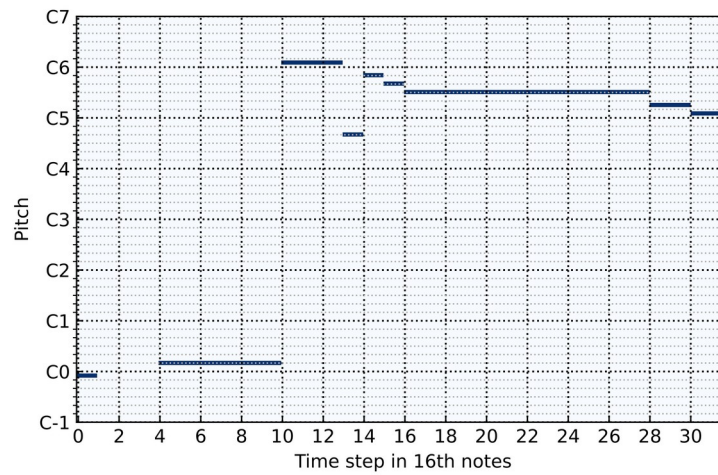
- LSTMs
  - forget gate
  - input gate
  - output gate



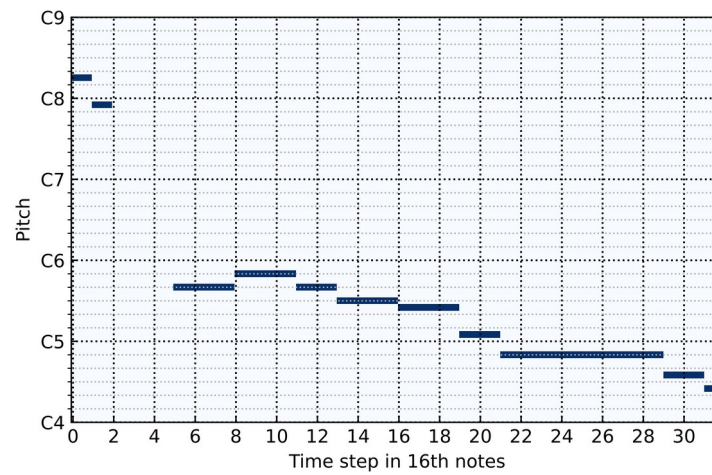
# Generated Excerpts



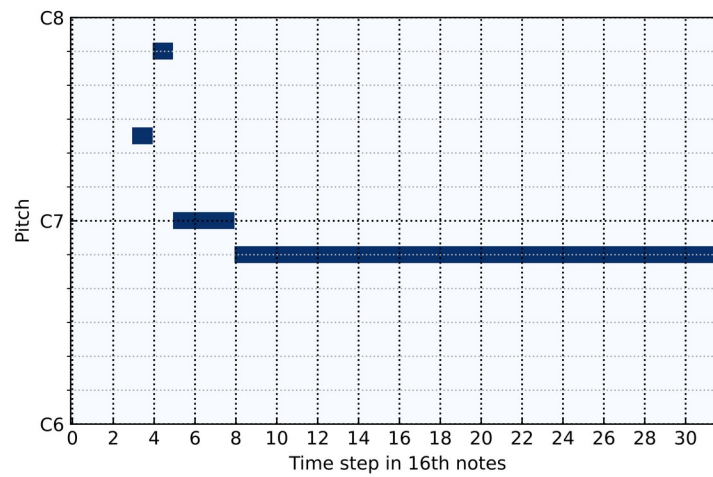
Sequence Nr. 30 (Figure 6.4)



Sequence Nr. 24 (Figure 6.7)

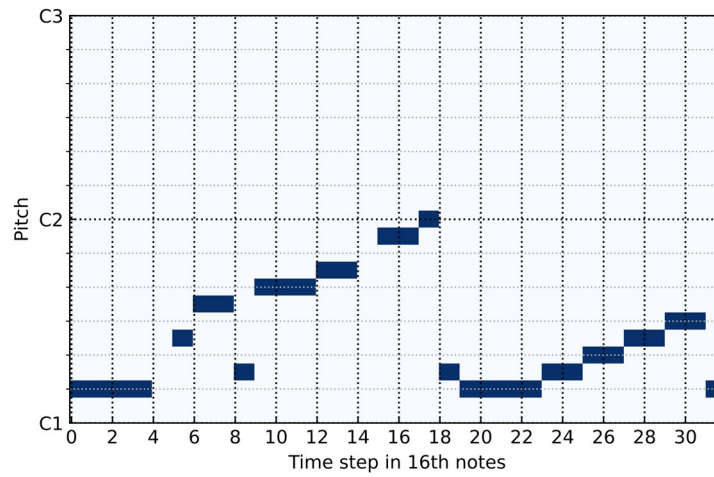


Sequence Nr. 13 (Figure 6.11)



Sequence Nr. 50 (Figure 6.12)

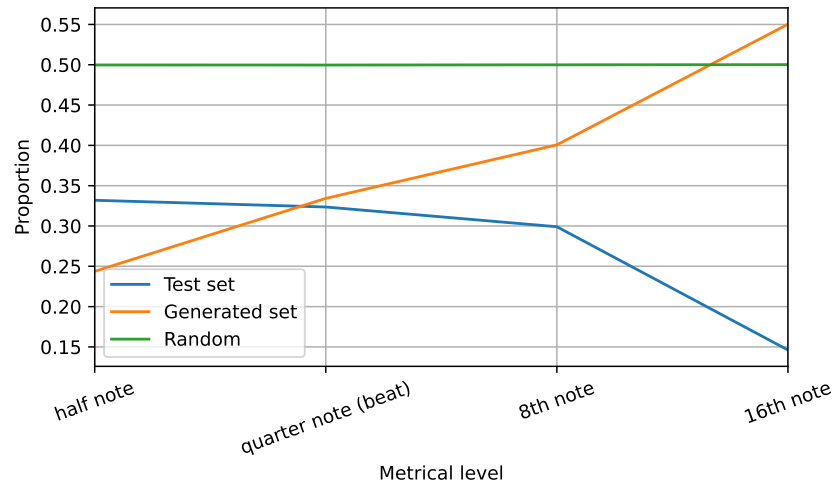
- chimes



Sequence Nr. 34 (Figure 6.15)

- rhythmically interesting

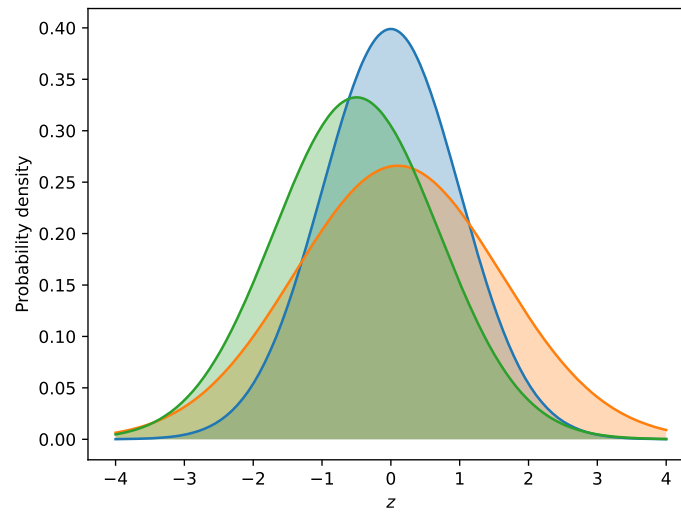
# Evaluation



In the generated set there are **more onsets on uneven 16<sup>th</sup> notes** than on even ones.

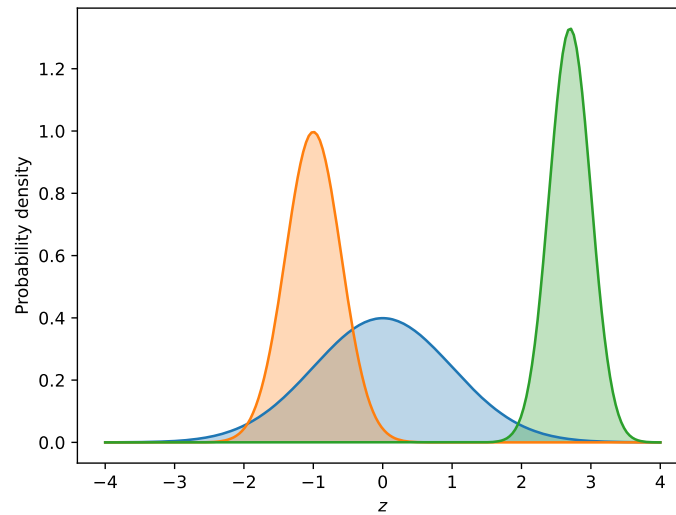


# Latent Space



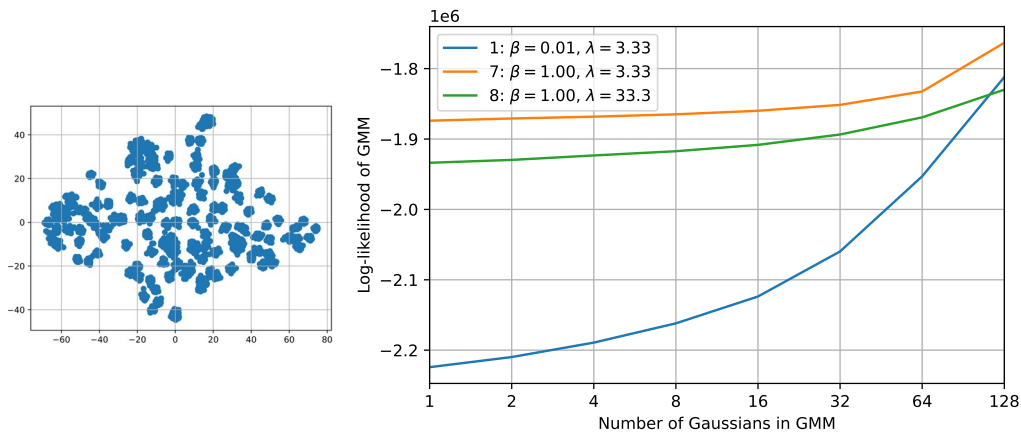
well-formed

- **Figure: desirable**



not well-formed

- **Figure: undesirable**



The structure of the latent spaces of Models 1, 7, and 8 was examined.

- grid search
  - ...
- structure of **latent space**