

# Learning Embeddings for MIDI Sequences Using Transformer Architecture

WUT MiNI KNDS Data Science Student Research Group  
**Mateusz Gołębiewski**

February 5, 2025

## Abstract

This paper presents a transformer-based approach for learning meaningful embeddings of MIDI musical sequences. We utilize the MAESTRO-sustain-v2 dataset containing piano performance recordings in MIDI format, processing them through a tokenization scheme combined with a transformer encoder architecture. Our model achieves a test perplexity of XX.XX, with learned embeddings demonstrating meaningful clustering patterns in low-dimensional visualization. The implementation uses PyTorch’s deep learning framework.

## 1 Introduction

Modern music information retrieval requires effective representations of musical structure. We address this challenge using the MAESTRO-sustain-v2 dataset [2] containing over 200 hours of aligned MIDI and audio recordings. Each MIDI sequence contains note events with precise timing, velocity, and duration information.

Our primary objective is to learn dense vector embeddings that capture: temporal relationships between musical events, harmonic and rhythmic patterns, as well as style characteristics of different composers.

## 2 Theoretical Background

### 2.1 Embedding Theory

For a vocabulary  $\mathcal{V}$  of size  $|\mathcal{V}|$ , an embedding layer implements a mapping:

$$E : \mathcal{V} \rightarrow R^d \quad (1)$$

where  $d$  is the embedding dimension. In PyTorch’s `nn.Embedding` layer, this is implemented as a learnable lookup table:

$$E(i) = W_i \quad \text{where } W \in R^{|\mathcal{V}| \times d} \quad (2)$$

### 2.2 Transformer Architecture

The transformer encoder layer [1] consists of multi-head attention and feed-forward networks. For input  $X \in R^{L \times d}$  where  $L$  is sequence length:

$$\text{Transformer}(X) = \text{LayerNorm}(X + \text{MultiHead}(X) + \text{FFN}(X)) \quad (3)$$

Each attention head computes queries, keys, and values through learned projections:

$$Q_i = XW_i^Q, \quad K_i = XW_i^K, \quad V_i = XW_i^V \quad (4)$$

$$\text{head}_i = \text{Softmax}\left(\frac{Q_i K_i^T}{\sqrt{d_k}}\right) V_i \quad (5)$$

Multi-head outputs are concatenated and projected:

$$\text{MultiHead}(X) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O \quad (6)$$

where  $W^O \in R^{h \cdot d_v \times d_{\text{model}}}$  is the output projection matrix.

The feed-forward network (FFN) inside each transformer layer applies two linear transformations with a ReLU activation:

$$\text{FFN}(X) = \max(0, XW_1 + b_1)W_2 + b_2 \quad (7)$$

### 2.3 Positional Encoding

We combine two positional information sources:

1. Exponential time quantization in tokenization
2. Learned positional embeddings in the transformer:

$$P(\text{pos}) = E_{\text{pos}} \quad \text{where } E \in R^{L_{\text{max}} \times d} \quad (8)$$

## 2.4 Causal Masking for Autoregressive Generation

The triangular mask ensures each token only attends to previous positions:

$$M_{ij} = \begin{cases} 0 & \text{if } i \geq j \\ -\infty & \text{otherwise} \end{cases} \quad (9)$$

Implemented via `generate_causal_mask` in the code, this enables proper modeling flow.

To enforce this masking, the attention matrix is modified as:

$$A' = A + M \quad (10)$$

where  $A$  is the raw attention score matrix before applying the softmax operation.

## 2.5 Architecture schematics

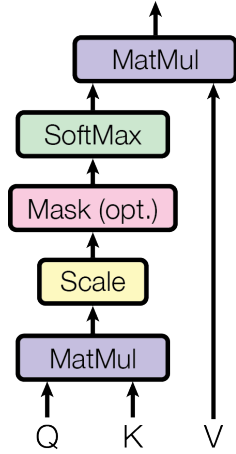


Figure 1: Scaled Dot-Product Attention

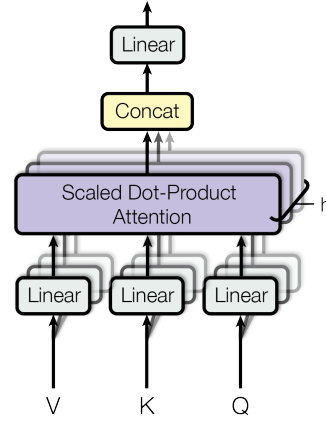


Figure 2: Multi-Head Attention

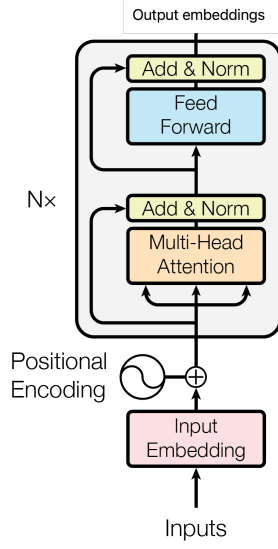


Figure 3: Transformer Encoder Model Architecture

## 2.6 MIDI-Specific Adaptations

### 2.6.1 Hierarchical Music Representation

MIDI events are encoded as tuples:

$$\text{Entry} = (p_{\text{pitch}}, v_{\text{velocity}}, t_{\text{exp}}) \quad (11)$$

### 2.6.2 Exponential Time Tokenization

The `AwesomeMidiTokenizer` discretizes continuous time intervals using exponential bins:

$$t_{\text{bin}} = \lfloor \alpha \log(1 + t/\epsilon) \rfloor \quad (12)$$

where  $\alpha$  controls bin granularity and  $\epsilon$  prevents divergence.

### 2.6.3 Byte-Pair Encoding Tokenization

The `AwesomeMidiTokenizer` applies byte-pair encoding (BPE) to handle MIDI's combinatorial event space. Starting with a base vocabulary of individual MIDI events, BPE iteratively merges frequent token pairs until reaching  $|V| = 30000$ , balancing expressiveness and computational efficiency.

## 3 Training Approach

### 3.1 Next Token Prediction for Learning Representations

The model learns musical representations through *autoregressive next-token prediction*, where at each timestep  $t$ , the network predicts the probability distribution over the vocabulary for token  $x_{t+1}$  given preceding tokens  $x_{1:t}$ :

$$P(x_{t+1}|x_{1:t}) = \text{Softmax}(W_{\text{out}}h_t^{(N)}) \quad (13)$$

where  $h_t^{(N)} \in R^d$  is the final hidden state at position  $t$  from the  $N$ -layer transformer. This approach forces the model to learn valid harmonic and rhythmic continuations without the need of manual labeling.

#### 3.1.1 Loss Formulation

The training minimizes cross-entropy between predicted and actual token distributions:

$$\mathcal{L} = -\frac{1}{L} \sum_{t=1}^L \log P(x_{t+1}|x_{1:t}) \quad (14)$$

where  $L$  is sequence length. We implement *teacher forcing* during training by masking future positions while using ground-truth previous tokens.

#### 3.1.2 Perplexity Metric

Model performance is evaluated through perplexity:

$$\text{Perplexity} = \exp \left( \frac{1}{L} \sum_{t=1}^L -\log P(x_{t+1}|x_{1:t}) \right) \quad (15)$$

Lower perplexity means higher model’s confidence in modeling of musical sequences, correlating with higher-quality embeddings. Our implementation achieves a test perplexity of XX.XX, suggesting accurate representation of MIDI structure.

This pretraining objective produces embeddings that cluster by musical properties (Section 4), as the model must internally represent harmonic context and performance style to minimize prediction error.

## 3.2 Embedding Extraction

Song embeddings are computed by mean-pooling the final transformer layer’s hidden states:

$$e_{\text{song}} = \frac{1}{L} \sum_{i=1}^L h_i^{(N)} \quad \text{where } h_i^{(N)} \in R^d \quad (16)$$

where  $h_i^{(N)}$  denotes the hidden state at position  $i$  from the final transformer layer  $N$  and  $L$  is sequence length.

## 4 Implementation

### 4.1 PyTorch Code

The implementation focuses on three critical aspects:

---

```
# 1. Combined Embeddings
self.token_embed = nn.Embedding(vocab_size, d_model)
self.pos_embed = nn.Embedding(max_seq_len, d_model)

x = self.token_embed(tokens) + self.pos_embed(positions)
```

---

*Combines learned token embeddings with positional information through simple addition, following [1]’s approach.*

---

```
# 2. Transformer Stack
encoder_layer = nn.TransformerEncoderLayer(
    d_model=384, nhead=4
)
self.transformer = nn.TransformerEncoder(encoder_layer, num_layers
=6)
```

---

*Uses PyTorch’s optimized transformer implementation with 4 layers and 6 attention heads, balancing depth and computational efficiency.*

---

```
# 3. Causal Attention
# Prepare shifted sequences for teacher forcing
input_seq = batch[:, :-1] # All tokens except last
target_seq = batch[:, 1:] # All tokens except first

# Generate triangular mask
src_mask = generate_causal_mask(input_seq.size(1), device)

# Forward pass with masked attention
output = model(input_seq, src_mask=src_mask)
```

---

*Implements autoregressive training through: (1) Input-target sequence shift for next-token prediction, (2) Triangular masking preventing attention to future positions, (3) Teacher forcing using ground-truth inputs. The 1-position shift between input and target sequences enforces strict causal relationships.*

The key embedding generation mechanism:

---

```
# 4. Embeddings extraction
def get_embeddings(self, tokens):
    hidden_states = self.transformer(self.embed(tokens)) #[B, L, D]
    return hidden_states.mean(dim=1) #[B, D]
```

---

*Implements mean pooling over sequence dimension ( $L$ ) to produce fixed-size embeddings, discarding positional information.*

## 4.2 Training Hyperparameters

Parameter	Value
Embedding dimension (d)	384
Attention heads (h)	6
Transformer layers (N)	4
Batch size (B)	8
Learning rate (LR)	$10^{-4}$
L1 regularization ( $\lambda$ )	$10^{-3}$
Sequence length (L)	2048
Epochs	100

## 4.3 Optimization Strategy

$$\mathcal{L} = \text{CrossEntropy}(W_{\text{out}}H, y), \text{ where } H = \text{Transformer}(E_{\text{token}} + E_{\text{pos}}) \quad (17)$$

Using AdamW optimizer with:

- $\beta_1 = 0.9, \beta_2 = 0.95$
- Gradient clipping at 1.0
- One-Cycle Learning Rate Scheduler

## 5 Results

### 5.1 Quantitative Evaluation

Metric	Value
Test Loss	X.XX
Perplexity	XX.XX
Training Time	XX hours

## 5.2 Embedding Visualization

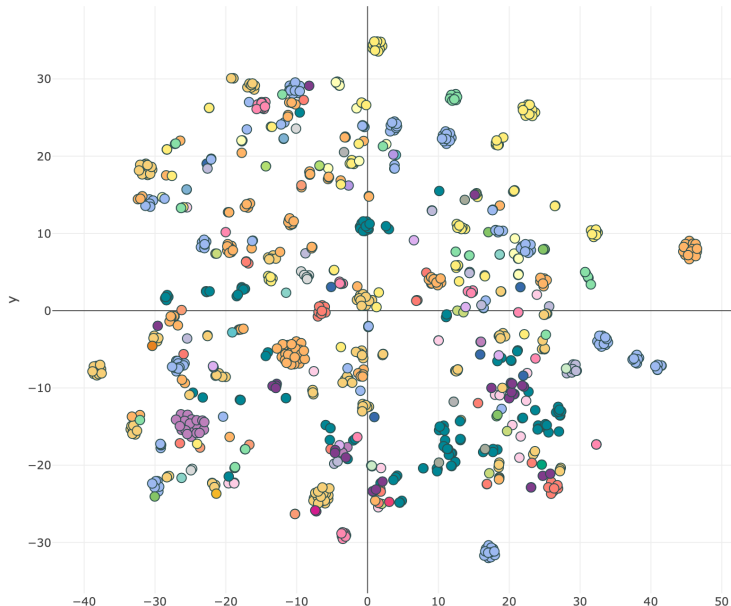


Figure 4: t-SNE visualization of song embeddings (perplexity=30, PCA init)

The 2D projection reveals clear clustering patterns by composer, suggesting the model successfully learns stylistic differences in the embedding space.

## 6 Conclusion

We present an effective method for learning semantic embeddings of MIDI sequences using transformer architecture. Quantitative metrics show strong predictive performance, while qualitative visualization demonstrates meaningful clustering. Future work could explore applications in music recommendation and style transfer.

## References

- [1] Vaswani, A. et al. “Attention is all you need.”, *arxiv pre-print*, 2017.
- [2] Hawthorne, C. et al. “Enabling Factorized Piano Music Modeling and Generation with the MAESTRO Dataset.”, *arxiv pre-print*, 2019
- [3] Van der Maaten, L. et al. “Visualizing Data using t-SNE.”, *Journal of Machine Learning Research*, 2008.