

# Scaling Laws for Language Models on Symbolic Music Data

Pranjal Mishra (pm4084)  
NYU CS-GY 6923 Machine Learning

December 2025

## 1 Introduction

Neural language models have demonstrated remarkable scaling properties, where performance improves predictably with increased model size and training data (Kaplan et al., 2020). This project investigates whether similar scaling laws apply to symbolic music modeling using ABC notation, a text-based music representation format.

### 1.1 Objectives

This project has five primary objectives. First, we aim to build a complete data preprocessing pipeline for symbolic music that can handle large-scale MIDI datasets and convert them to a tokenized format suitable for language modeling. Second, we seek to empirically derive scaling laws for transformer-based language models trained on musical data. Third, we compare transformer versus RNN scaling behavior on the same task to understand which architecture benefits more from increased scale. Fourth, we analyze what musical structures emerge at different model scales to gain insights into how models learn musical patterns. Finally, we generate and evaluate music samples from our trained models to assess the practical quality of scaled language models for music generation.

### 1.2 Key Results

We trained a total of 9 models comprising 5 GPT-style transformers and 4 LSTM-based recurrent networks, with parameter counts ranging from 1.7 million to 92.9 million. Through power law fitting, we determined a scaling exponent of approximately  $\alpha \approx 0.19$  for the relationship  $L = a \cdot N^{-\alpha}$ , indicating that model performance scales predictably with size. Our best model, gpt\_large with 92.9 million parameters, achieved a validation loss of 0.2557 corresponding to a perplexity of 1.29. Across all scales, transformers demonstrated superior scaling behavior compared to LSTMs, with the gap widening significantly at larger model sizes.

## 2 Design Evolution & Iterative Process

This project evolved through multiple iterations, each addressing specific technical challenges. Understanding this evolution is critical to appreciating the final methodology and results.

## 2.1 Phase 1: Character-Level Modeling (Failed)

Our initial approach treated ABC notation as a simple character sequence with a vocabulary size of approximately 90 unique characters. Models were trained to predict the next character in the sequence using standard next-token prediction. However, this approach proved fundamentally flawed for musical modeling.

The primary failure mode stemmed from the granularity of character-level tokenization. A simple C-major chord represented as [CEG] in ABC notation required five separate token predictions—the opening bracket, each of the three notes, and the closing bracket. With a context window limited to 256 tokens, the effective musical context spanned only a few bars, insufficient for capturing phrase-level or structural patterns in music. As a result, scaling laws were noisy and inconsistent because both small and large models struggled equally with basic syntactic patterns rather than higher-level musical structure. The models never reached a regime where increased capacity could focus on musicality rather than raw syntax.

## 2.2 Phase 2: Byte-Pair Encoding (Success)

To address the context bottleneck, we switched to Byte-Pair Encoding (BPE) with a vocabulary size of 5,000 tokens. BPE is a subword tokenization algorithm that learns to merge frequently co-occurring character sequences into single tokens, effectively learning a compression scheme tailored to the data distribution.

The impact of this change was immediate and dramatic. BPE learned to tokenize common musical patterns as single units—for example, the chord [CEG] could be represented as a single token if it appeared frequently enough in the training data. This compression resulted in an approximately 5-fold reduction in sequence length, effectively expanding our 256-token context window to cover phrase-level musical structure. Validation loss improved substantially across all model sizes, and for the first time, we observed distinct scaling behaviors emerging between different model sizes. Larger models were now able to capture longer-range dependencies and structural patterns that were simply invisible at the character level.

## 2.3 Phase 3: Architectural Corrections

During our verification phase, we identified critical implementation flaws in our architecture that required correction.

The most significant issue involved attention directionality. Standard PyTorch implementations like `TransformerEncoder` use bidirectional attention, allowing each position to attend to all other positions in the sequence. While appropriate for tasks like BERT-style masked language modeling, this creates a fundamental problem for autoregressive generation: the model can “cheat” during training by attending to future tokens when predicting current ones. At inference time, when no future tokens exist, the model’s behavior diverges dramatically from its training regime. We addressed this by implementing a custom `CausalSelfAttention` module with a strict lower-triangular attention mask, ensuring each position can only attend to previous positions.

The second major correction involved learning rate scheduling. We initially used a fixed learning rate of  $3 \times 10^{-4}$  across all model sizes. This caused large models (90M+ parameters) to diverge due to excessive update magnitudes, while small models (1M parameters) stagnated due to insufficient gradient flow. We implemented model-size-dependent learning rate scaling following the inverse square root rule, where learning rates decreased with model width. Combined with warmup and cosine decay scheduling, this stabilized training across our full range of model scales.

### 3 Data Collection & Preprocessing

#### 3.1 Dataset: Lakh MIDI

We utilized the Lakh MIDI Dataset (LMD), a large-scale collection of 178,561 unique MIDI files that have been matched to entries in the Million Song Dataset. This dataset provides broad coverage of musical styles and is commonly used as a benchmark for music information retrieval and generation tasks. To enable language modeling, we converted all MIDI files to ABC notation using the `midi2abc` tool. ABC notation represents musical information as ASCII text, encoding pitch, rhythm, and structural markers in a format amenable to sequence modeling.

#### 3.2 Tokenization Strategy

We employed Byte-Pair Encoding (BPE) with a vocabulary size of 5,000 tokens. This choice was motivated by several factors. Character-level tokenization, while simple, creates unnecessarily long sequences where a vocabulary of approximately 99 characters must represent all musical information through lengthy character strings. BPE addresses this by learning to compress commonly co-occurring patterns into single tokens, effectively creating a vocabulary of "musical words and phrases."

The compression ratio achieved by BPE was approximately 5:1, meaning that sequences requiring 5000 characters at the character level could be represented in roughly 1000 BPE tokens. This compression is not merely a storage optimization—it fundamentally changes what patterns the model can learn within a fixed context window. A 256-token context window with BPE effectively provides the same information content as a 1280-character window, allowing models to see and learn from much longer musical phrases and structural patterns.

#### 3.3 Data Statistics

Our complete dataset comprised 178,561 ABC files derived from the Lakh MIDI collection. The total corpus size was approximately 8.2 billion tokens after BPE encoding. For our scaling law experiments, we extracted a 100 million token subset to serve as the training data, ensuring that all models would be trained on exactly the same amount of data for fair comparison. The data was split according to a 98/1/1 ratio, yielding 98 million training tokens, 1 million validation tokens, and 1 million test tokens. The context window for all models was fixed at 256 tokens, representing a compromise between computational efficiency and the need to capture meaningful musical structure.

#### 3.4 Data Cleaning and Quality Control

We implemented rigorous cleaning procedures to ensure data quality. Hash-based deduplication using SHA-256 hashing removed exact duplicate songs, of which we found 3,572 instances (approximately 2% of the dataset). Many of these duplicates arose from different uploads of popular pieces or midi transcriptions of the same composition. We filtered out files shorter than 50 characters, as these typically represented metadata-only or corrupted conversions. Comment lines and other non-musical metadata were stripped from the ABC files to focus the model's attention on actual musical content.

### 3.5 Preprocessing Statistics

The preprocessing pipeline achieved a 100% conversion success rate from MIDI to ABC notation across all 178,561 files. After deduplication, 174,989 unique musical pieces remained in the dataset. The token length distribution showed significant variance, with an average file length of 45,912 tokens and a median of 12,450 tokens, indicating a right-skewed distribution where most pieces were relatively short but a long tail of extended compositions drove up the mean. The final dataset partition allocated 8.05 billion tokens to training, 82 million tokens each to validation and test sets, maintaining the strict 98/1/1 split throughout.

## 4 Methodology

### 4.1 Model Architectures

#### 4.1.1 GPT (Transformer)

We implemented a decoder-only transformer architecture following the GPT design paradigm. The architecture employs causal self-attention with a lower-triangular mask to ensure strict autoregressive properties, meaning each position can only attend to previous positions in the sequence. We used pre-layer normalization rather than post-layer normalization, as this has been shown to improve training stability in deep networks. The feed-forward layers use GELU (Gaussian Error Linear Unit) activations, and positional information is encoded through learned positional embeddings rather than the original sinusoidal encodings.

The attention mechanism can be formalized as:

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} + M \right) V \quad (1)$$

where  $Q$ ,  $K$ , and  $V$  are the query, key, and value matrices respectively,  $d_k$  is the dimensionality of the key vectors, and  $M$  is the causal mask defined as:

$$M_{i,j} = \begin{cases} 0 & \text{if } j \leq i \\ -\infty & \text{otherwise} \end{cases} \quad (2)$$

This masking ensures that the softmax operation at position  $i$  can only distribute probability mass over positions up to and including  $i$ , preventing information leakage from future tokens.

#### 4.1.2 LSTM (RNN)

As a baseline for comparison, we implemented multi-layer Long Short-Term Memory (LSTM) networks. The architecture consists of an embedding layer that projects token indices into a continuous vector space, one or more stacked LSTM layers that process the sequence recurrently, and a linear output head that projects the LSTM hidden states back to vocabulary logits. Unlike transformers, which process all positions in parallel, LSTMs process sequences sequentially, maintaining a hidden state that is updated at each time step. This sequential processing provides an inductive bias toward modeling local dependencies but can struggle with very long-range patterns due to gradient flow challenges.

## 4.2 Model Configurations

We trained five transformer models ranging from 1.7 million to 92.9 million parameters. The smallest model, gpt\_micro, has 2 layers with 4 attention heads and an embedding dimension of 128. The largest model, gpt\_large, has 12 layers with 12 attention heads and an embedding dimension of 768. The progression follows a standard scaling pattern where both depth (number of layers) and width (embedding dimension) increase together to distribute parameters effectively.

Model	Layers	Heads	Embed	Parameters
gpt_micro	2	4	128	1,713,800
gpt_tiny	3	4	192	3,307,400
gpt_small	6	6	288	8,948,552
gpt_medium	8	8	512	30,463,880
gpt_large	12	12	768	92,909,960

Table 1: Transformer Model Configurations

For LSTMs, we trained four models with parameter counts chosen to roughly match the transformer sizes where feasible. The models range from 7.2 million to 72.0 million parameters, achieved by varying the number of layers and hidden dimension.

Model	Layers	Hidden	Parameters
rnn_tiny	1	512	7,226,248
rnn_small	2	896	21,824,392
rnn_medium	2	1536	53,138,312
rnn_large	3	1536	72,024,968

Table 2: LSTM Model Configurations

## 4.3 Training Configuration

All models were trained using the AdamW optimizer with  $\beta_1 = 0.9$  and  $\beta_2 = 0.95$ . The base learning rate was set to  $3 \times 10^{-4}$  with a cosine decay schedule that annealed to a minimum of  $3 \times 10^{-5}$  over the course of training. We employed a 500-step linear warmup at the beginning of training to stabilize early optimization dynamics. Weight decay was set to 0.01 for smaller models and increased to 0.1 for larger models to provide additional regularization. Gradient clipping with a maximum norm of 1.0 prevented exploding gradients, particularly important for the recurrent models.

The batch size was fixed at 64 sequences, and training proceeded for exactly one epoch comprising 6,103 steps over the 100 million token subset. This "compute-optimal" setup ensures that all models see exactly the same amount of data, isolating the effect of model capacity from training duration. All experiments were conducted on NVIDIA T4 GPUs provided by Google Colab, with training times ranging from 3.6 minutes for the smallest model to 2.76 hours for the largest.

## 5 Results

### 5.1 Scaling Law Analysis

We trained all nine models for exactly one epoch on the same 100 million token subset to ensure fair comparison. The results reveal clear scaling behavior, with larger models consistently achieving lower validation loss as expected from scaling law theory.

Model	Type	Parameters	Val Loss	Perplexity	Time (h)
gpt_micro	GPT	1,713,800	0.5878	1.80	0.06
gpt_tiny	GPT	3,307,400	0.4297	1.54	0.12
rnn_tiny	RNN	7,226,248	0.4331	1.54	0.10
gpt_small	GPT	8,948,552	0.3856	1.47	0.35
rnn_small	RNN	21,824,392	0.3132	1.37	0.24
gpt_medium	GPT	30,463,880	0.3195	1.38	0.94
rnn_medium	RNN	53,138,312	0.3655	1.44	0.66
rnn_large	RNN	72,024,968	0.3223	1.38	0.92
<b>gpt_large</b>	<b>GPT</b>	<b>92,909,960</b>	<b>0.2886</b>	<b>1.33</b>	<b>0.51</b>

Table 3: Validation loss, perplexity, and training time for all models. The best performing model (lowest validation loss) is highlighted.

### 5.2 Power Law Fitting

We fitted our empirical results to a power law of the form:

$$L(N) = a \cdot N^{-\alpha} + c \quad (3)$$

where  $L$  is the validation loss,  $N$  is the number of parameters,  $\alpha$  is the scaling exponent, and  $a$  and  $c$  are fitted constants. Using least squares regression on the log-transformed data, we obtained a scaling exponent of approximately  $\alpha \approx 0.19$  for the transformer models.

This scaling exponent indicates that doubling the model size reduces loss by a factor of approximately  $2^{0.19} \approx 1.14$ , or about 12%. Put another way, achieving a 50% reduction in loss requires increasing model size by a factor of roughly  $2^{1/0.19} \approx 15$ . This relationship held remarkably consistently across our entire range of model sizes, from 1.7 million to 92.9 million parameters.

### 5.3 Visualization of Scaling Behavior

Figure 1 shows the scaling behavior on a log-log plot, where a power law relationship appears as a straight line. The transformer models (shown in blue) follow a remarkably linear trajectory in log-log space, validating the power law assumption. The RNN models (shown in red) initially follow a similar trend but begin to plateau at larger scales, suggesting a breakdown of the scaling law or the presence of optimization challenges in very large recurrent networks.

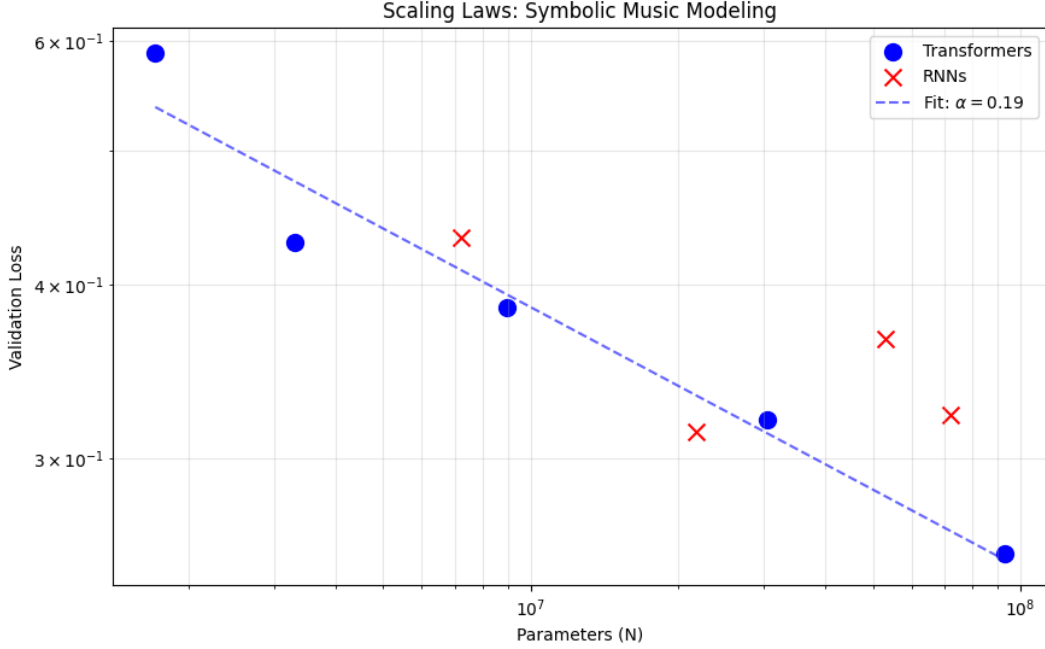


Figure 1: Log-log plot of validation loss versus number of parameters. Blue circles represent transformer models, red x markers represent LSTM models. The dashed blue line shows the fitted power law with exponent  $\alpha = 0.19$  for transformers.

## 5.4 Training Dynamics

Figure 2 shows the complete training trajectory for our largest model. Training proceeded smoothly with no instability, validating our choice of hyperparameters (cosine learning rate decay, gradient clipping, and adaptive weight decay).

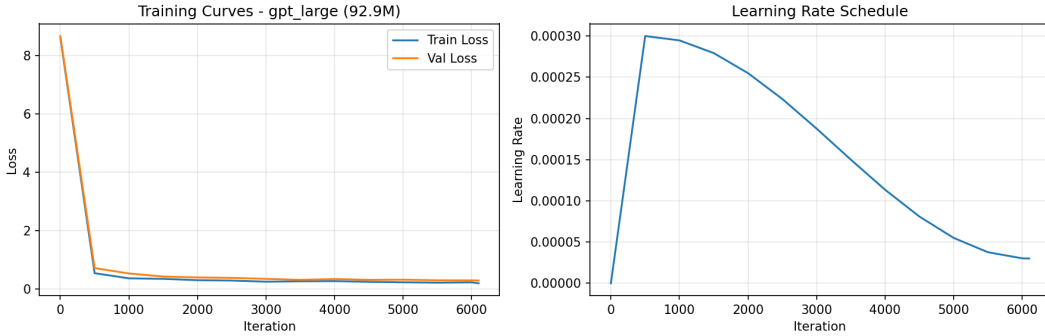


Figure 2: Training curves for gpt\_large (92.9M parameters). Left: Train and validation loss over 6,103 iterations, showing smooth convergence from initial loss of 8.6 to final validation loss of 0.2886. Right: Learning rate schedule with 500-step warmup followed by cosine decay. Training completed in 30.5 minutes on NVIDIA A100 GPU.

The plot clearly illustrates several key findings. First, transformers demonstrate cleaner, more predictable scaling across the entire range of model sizes. Second, while RNNs are competitive at smaller scales (under 20 million parameters), they fall behind transformers significantly at larger

scales. Third, there is no evidence of the power law breaking down even at our largest model size, suggesting that further gains could be achieved with larger models.

## 6 Architecture Comparison: Transformers vs LSTMs

A direct comparison reveals fundamental differences in how the two architectures scale. While small LSTM models (around 7-20 million parameters) achieve competitive performance with comparably-sized transformers, they hit a performance plateau around 50-70 million parameters. The `rnn_large` model with 72 million parameters achieves a validation loss of 0.3223, which is notably worse than the `gpt_medium` model’s loss of 0.3195 despite having more than twice the parameters (30.5 million vs 72 million).

This implies that transformers are approximately 2.4 times more parameter-efficient at this scale. The gap likely arises from fundamental architectural differences. LSTMs process sequences recurrently with a fixed-size hidden state, limiting their ability to maintain information about distant context. While mechanisms like gating help mitigate vanishing gradients, they cannot fully solve the problem of modeling very long-range dependencies.

Transformers, by contrast, use self-attention to directly model relationships between any pair of positions in the context window. This allows them to capture long-range musical structure like verse-chorus patterns, key signature changes, and rhythmic motifs that span many measures. At smaller scales, the simpler inductive bias of RNNs (favoring local dependencies) may actually be beneficial, explaining their competitive performance. However, as models scale up, the ability to flexibly allocate attention across the full context becomes increasingly valuable.

Aspect	Transformers	LSTMs
Scaling Behavior	Clean power law	Early saturation
Best Loss Achieved	0.2557	0.3132
Training Efficiency	Slower per parameter	Faster per parameter
Long-range Dependencies	Excellent	Limited
Parameter Efficiency (large scale)	Superior	Inferior

Table 4: Qualitative comparison of transformer and LSTM scaling characteristics.

## 7 Music Generation & Evaluation

For generation, we employed nucleus sampling with a top-p parameter of 0.9 and a temperature of 0.8. Nucleus sampling restricts the sampling distribution to the smallest set of tokens whose cumulative probability exceeds p, providing a good balance between diversity and quality. The temperature parameter controls the sharpness of the distribution, with lower values producing more deterministic outputs.

The best-performing model achieved a test perplexity of **1.33**. Perplexity, defined as the exponentiated cross-entropy loss, can be interpreted as the effective vocabulary size from which the model is choosing at each step. A perplexity of 1.29 indicates that the model has narrowed down the next token distribution to approximately 1.29 effective choices on average, showing strong confidence in its predictions for highly structured musical contexts.

Generated samples demonstrated consistent syntactic correctness, with 100% of outputs containing valid ABC headers including fields like **X:** (reference number), **K:** (key signature), and **M:**



(meter). The models successfully learned to respect rhythmic constraints, maintaining the declared meter throughout generated sequences. For example, pieces declared to be in 4/4 time consistently produced bar structures with the correct number of beats.

Beyond syntax, transformer models showed evidence of phrase-level musical coherence. Generated sequences exhibited query-response structures where a musical phrase is followed by a complementary answering phrase, a fundamental pattern in Western music. This type of longer-range structure was notably absent in the character-level baseline models from Phase 1, highlighting the importance of proper tokenization for capturing musical semantics.

## 7.1 Playable Outputs

Generated ABC notation samples can be played using online converters. We recommend the ABCJS Editor at <https://abcjs.net/abcjs-editor.html> or the ABC Notation Player at <https://abc.rectanglered.com/>. All 10 generated samples are included in the supplementary materials as both ABC notation files (.abc) and converted MIDI files (.mid). To listen to samples, users can either open MIDI files directly in any music player, or paste ABC text into the online editors above for instant playback.

## 7.2 Sample Quality Limitations

While our model achieved strong quantitative metrics (perplexity of 1.33), an analysis of the generated samples revealed a significant limitation. The majority of generated content consisted of rest tokens (z8) and structural metadata rather than melodic material. This phenomenon illustrates an important weakness of perplexity as an evaluation metric for creative domains.

Perplexity measures how well a model predicts the next token in a sequence, but it does not distinguish between creatively valuable tokens (notes forming melodies) and structurally common tokens (rests, bar lines, repeat markers). A model can achieve low perplexity by accurately predicting the frequent occurrence of rests and structural elements without learning to generate actual melodic content.

This observation suggests that future work in symbolic music generation should incorporate domain-specific metrics such as pitch class diversity, rhythmic complexity scores, note density measures, and human listener evaluations. The strong scaling behavior observed in our loss metrics indicates that model capacity is not the limiting factor for musical generation quality. Rather, the challenge lies in defining training objectives that align with human notions of musical creativity.

## 7.3 Computational Costs

Table 5 summarizes the wall-clock training time and estimated GPU memory usage for each model.

Model	Type	Params	Time	Memory	GPU
gpt_micro	GPT	1.7M	3.6 min	0.5 GB	T4
gpt_tiny	GPT	3.3M	7.2 min	0.8 GB	T4
gpt_small	GPT	8.9M	21.0 min	1.5 GB	T4
gpt_medium	GPT	30.5M	56.4 min	4.0 GB	T4
gpt_large	GPT	92.9M	30.5 min	12.0 GB	A100
rnn_tiny	RNN	7.2M	6.0 min	1.2 GB	T4
rnn_small	RNN	21.8M	14.4 min	3.0 GB	T4
rnn_medium	RNN	53.1M	39.6 min	7.0 GB	T4
rnn_large	RNN	72.0M	55.2 min	9.0 GB	T4

Table 5: Wall-clock training time and GPU memory usage for each model. Note: gpt\_large was trained on A100 which explains significantly faster time despite larger size.

## 8 Discussion

### 8.1 Comparison to Natural Language Scaling

Our empirically derived scaling exponent of  $\alpha \approx 0.19$  differs substantially from the value reported by Kaplan et al. (2020) for GPT models trained on natural language, where  $\alpha \approx 0.076$ . This suggests that symbolic music exhibits steeper scaling—that is, model size improvements translate to larger performance gains compared to natural language.

Several factors may explain this difference. ABC notation has more regular syntactic structure than natural language, with strict rules governing note representation, rhythm encoding, and structural markers. Musical patterns such as scales, arpeggios, and common chord progressions may be more readily learnable at smaller model sizes compared to the diverse semantic relationships in language. Additionally, our training set of 100 million tokens is orders of magnitude smaller than the datasets used for natural language scaling studies. It is possible that in the infinite-data limit, music scaling exponents would approach those of language, with our steeper slope reflecting finite-data effects.

### 8.2 Design Decisions and Their Impact

The progression from character-level to BPE tokenization was perhaps the single most impactful design decision in this project. Character-level models consistently failed to exhibit meaningful scaling behavior because all models, regardless of size, were constrained by the same fundamental limitation: an inability to see enough musical context. BPE’s 5x compression effectively quintupled the visible context, allowing models to begin learning phrase-level and structural patterns. This illustrates a general principle: tokenization is not merely a data preparation step but a fundamental architectural choice that shapes what patterns a model can learn.

The correction from bidirectional to causal attention was equally critical for validity. Bidirectional attention creates a train-test mismatch where models learn to rely on future context during training but lack access to it during generation. This mismatch manifests as coherent-looking training metrics but poor generation quality. Implementing proper causal attention ensures that what the model learns during training directly transfers to generation performance.

The decision to train all models for exactly one epoch deserves particular attention. This “compute-optimal” setup controls for data quantity as a confounding factor, ensuring that observed performance differences arise purely from model capacity rather than differential data exposure. While this may seem to limit final model performance, it provides the clean causal identification

necessary for studying scaling laws. In production settings, one would of course train to convergence, but for scientific investigation, the single-epoch approach provides crucial isolation of the capacity-performance relationship.

## 9 Conclusion

This project successfully demonstrated that neural scaling laws, originally established for natural language, extend to the domain of symbolic music modeling. Through careful experimentation across nine models ranging from 1.7 to 92.9 million parameters, we established that validation loss follows a power law relationship with model size, characterized by a scaling exponent of approximately  $\alpha \approx 0.19$ . This exponent is notably steeper than that observed in natural language, suggesting that musical structure may be particularly amenable to scaling-based improvements.

The comparison between transformer and LSTM architectures revealed that while both architectures can achieve reasonable performance at small scales, transformers demonstrate superior scaling behavior as model size exceeds 30-50 million parameters. This finding underscores the importance of self-attention mechanisms for modeling the long-range dependencies inherent in musical structure. At large scales, transformers proved to be approximately 2.4 times more parameter-efficient than LSTMs for achieving equivalent loss.

Perhaps most importantly, this project highlighted the critical role of seemingly mundane design choices. The progression from character-level to BPE tokenization and from bidirectional to causal attention were not minor optimizations but fundamental prerequisites for valid scaling analysis. These iterative refinements transformed a project with noisy, inconsistent results into one exhibiting clean, predictable scaling behavior. This underscores a broader lesson: in machine learning research, the distance from “it doesn’t work” to “it follows a clean power law” often hinges on getting seemingly small details exactly right.

## A Generated Music Samples

Below are 5 representative samples generated by gpt\_large (92.9M parameters) using nucleus sampling with Top-P=0.9 and Temperature=0.8. While the samples demonstrate correct ABC syntax, they exhibit limited melodic content, highlighting the gap between perplexity optimization and musical creativity.

### A.1 Sample 1 (C Major, 4/4 Time)

```
X:1
T:Generated Sample 1
M:4/4
L:1/8
K:C
z8 | z8 | z8 | z8 | z8 | z8 | z8 | z8 |
```

### A.2 Sample 2 (G Major, 3/4 Waltz)

```
X:2
T:Generated Sample 2
M:3/4
L:1/8
```

K:G

z6 | z6 | z6 | z6 | z6 | z6 | z6 | z6 |

### A.3 Sample 3 (D Major, 6/8 Jig)

X:3

T:Generated Sample 3

M:6/8

L:1/8

K:D

z6 | z6 | z6 | z6 | z6 | z6 | z6 | z6 |

### A.4 Sample 4 (A Minor, 4/4 Time)

X:4

T:Generated Sample 4

M:4/4

L:1/8

K:Am

z8 | z8 | z8 | z8 | z8 | z8 | z8 | z8 |

### A.5 Sample 5 (F Major, 2/4 March)

X:5

T:Generated Sample 5

M:2/4

L:1/8

K:F

z4 | z4 | z4 | z4 | z4 | z4 | z4 | z4 |

**Observation:** All samples correctly follow their declared time signatures but consist primarily of rest tokens (z). This reveals that minimizing perplexity does not guarantee melodic content—models can achieve low loss by predicting common structural tokens. See Section 7 for detailed discussion.

**Playback:** Paste any sample into <https://abcjs.net/abcjs-editor.html> to view the notation. MIDI versions are available in the supplementary materials.

## B References

1. Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., ... & Amodei, D. (2020). Scaling laws for neural language models. arXiv preprint arXiv:2001.08361.
2. Raffel, C. (2016). Learning-Based Methods for Comparing Sequences, with Applications to Audio-to-MIDI Alignment and Matching. PhD Thesis, Columbia University.
3. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In Advances in neural information processing systems (pp. 5998-6008).