# LZMidi: Compression-Based Symbolic Music Generation

Connor Ding, Abhiram Gorle, Sagnik Bhattacharya, Divija Hasteer, Naomi Sagan, Tsachy Weissman

Department of Electrical Engineering, Stanford University

Stanford, CA, 94305

czsding, abhiramg, sagnikb, dhasteer, nasagan, tsachy@stanford.edu

*Abstract*—Recent advances in symbolic music generation primarily rely on deep learning models such as Transformers, GANs, and diffusion models. While these approaches achieve high-quality results, they require substantial computational resources, limiting their scalability. We introduce LZMidi, a lightweight symbolic music generation framework based on a Lempel-Ziv (LZ78)-induced sequential probability assignment (SPA). By leveraging the discrete and sequential structure of MIDI data, our approach enables efficient music generation on standard CPUs with minimal training and inference costs. Theoretically, we establish universal convergence guarantees for our approach, underscoring its reliability and robustness. Compared to state-of-the-art diffusion models, LZMidi achieves competitive Fréchet Audio Distance (FAD), Wasserstein Distance (WD), and Kullback-Leibler (KL) scores, while significantly reducing computational overhead—up to 30× faster training and 300× faster generation. Our results position LZMidi as a significant advancement in compression-based learning, highlighting how universal compression techniques can efficiently model and generate structured sequential data, such as symbolic music, with practical scalability and theoretical rigor.

## I. INTRODUCTION

Deep learning–based generative models have achieved remarkable success in text, image, and audio synthesis. However, their substantial computational demands—particularly in sampling procedures such as those employed in diffusion-based models—pose significant challenges for practical deployment because of high latency and the reliance on specialized hardware. Recent research has explored more computationally tractable alternatives that maintain competitive output quality. [22] introduces a learning framework based on universal sequential probability assignments (SPAs) derived from the celebrated Lempel-Ziv (LZ78) [26] compression. Using LZ parsing under *stationary* and *ergodic* assumptions, the proposed approach represents sequences efficiently within a tree-based structure. This methodology offers strong theoretical guarantees on runtime, memory usage and has demonstrated promising results in text generation, achieving low-latency training and sampling.

Our present work focuses on the generation of symbolic music, which refers to the task of generating music in a structured, discrete format, typically represented as MIDI or other symbolic encodings rather than raw audio waveforms. With its discrete structure and finite alphabet, symbolic music is well-suited to LZ-based SPAs. In this work, we induce an LZ78-based SPA on symbolic music from the Lakh MIDI dataset and use this as a tool for symbolic music generation. Empirical evaluations indicate that LZ78-based SPA produces music of excellent perceptual quality—quantified using Fréchet Audio Distance (FAD)—while significantly reducing both training time and sampling overhead.

Notably, our proposed framework operates efficiently on standard CPUs, eliminating the need for energy-intensive GPUs. This not only enhances accessibility for researchers and practitioners with limited computational resources but also aligns with sustainability goals by reducing environmental impact. Collectively, these results position LZMidi as a compelling, resource-efficient alternative to deep learning approaches in symbolic music generation.

## II. RELATED WORK

### A. Deep Learning for MIDI Generation

State-of-the-art symbolic music generation predominantly employs deep learning techniques, with transformer-based models and diffusion frameworks leading recent advancements. Early influential contributions within the Magenta suite, such as MusicVAE [20], provided hierarchical VAEs for capturing long-term musical structure, while Music Transformer [12] introduced relative attention mechanisms to handle long-range dependencies. Transformer-GANs [17] integrate Transformers with GAN architectures to enhance sequence coherence and mitigate exposure bias.

Concurrently, diffusion-based works, such as [16] and [19], advanced symbolic music generation through iterative refinement procedures in continuous and discrete latent spaces, respectively. Although effective, these approaches remain computationally intensive. Recent efforts, including fast diffusion GANs [24], attempt to accelerate generation by reducing denoising steps. Motivated by these computational constraints, our work introduces compression-based universal sequential probability assignments (SPA) as a lightweight, scalable alternative achieving comparable generative quality with significantly reduced training and sampling costs.

### B. Universal Information Processing for MIDI data

Prior to the deep learning era, universal algorithms such as Lempel-Ziv (LZ77, LZ78) [25], [26] and Context-Tree

Weighting (CTW) [23] inspired extensive research in symbolic music processing. This includes (1) identifying themes and patterns in musical data, (2) compression-based similarity and perceptual measures, and (3) sequence-prediction algorithms for music generation [7], [18], [21]. In particular, in the context of symbolic music generation, [4] employs CTW to model musical event probabilities based on prior context, effectively capturing hierarchical and temporal dependencies. This approach directly inspires our use of LZ78-based sequential probability assignment. Additionally, [14] also hints at the possible efficacy of LZ78 for universal sequence modeling, which further motivates our efforts.

*C. Compression for Learning*

Compression serves as a proxy for learning because an algorithm's ability to compress data reflects its capacity to capture underlying structure. [8] argues that language modeling is equivalent to compression, as minimizing the expected negative log-likelihood minimizes the expected code length. In essence, a model with lower perplexity compresses text more efficiently, demonstrating its grasp of statistical regularities. Moreover, off-the-shelf compressors can function as probabilistic sequence models by selecting the token that minimizes compressed file size. Similarly, compressor-based methods, which operate without explicit model parameters, have demonstrated competitive performance in classification tasks, highlighting that compression itself can leverage inherent redundancies to approximate learned representations without conventional model training [13]. [1] Together, these results imply that minimizing redundancy is tantamount to learning the data's probabilistic structure

Further, Merhav and Weinberger [15] examine "universal simulation", a closely related concept demonstrating that sampling uniformly from the type class of a training sequence yields samples exactly matching the source distribution, thus minimizing dependency (mutual information) with the training data. They quantify this dependency as the "price of universality", reflecting the statistical cost incurred when generating samples without knowledge of the source distribution. Our LZMidi method is directly motivated by these foundational concepts: it leverages universal compression (via LZ78-based SPA) to efficiently approximate the underlying statistical structure of symbolic music. This allows LZMidi to effectively capture the intrinsic redundancy and repetitive structure in musical sequences [5], providing a resource-efficient alternative for symbolic music generation.

## III. Theoretical Background

In this section, we introduce the LZ78-based sequential probability assignment (SPA) from [22] and outline a theoretical justification of its application to symbolic music. Consider a finite, individual sequence $x^n = (x_1, x_2, \ldots, x_n)$, where

---

[1] However, we note that the validity of such compressor-based classification methods has sparked some debate regarding their empirical accuracy evaluation criteria and whether such accuracy comparisons might be inherently biased or overly optimistic.

---

each symbol $x_i$ takes a value from a finite alphabet $\mathcal{X}$. The tree-based Lempel-Ziv algorithm (LZ78) [26] is a universal compression algorithm where the number of bits expended per source symbol converge to the entropy rate, while also inducing a sequential probability assignment. Following this idea, [22] derives a practical sequential probability assignment: for a given sequence $x^{t-1} = (x_1, x_2, \ldots, x_{t-1})$, to predict the next symbol $x_t$, we can derive the probability of observing symbol $a$ given context $x^{t-1}$ shown below. We adopt the notations from [22]:

$$q^{LZ,\gamma}(a|x^{t-1}) \triangleq \frac{N_{LZ}(a|x^{t-1}) + \gamma}{\sum_{a' \in \mathcal{X}} N_{LZ}(a'|x^{t-1}) + \gamma|\mathcal{X}|}. \quad (1)$$

The probability assignment for symbol $a$ given context $x^{t-1}$ identifies the fraction of times symbol $a$ followed the context as $x_t$; it includes a perturbation term directed by $\gamma$ to tune how much the probability assignment should respect the empirical distribution from training data. Further details regarding the outline of the LZ78-induced SPA are provided in Appendix A. To justify our use of the LZ tree model, we present the following theorem, which establishes that an LZ tree trained on a sufficiently large dataset will closely approximate the true data distribution:

**Theorem III.1** (Universal Convergence of LZ78-SPA). *Let $P$ be the law of a process with components taking values in a finite alphabet $\mathcal{X}$, and let $Q^m$ be the LZ78-based sequential probability assignment (SPA) constructed using $m$ i.i.d training sequences from $P_{X^n}$. Then, for any fixed $n$,*

$$D\big(P_{X^n} \,\big\|\, Q^m_{X^n}\big) \xrightarrow[m \to \infty]{a.s.} 0,$$

*where $D(\cdot\|\cdot)$ denotes the Kullback–Leibler divergence.*

*Proof Sketch.* By construction, each node in the LZ78 tree tracks the empirical frequency of symbols following a particular context, and since every context with nonzero probability is visited infinitely often, the frequency with which a symbol $a$ appears converges to the true conditional probability $P(a|\text{context})$. Consequently, the LZ78-SPA, which returns $q(a|\text{context})$, an estimate of $P$ via empirical frequencies, assigns probabilities that approximate $P$ increasingly accurately with more training data. As a result, when the assigned probabilities agree with the source distribution for all positively-probable contexts, the relative entropy $D(P_{X^n}\|Q^m_{X^n})$ converges to 0 almost surely as $m \to \infty$.

A full, detailed proof is provided in Appendix B. Here, we emphasize that the key idea relies on the law of large numbers, which can be invoked because each relevant context is visited infinitely often. This allows us to conclude that local (node-wise) empirical distributions converge to the true underlying source probabilities.

**Remark III.2.** While standard universal compression theory typically relies on ergodicity for asymptotic guarantees, our practical setting employs fixed-length sequences, making

explicit ergodicity assumptions less critical for our empirical results.

## IV. METHODS

### A. Data Structure

We use the **Lakh MIDI Dataset (LMD)**, containing 648,574 samples, each with 256 notes drawn from the alphabet $\mathcal{X} = \{0, 1, \ldots, 89\}$, to train our LZ-based model for symbolic music generation. Here, 0 represents a rest, 1 denotes consecutive note continuation, and 2–89 correspond to actual pitch values. Figure 1 illustrates a sample MIDI sequence. The dataset's note distribution (Figure 2) highlights the dominance of rests (0s) and continuations (1s), while Figure 3 shows that actual note pitches are clustered around the mid-range. To build the alphabet for the LZ model, we simply treat each individual note as a symbol within the alphabet $\mathcal{X} = \{0, 1, 2, \ldots, 89\}$, allowing us to traverse the tree and update the SPA sequentially for each note in the dataset.
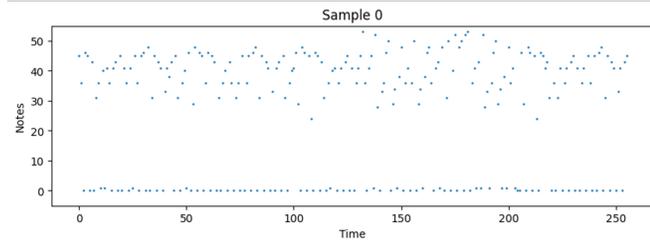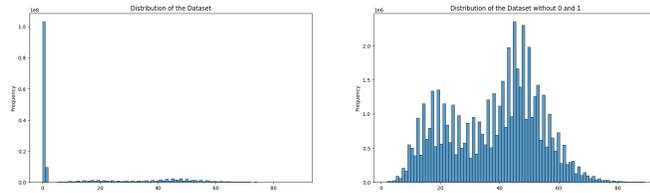


**Figure 1:** Sample Midi File



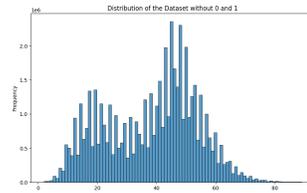**Figure 2:** Data Distribution including 0 & 1



**Figure 3:** Data Distribution without 0 & 1

### B. Baseline Model

For comparison, we use the absorbing state discrete diffusion denoising model (**ASD3PM**) by Plasser et al. [19], a discrete diffusion-based generative model for symbolic music. ASD3PM outperforms transformer-based autoregressive models [6] and prior continuous diffusion approaches [16]. ASD3PM follows an iterative denoising process, progressively refining a noisy symbolic sequence. Unlike continuous diffusion models that rely on latent representations, ASD3PM directly models discrete token transitions, making it well-suited for symbolic data like MIDI. The forward process masks tokens using a transition matrix $Q_t$, introducing an absorbing state for corrupted inputs, while the reverse process predicts $p_\theta(x_0|x_t)$ using a neural network. Inspired by [3], it simplifies training by parameterizing loss directly on the original data $x_0$ and dynamically adjusting diffusion steps for flexible sampling. Its hierarchical architecture combining convolutional and transformer layers for refinement helps it outperform larger latent-space diffusion models [16].

The loss function minimizes the evidence lower bound (ELBO) on the likelihood of the original data:

$$\mathcal{L}_{\text{ELBO}} = \mathbb{E}_{q(x_0)} \left[ \sum_{t=1}^{T} \frac{T-t-1}{T} \mathbb{E}_{q(x_t|x_0)} \left[ \log p_\theta(x_0|x_t) \right] \right] \tag{2}$$

where $T$ is the number of diffusion steps. This formulation ensures efficient training by focusing on intermediate diffusion steps. Its hierarchical architecture improves efficiency by reducing trainable parameters while maintaining fidelity and diversity.

Note that while ASD3PM is a state-of-the-art discrete diffusion model, it was originally designed for more elaborate tasks (e.g., polyphonic generation with flexible infilling). Our LZMidi focuses on simpler unconditional generation.

### C. Metrics

We will evaluate the quality of the generated music through the following set of metrics:

*1) Framewise Self-Similarity Metrics (Consistency and Variance):* To evaluate statistical similarity between generated and original sequences, we adopt the overlapping area (OA) metric from [16], which quantifies local pitch and duration distributions. Using a sliding 4-measure window (with a 2-measure hop), we fit Gaussian PDFs to pitch ($p(k)$) and duration ($d(k)$) distributions. The overlapping area (OA) between adjacent windows is defined as:

$$\text{OA}(k, k+1) = 1 - \text{erf}\left(\frac{c - \mu_1}{\sqrt{2}\,\sigma_1}\right) + \text{erf}\left(\frac{c - \mu_2}{\sqrt{2}\,\sigma_2}\right) \tag{3}$$

where $c$ is the intersection of the two Gaussian PDFs, and $(\mu_1, \sigma_1), (\mu_2, \sigma_2)$ are the respective means and standard deviations.

From the overlapping areas for pitch ($\text{OA}_P$) and duration ($\text{OA}_D$), we compute the *consistency* (C) and *variance* (Var) as follows:

$$\text{C} = \max\left(0, 1 - \frac{|\mu_{\text{OA}} - \mu_{\text{GT}}|}{\mu_{\text{GT}}}\right) \tag{4}$$

$$\text{Var} = \max\left(0, 1 - \frac{|\sigma_{\text{OA}}^2 - \sigma_{\text{GT}}^2|}{\sigma_{\text{GT}}^2}\right) \tag{5}$$

where $\mu_{\text{OA}}, \sigma_{\text{OA}}^2$ and $\mu_{\text{GT}}, \sigma_{\text{GT}}^2$ are the means and variances of generated and ground-truth samples, respectively.

*Consistency* measures alignment with ground truth, while *variance* reflects diversity. Higher consistency suggests realistic sequence structure, while balanced variance prevents mode collapse. However, strong OA scores alone **do not** guarantee perceptual quality, necessitating complementary evaluation (e.g., FAD). This motivates us to experiment with alternative qualitative metrics described below.

*2) Fréchet Audio Distance (FAD):* Fréchet Audio Distance (FAD), inspired by the Fréchet Inception Distance (FID) [11], quantifies how closely the statistical distribution of generated audio aligns with real data. It computes the Fréchet distance between feature embeddings extracted from a pre-trained model (e.g., VGGish [9]). Lower FAD scores indicate better perceptual similarity, making it a robust metric for evaluating generative quality.

*3) KL-Divergence:* Kullback-Leibler (KL) divergence measures the discrepancy between the probability distributions of real and generated data. While lower KL values suggest better alignment, it primarily favours distributional similarity over perceptual quality. Given that KL can favor models producing mode-collapsed outputs, we emphasize FAD as the primary metric for evaluating generation fidelity.

*4) Wasserstein Distance:* While consistency, variance, FAD, and KL divergence target audio-based evaluations, we also require a metric that directly examines numerical sequence distributions—especially for hyperparameter tuning without relying on costly neural-network inferences. We therefore use the Wasserstein Distance (WD) [2], [10], which measures the minimal cost of transforming one distribution into another. In our setup, we compare feature distributions from real and generated data; with a lower WD indicating a closer match between generated outputs and the ground truth.

## V. Experiments

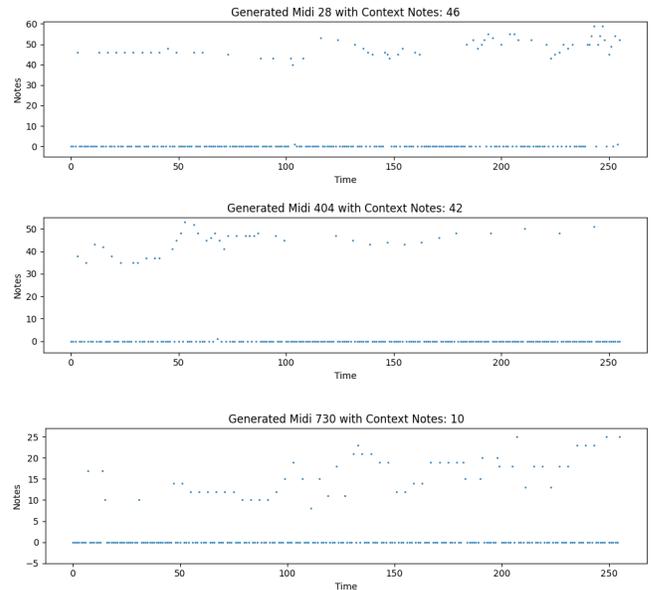### A. Training and Generation Setup

We trained the LZMidi model using the provided implementation of LZ-based SPA from [22]. The *Lakh MIDI Dataset* was split into an **80/20** train-test split. For training, we iterated through all samples in the training set, updating the LZ tree with each sequence. For evaluation, we generated 1,000 samples for each value of block length. All experiments were conducted on a CPU (Apple M1 Chip, 2021 MacBook). Following are some key parameters we can fine-tune for our training and generation of the LZ:

1) **Dirichlet Parameter** ($\gamma$): This parameter, used in the computation of the sequential probability assignment (Equation 1), determines the proximity of the SPA to the empirical distribution. A smaller $\gamma$ results in the SPA being closer to the empirical distribution.
2) **Top-$K$**: The number of allowed symbols from which the model predicts.
3) **Temperature**: This parameter controls the randomness of the generated output by adjusting the probabilities of predicted symbols. A value approaching 0 samples the most likely outcome; a value approaching 1 samples directly from the SPA; a value approaching $\infty$ samples from a uniform distribution over the symbols.
4) **Minimum Context**: The minimum context length that the SPA maintains during prediction. We set this value to 64 in our experiment.

We performed a hyperparameter sweep using Optuna [1] to find the $\gamma$, Top-$K$, and temperature that optimizes Wasserstein distance. We suggest a categorical selection of the

hyperparameters to Optuna for the hyperparameter sweep with the Wasserstein distance as the objective to minimize. We observed that $\gamma$ and temperature affect the generation quality the most, with a smaller $\gamma \approx 5 \times 10^{-5}$ and a larger temperature $T \approx 0.8$ being optimal in terms of the Wasserstein distance. For the final generation, we choose $\gamma = 5 \times 10^{-5}, T = 0.8$, and $K = 8$ to be our selection of hyperparameters.

To generate a sample, we randomly select a symbol from the alphabet as seed data, corresponding to a direct child node of the LZ root, and generate a sequence of length 256. The generated sequence is then mapped from its integer representation to note values and post-processed to ensure that no '1' follows '0'. Finally, the sequence is plotted and converted into both MIDI and WAV formats for metric computation and audio analysis.



**Figure 4:** MIDI plots for Generated Samples using the LZMidi Model.

Note that we independently trained the ASD3PM baseline due to significant differences in sequence length—our setup uses sequences of 256 tokens, whereas Plasser et al. [19] trained on 1024-token sequences—making direct use of these results unsuitable for a fair comparison. [19] reports a 24-hour training duration on 4x NVIDIA 2080 Ti GPUs. Due to computational constraints, we fixed our training time to approximately one hour. We include the metrics of the fully-trained model in the Appendix C.

### B. Results

We computed the consistency and variance metrics by comparing generated samples with 1000 randomly sampled sequences from both the training and test sets (in Table I). LZMidi exhibits high consistency and variance for both pitch and duration, closely matching the dataset statistics. In fact, variance is greater in all cases for LZMidi.
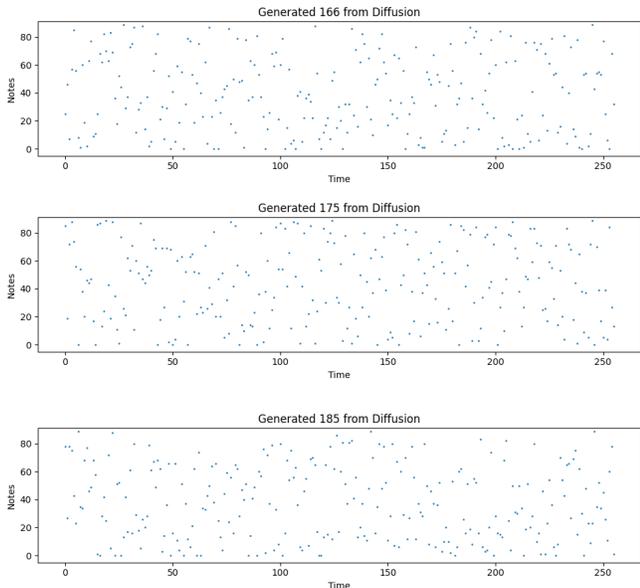
**Figure 5:** MIDI plots for Generated Samples using the D3PM Model.

**Table I:** Consistency and Variance

| | Training Set | | | | Test Set | | | |
|---|---|---|---|---|---|---|---|---|
| | Pitch | | Duration | | Pitch | | Duration | |
| | C | Var | C | Var | C | Var | C | Var |
| LZMidi | 0.97 | **0.92** | 0.97 | **0.93** | 0.97 | **0.93** | 0.97 | **0.94** |
| ASD3PM | **0.98** | 0.85 | **0.99** | 0.87 | **0.98** | 0.86 | **0.99** | 0.87 |

We evaluate the quality of the generated MIDI samples using the three aforementioned metrics: Wasserstein Distance (WD), Fréchet Audio Distance (FAD), and Kullback-Leibler (KL) Divergence. As shown in Table II, we compare our generated data to both the training and testing datasets. LZMidi achieves a much lower Wasserstein Distance (WD) and Fréchet Audio Distance (FAD) in both the training and test sets, indicating superior fidelity and distributional alignment.

**Table II:** WD, FAD and KL Divergence metrics

| | Training | | | Test | | |
|---|---|---|---|---|---|---|
| | WD | FAD | KL | WD | FAD | KL |
| LZMidi | **8.57** | **0.69** | **1.42** | **8.39** | **0.64** | **1.37** |
| ASD3PM | 27.91 | 4.22 | 2.29 | 27.96 | 4.05 | 2.26 |

### C. Training, Generation Time, and Memory Usage

We record the training time, generation time (per sample), and memory usage of LZMidi to underscore its computational advantages over deep learning–based methods. Table III summarizes these metrics for different $L$ settings. The metrics for our models are substantially lower than those of D3PM. For instance, while our training time is fixed at approximately one hour (reflecting our computational constraints), [16] report a 6.5-hour training duration on an Nvidia Tesla V100 GPU.

### D. Floating point operations (FLOPS) for ASD3PM Baseline

We analyze the computation (FLOPS) required for sequential MIDI generation using the ASD3PM baseline. We set the

**Table III:** Training Time, Generation Time, Memory Usage

| | Training Time (s) | Generation Time (s/sample) | Model Size (MB) |
|---|---|---|---|
| LZMidi | **107.7** | **0.016** | **287.1** |
| ASD3PM | 3480 | 5.4 | 306.2 |

sequence length for the generated MIDI files to 256 timesteps as in the experiments. To evaluate the computational efficiency of our diffusion-based baseline model for symbolic music generation, we analyzed the floating-point operations per second (FLOPS) across different layers of the network during training and inference. Table IV provides a breakdown of the total FLOPS and the corresponding CPU utilization metrics for key operations.

**Table IV:** FLOPs and CPU Utilization of the Diffusion Baseline

| Operation | Calls | Total FLOPs (MFLOPs) | Self CPU % | CPU Total Time (ms) |
|---|---|---|---|---|
| `aten::addmm` | 145 | 620,622.774 | 37.36% | 736.399 |
| `aten::bmm` | 48 | 12,884.902 | 6.91% | 126.123 |
| `aten::mul` | 1 | 377.487 | 1.08% | 21.089 |
| `aten::add` | 97 | 203.424 | 1.25% | 14.255 |
| `DataParallel::forward` | 50 | 104.858 | 2.14% | 18.332 |
| `aten::expand` | 243 | 0.00 | 0.56% | 55.675 |
| `aten::reshape` | 245 | 0.00 | 0.88% | 56.436 |

## VI. CONCLUSION

In this work, we introduced a novel approach to symbolic music generation using LZ78-based sequential probability assignment (SPA). Our method generates high-fidelity musical samples while substantially reducing computational requirements relative to state-of-the-art deep learning methods. Our experiments show that LZMidi effectively captures musical patterns, maintaining diversity and consistency metrics that rival or exceed those of diffusion-based models. Our experiments suggest that LZMidi rivals the diffusion baseline on all the aforementioned qualitative metrics (WD, FAD, KL) during both training and testing. Moreover, the method exhibits significant computational benefits—lower training/generation times and reduced memory usage—thus enabling efficient symbolic music generation on CPUs.

Future work will extend this study by training on longer sequences (e.g., 64-bar samples), as well as polyphonic sequences, and comparing results against fully trained diffusion models to assess scalability.

**Remark VI.1.** Some of our generated samples (corresponding to the plots in 4) are attached here for the reader's listening.

### REFERENCES

[1] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework, 2019.

[2] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. In *Proceedings of the 34th International Conference on Machine Learning*. PMLR, 2017.

[3] Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021.

[4] Ron Begleiter, Ran El-Yaniv, and Golan Yona. On prediction using variable order markov models. *Journal of Artificial Intelligence Research*, 22:385–421, 2004.

[5] Anton Chen and Ross Greer. Quantifying repetition in symbolic music using lempel-ziv compression. *2023 IEEE 14th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEM-CON)*, pages 0387–0393, 2023.

[6] Kristy Choi, Curtis Hawthorne, Ian Simon, Monica Dinculescu, and Jesse Engel. Encoding musical style with transformer autoencoders, 2020.

[7] Darrell Conklin and Ian H Witten. Multiple viewpoint systems for music prediction. *Journal of New Music Research*, 24(1):51–73, 1995.

[8] Grégoire Delétang and Others. Language modeling is compression. In *Proceedings of ICLR 2024*, 2024.

[9] Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 776–780, 2017.

[10] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C. Courville. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems*, 2017.

[11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium, 2018.

[12] Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M. Dai, Matthew D. Hoffman, Monica Dinculescu, and Douglas Eck. Music transformer, 2018.

[13] Zhiying Jiang and Others. 'low-resource' text classification: A parameter-free classification method with compressors. In *Findings of ACL 2023*, 2023.

[14] G. Langdon. A note on the ziv-lempel model for compressing individual sequences (corresp.). *IEEE Transactions on Information Theory*, 29(2):284–287, 1983.

[15] N. Merhav and M.J. Weinberger. On universal simulation of information sources using training data. *IEEE Transactions on Information Theory*, 50(1):5–20, 2004.

[16] Gautam Mittal, Jesse Engel, Curtis Hawthorne, and Ian Simon. Symbolic music generation with diffusion models, 2021.

[17] Aashiq Muhamed, Liang Li, Xingjian Shi, Suri Yaddanapudi, Wayne Chi, Dylan Jackson, Rahul Suresh, Zachary C. Lipton, and Alex J. Smola. Symbolic music generation with transformer-gans. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(1):408–417, May 2021.

[18] Marcus Pearce and Daniel Müllensiefen. Compression-based modelling of musical similarity perception. *Journal of New Music Research*, 46(2):135–155, 2017.

[19] Matthias Plasser, Silvan Peter, and Gerhard Widmer. Discrete diffusion probabilistic models for symbolic music generation, 2023.

[20] Adam Roberts, Jesse Engel, Colin Raffel, Curtis Hawthorne, and Douglas Eck. A hierarchical latent vector model for learning long-term structure in music, 2019.

[21] Pierre-Yves Rolland, Jean-Gabriel Ganascia, and G Raskinis. Modeling temporal knowledge for harmonic analysis: A first step towards musical style recognition. *Computers and the Humanities*, 35(1):65–91, 2001.

[22] Naomi Sagan and Tsachy Weissman. A family of lz78-based universal sequential probability assignments, 2024.

[23] F.M.J. Willems, Y.M. Shtarkov, and T.J. Tjalkens. The context-tree weighting method: basic properties. *IEEE Transactions on Information Theory*, 41(3):653–664, 1995.

[24] Jincheng Zhang, György Fazekas, and Charalampos Saitis. Composer style-specific symbolic music generation using vector quantized discrete diffusion models, 2024.

[25] J. Ziv and A. Lempel. A universal algorithm for sequential data compression. *IEEE Transactions on Information Theory*, 23(3):337–343, 1977.

[26] J. Ziv and A. Lempel. Compression of individual sequences via variable-rate coding. *IEEE Transactions on Information Theory*, 24(5):530–536, 1978.

*A. LZ78-Based Sequential Probability Assignment for Symbolic Music Generation*

In this work, we leverage the LZ78-based Sequential Probability Assignment (SPA) [22] for modeling symbolic music sequences. The model employs the incremental parsing algorithm of LZ78 to construct a prefix tree representation of an input training corpus of note events. Unlike conventional deep learning approaches that require extensive parameter optimization, the LZ78-SPA model learns symbol-by-symbol probability assignments incrementally and efficiently.

*Preliminaries and Notation:* Let $\{x_t\}_{t=1}^n$ be a sequence of discrete musical symbols drawn from a finite alphabet $\mathcal{A}$, where $|\mathcal{A}| = A < \infty$. Each symbol $x_t \in \mathcal{A}$ may represent a specific note or a rest indicator. The LZ78 parsing procedure incrementally partitions the sequence $x^n = x_1 x_2 \cdots x_n$ into a set of phrases (nodes) that form a growing prefix tree.

We denote by $Z(x^t)$ the set of phrases (including the empty root phrase) in the LZ78 parsing of the prefix $x^t$. Each phrase is a node in the prefix tree. For a given symbol $x_t$, let $z(x_t)$ represent the phrase (node) associated with $x_t$, and let $z_c(x_{t-1})$ represent the LZ78 context of $x_t$, i.e., the phrase prefix excluding the current symbol $x_t$.

*Universal Sequential Probability Assignment:* We define a sequential probability assignment (SPA) as a family of conditional distributions:

$$q = \{q_t(\cdot \mid x^{t-1}) : t \geq 1\}, \quad q_t(a|x^{t-1}) \in M(\mathcal{A}), \tag{6}$$

where $q_t(a|x^{t-1})$ is the probability of symbol $a \in \mathcal{A}$ given the observed sequence $x^{t-1}$.

Following [22], our LZ78-SPA conditions the probability assignment for $x_t$ on its LZ78 context $z_c(x_{t-1})$. To introduce smoothness and avoid zero-probability assignments for unseen events, we employ a Dirichlet prior with parameter $\gamma > 0$. This prior acts as an additive smoothing factor. Formally, define $N(a|x^{t-1}, z_c(x_{t-1}))$ as the count of symbol $a$ for phrases with prior context $(z_c(x_{t-1})$.

The LZ78-based SPA with a Dirichlet prior is given by:

$$q^{LZ,\gamma}(a|x^{t-1}) \triangleq \frac{N_{LZ}(a|x^{t-1}) + \gamma}{\sum_{a' \in \mathcal{X}} N_{LZ}(a'|x^{t-1}) + \gamma|\mathcal{X}|}. \tag{7}$$

Equivalently, since $\sum_{b \in \mathcal{A}} N(b \mid x^{t-1}, z_c(x_{t-1}))$ is the number of times we have visited the context $z_c(x_{t-1})$ in the prefix tree, this can be written as:

$$q_t(a|x^{t-1}) = \frac{N(a|x^{t-1}, z_c(x_{t-1})) + \gamma}{N(\cdot|x^{t-1}, z_c(x_{t-1})) + \gamma A},$$

where

$$N(\cdot|x^{t-1}, z_c(x_{t-1})) = \sum_{b \in \mathcal{A}} N(b|x^{t-1}, z_c(x_{t-1})).$$

*Training Procedure:* Training the LZ78-SPA consists of two steps:

1) **LZ78 Parsing:** Given the training corpus of symbolic music, we parse the entire training set using the LZ78 algorithm. This involves incrementally building a prefix tree, adding a new branch whenever a previously unseen context-symbol combination is encountered. The result is a tree structure $Z(x^N)$ where $N$ is the length of the entire training set.
2) **Count Aggregation and Prior Incorporation:** For each node (phrase) in the LZ78 tree, we record the counts $N(a \mid x^{t-1}, z_c(x_{t-1}))$ for all symbols $a \in \mathcal{A}$ that follow the context $z_c(x_{t-1})$. After parsing, these counts are fixed and used with the Dirichlet prior parameter $\gamma$ to define the SPA.

*Loss Function (Log Loss):* The quality of a sequential probability assignment is typically measured by the log loss. For a given test sequence $x^n$, the log loss under the LZ78-SPA model is:

$$L(x^n) = -\frac{1}{n} \sum_{t=1}^n \log q_t(x_t \mid x^{t-1}).$$

This log loss corresponds to the negative log-likelihood per symbol and provides a measure of how well the learned model predicts the sequence. Lower values of $L(x^n)$ indicate better predictive performance.

In practice, the LZ78-SPA often achieves asymptotically optimal log loss behavior when compared to Markovian or finite-state models [22]. For symbolic music generation, this translates into capturing the underlying repetitive and hierarchical patterns of musical structure while requiring significantly fewer computational resources than typical deep learning approaches.

*B. Proof of the Theorem*

**Definition A.1.** Let $Q^m$ be the probability model induced by an LZ78 tree built with $m$ equal-length realizations $X^n \overset{\text{i.i.d.}}{\sim} P_{X^n}$ sampled from source $P$ over alphabet $\mathcal{A}$. Let $X^{(i),n}$ denote the $i^{\text{th}}$ such sequence generated. No assumptions are placed on $P$.

**Remark A.2.** In this setting, the depth of the LZ78 tree is upper-bounded by $n$.

**Definition A.3** (Symbol counts). $\mathcal{C}(a|x^n)$ is the number of times that symbol $a$ appears in $x^n$.

**Theorem A.4.** *Assume that the SPA at each node of the LZ78 tree, $q$, satisfies $q(a|y^m) - \frac{\mathcal{C}(a|y^m)}{m} \to 0$, for all individual sequences $\mathbf{y}$. Then as $m \to \infty$,*
$$D(P_{X^n}\|Q^m_{X^n}) \xrightarrow{\text{a.s.}} 0.$$

*Proof.* First, we define some additional notation:

- When parsing the $m^{\text{th}}$ sample at index $t$, denote the current node of the LZ78 tree by $z^m_t$. The subsequence of symbols seen at $z^m_t$ until time $t$ is denoted $\mathcal{S}(z^m_t, m, t)$. Denote the length of this subsequence by $\ell(z^m_t, m, t)$.

- The LZ78 SPA at sample $m$, step $t$ is denoted $q(\cdot|\mathcal{S}(z^m_t, m, t))$. The SPA for a node that has not yet see data is denoted $q(\cdot)$.

· Consider any $Y^n \in \mathcal{A}^n$ such that $P_{X^n}(Y^n) > 0$.
$$\log \frac{1}{Q^m_{X^n}(Y^n)} = \sum_{t=1}^{n} \log \frac{1}{q(Y_t|\mathcal{S}(z^m_t, m, t))}.$$

**Fact A.5.** *Fix $t \geq 1$, $Y^{t-1} \in \mathcal{A}^n$ such that $P_{X^{t-1}}(Y^{t-1}) > 0$. Then, $\forall a \in \mathcal{A}$,*
$$\frac{\mathcal{C}(a|\mathcal{S}(z^m_t, m, t)))}{\ell(z^m_t, m, t)} \xrightarrow{\text{a.s.}} P_{X_t|X^{t-1}}(a|Y^t) \quad \text{as } m \to \infty.$$

*Proof.* Fix some $m > 0$, and define
$$\mathcal{C}_m(Y^{t-1}) = \sum_{i=1}^{m} \mathbf{1}\left\{X^{(i),t-1} = Y^{t-1}\right\}.$$

We can bound $\ell(z^m_t, m, t)$ by a constant plus $\mathcal{C}_m(Y^{t-1})$ on both sides. Note that, for $\mathcal{C}_m(Y^{t-1}) \geq t$, no returns to the root occur before reaching $z^m_t$, so $z^m_t$ is a depth-$(t-1)$ corresponding directly to the prefix $Y^{t-1}$. Otherwise, $z^m_t$ is some node encountered after a return to the root, in which case we cannot say much about $\ell(z^m_t, m, t)$ relative to $\mathcal{C}_m(Y^{t-1})$.

By the law of large numbers, $\frac{1}{m}\mathcal{C}_m(Y^{t-1}) \xrightarrow{\text{a.s.}} P_{X^{t-1}}(Y^{t-1}) > 0$ as $m \to \infty$. Therefore, there almost surely exists some $M$ such that $\mathcal{C}_M(Y^{t-1}) \geq t$. From this point, consider $m > M$.

For a lower bound on $\ell(z^m_t, m, t)$, the node $z^m_t$ was visited for all but maybe the first $t$ times $Y^t$ was seen. For an upper bound, we consider all the possible times $z^m_t$ was visited that do not correspond to $Y^{t-1}$, *i.e.*, times $z^m_t$ was visited after a return to the root. There is exactly one return to the root for every leaf of the tree, so the number of extra visits is bounded. So, $\forall m > M$, almost surely $\ell(z^m_t, m, t) = \mathcal{C}_m(Y^{t-1}) + O(1)$.

By the same logic, $\mathcal{C}(a|\mathcal{S}(z^m_t), m, t)) = \mathcal{C}_m(Y^{t-1} \frown a) + O(1)$, where $\frown$ represents sequence concatenation.

By the law of large numbers, $\frac{1}{m}\mathcal{C}_m(Y^{t-1}) = P_{X^{t-1}}(Y^{t-1}) + o_p(1)$, and analogously for $\frac{1}{m}\mathcal{C}_m(Y^{t-1} \frown a)$. As a result,
$$\frac{\mathcal{C}(a|\mathcal{S}(z^m_t, m, t)))}{\ell(z^m_t, m, t)} = \frac{\frac{1}{m}\left(\mathcal{C}_m(Y^{t-1} \frown a) + O(1)\right)}{\frac{1}{m}\left(\mathcal{C}_m(Y^{t-1}) + O(1)\right)} = \frac{P_{X^t}(Y^{t-1} \frown a) + o_p(1)}{P_{X^{t-1}}(Y^{t-1}) + o_p(1)} \xrightarrow{\text{a.s.}} P_{X_t|X^{t-1}}(a|Y^{t-1}),$$

by Slutsky's theorem. $\square$

**Corollary A.6.** *We know that, almost surely, for sufficiently large $m$ $\ell(z_t^m, m, t) = \mathcal{C}_m(Y^{t-1}) + O(1)$. So, by the law of large numbers, $\ell(z_t^m, m, t)$ almost surely grows unbounded, for any $Y^{t-1}$ with non-zero measure under $P$.*

**Corollary A.7.** *By assumption, $\forall Y^{t-1} \in \mathcal{A}^{t-1}$ with nonzero measure under $P$,*

$$q(a|\mathcal{S}(z_t^m, m, t)) \to \frac{\mathcal{C}(a|\mathcal{S}(z_t^m, m, t)))}{\ell(z_t^m, m, t)}, \quad as \quad \ell(z_t^m, m, t) \to \infty.$$

*Applying the fact $\ell(z_t^m, m, t) \xrightarrow{a.s.} \infty$, along with Fact A.5, as $m \to \infty$,*

$$q(a|\mathcal{S}(z_t^m, m, t)) \xrightarrow{a.s.} P_{X_t|X^{t-1}}(a|Y^{t-1}).$$

By Corollary A.7, the continuous mapping theorem, and Slutsky's theorem (as $n < \infty$ is a fixed quantity), for any fixed $Y^n$ with nonzero measure under $P$,

$$\log \frac{1}{Q_{X^n}^m(Y^n)} = \sum_{t=1}^{n} \log \frac{1}{q(Y_t|\mathcal{S}(z_t^m, m, t))} \xrightarrow{a.s.} \sum_{t=1}^{n} \log \frac{1}{P_{X_t|X^{t-1}}(Y_t|Y^{t-1})} = \log \frac{1}{P_{X^n}(Y^n)}.$$

Applying this to the relative entropy,

$$D(P_{X^n}\|Q_{X^n}^m) = \mathbb{E}\left[\log P_{X^n}(X^n)\right] + \mathbb{E}\left[\log Q_{X^n}^m(X^n)\right] = \mathbb{E}\left[\log P_{X^n}(X^n)\right] + \sum_{Y^n : P_{X^n}(Y^n) > 0} P_{X^n}(Y^n) \log \frac{1}{Q_{X^n}^m(Y^n)}.$$

Applying Slutsky's theorem (using the fact that the summation has a fixed, finite number of terms),

$$\sum_{Y^n : P_{X^n}(Y^n) > 0} P_{X^n}(Y^n) \log \frac{1}{Q_{X^n}^m(Y^n)} \xrightarrow{a.s.} \mathbb{E}\left[\log P_{X^n}(X^n)\right],$$

and, therefore,

$$D(P_{X^n}\|Q_{X^n}^m) \xrightarrow{a.s.} \mathbb{E}\left[\log P_{X^n}(X^n)\right] - \mathbb{E}\left[\log P_{X^n}(X^n)\right] = 0.$$

$\square$

Additionally, in our future work, we intend to compare LZMidi to other symbolic music baselines (e.g., n-gram or transformer models) for a more direct comparison in unconditional settings.

### C. Fully-trained ASD3PM results

In addition to the partially-trained diffusion baseline (ASD3PM) results presented in Section V-B (Plasser et al. [19]), we provide here the results for the fully-trained ASD3PM diffusion model. We additionally emphasized several key experimental differences between their setup and ours in 1) Sequence Length (256 vs. 1024 tokens) and 2) Hardware & Training time.

Due to these differences, the following results are **not** directly comparable but instead provide additional context regarding the performance achieved by a longer-trained diffusion approach with extended sequences:

| Setting | Unconditional | | | |
|---|---|---|---|---|
| | Pitch | | Duration | |
| **Metric** | C | Var | C | Var |
| **ASD3PM** | 0.992 | 0.920 | 0.993 | 0.937 |

**Table V:** Fully-trained metrics for Melody 64 bar setting.

These fully-trained results serve to contextualize the partially-trained ASD3PM baseline discussed in the main text.