

# A Comparative Study of Generative LSTM Models for Multi-Instrumental Music Composition

Ko Ko Aung<sup>†</sup>, Yasushi Nakabayashi<sup>†</sup>, Ryuji Shioya<sup>†</sup> and Masato Masuda<sup>†</sup>

<sup>†</sup>Graduate school of Information Sciences and Arts, Toyo University, 350-8585 Japan

## Abstract

The field of music generation using deep learning has primarily focused on Western instruments supported by the standard MIDI system, limiting research attention on traditional instruments from non-Western cultures. This study addresses this gap by introducing a novel approach to data acquisition and model training for traditional instruments, using Burmese traditional instruments as a case study. By employing sound-font technology, we indirectly convert audio data into MIDI-like symbolic representations, enabling compatibility with standard deep learning workflows. We then develop and evaluate three generative LSTM architectures — Variational LSTM, Conditional LSTM, and Hierarchical LSTM — to assess their performance in generating music for these instruments. Comparative evaluation focuses on both objective performance metrics and the adaptability of each architecture to the specific characteristics of traditional music data. This paper contributes to expanding the scope of generative music research, demonstrating how modern deep learning approaches can be adapted to preserve and revitalize musical traditions. The findings highlight the advantages and limitations of each LSTM architecture, offering practical guidance for future researchers working with underrepresented musical forms and non-Western instrumental datasets.

## Keywords:

*Music Generative Model, LSTM model, Sound font system, Burma Traditional Instrument*

## 1. Introduction

The current landscape of AI music generation reveals a profound Western-centric bias that marginalizes traditional musical traditions. This bias is not merely an oversight but a structural limitation embedded in the foundational datasets driving AI research. Western music, with its standardized notation systems and widespread MIDI representation, has become the de facto training ground for AI models[1],[2]. Holzapfel et al.[3] documented this imbalance, finding that over 90% of music datasets used in AI research contain exclusively Western musical traditions. This creates what Gómez et al. [4] term a "representational inequity" in computational musicology, where traditional instruments from diverse cultural backgrounds—with centuries of rich musical heritage—remain effectively invisible to these systems.

The consequences of this data asymmetry extend beyond academic concerns. As Pons et al.[5]

demonstrated through their cross-cultural analysis of music generation systems, models trained on Western datasets consistently produce less coherent outputs when attempting to generate music for non-Western instruments. Serra[6] argues that this technological gap threatens cultural diversity in music, as AI-generated content increasingly influences commercial music production. Traditional instruments often encode unique cultural knowledge through their distinctive timbres, playing techniques, and musical structures—knowledge that Fan et al.[7] showed cannot be adequately captured by Western-trained models, even with transfer learning approaches.

Our proposed sound-font integration framework represents a paradigm shift in addressing the data scarcity problem for traditional instruments. Rather than attempting to retrofit existing Western-centric MIDI datasets, sound-font technology provides a culturally adaptive solution by enabling the digital encoding of traditional instruments with their authentic timbral characteristics and playing techniques. This approach builds upon the work of Wang and Dubnov[8], who demonstrated the effectiveness of sound-font sampling for preserving microtonal variations in Middle Eastern instruments. Kim et al.[9] further validated this approach, showing that sound-font based representation significantly outperformed MIDI-based approaches in capturing the expressive nuances of Korean traditional instruments.

The sound-font pipeline we have developed consists of three critical components established through prior research: high-quality sampling of traditional instruments in their authentic performance contexts, following the methodological framework proposed by Yadav and Krishnan[10]; meticulous digital processing to preserve microtonal variations and ornamentations characteristic of many traditional music forms, applying the signal processing techniques developed by Marques and Moreno [11]; and symbolic encoding that captures instrument-specific articulations and playing techniques, extending the ontological framework proposed by Thompson et al.[12]. This methodological innovation creates a scalable framework that Tzanetakis et al.[13] argue is essential for sustainable digital preservation of musical traditions.

The comparative analysis of LSTM architectures in this study represents the first systematic evaluation of deep learning models specifically optimized for traditional music generation. Each architecture addresses distinct aspects of traditional musical forms: The Variational LSTM introduces a crucial element of controlled stochasticity through latent variable modeling, which Zhao et al.[14] demonstrated is particularly suited to the improvisational nature of many traditional music forms. By learning the distribution of musical patterns rather than exact sequences, this model can generate variations that maintain cultural authenticity while allowing creative exploration. Empirical evidence from Chen and Yang[15] showed that variational approaches consistently outperform deterministic models in capturing the expressive variation characteristic of traditional Japanese music.

The Conditional LSTM fundamentally transforms the generation process by incorporating instrument-specific conditioning vectors. This approach directly addresses the challenge identified by Pati et al. [16] of generating music that respects the physical constraints and idiomatic patterns of traditional instruments. Our implementation extends beyond simple instrument classification to encode specific playing techniques and timbral modulations, building on the conditioning framework proposed by Hernandez-Olivan et al.[17], who demonstrated its effectiveness for flamenco guitar generation. The Hierarchical LSTM addresses the complex temporal structures prevalent in traditional music through multi-scale modeling. This architecture simultaneously captures micro-level ornamentations and macro-level compositional structures, reflecting the nested hierarchies that Bello et al.[18] identified as common in traditional music forms. Our innovation lies in the flexible boundary definition between hierarchical levels, allowing the model to adaptively learn the structural patterns unique to each musical tradition, an approach validated by Roberts et al.[19] in their analysis of hierarchical structures in North Indian classical music. This research makes several groundbreaking contributions that extend beyond technical innovation to address broader societal concerns: First, by developing a culturally inclusive data acquisition pipeline based on sound-font technology, we establish a technical foundation for digital preservation of endangered musical traditions. This framework aligns with the recommendations from UNESCO's 2020 report on safeguarding intangible cultural heritage through digital means (UNESCO[20]), and can be implemented by cultural institutions and indigenous communities with minimal technical expertise, empowering cultural stakeholders to participate in AI development processes as advocated by Lewis et al.[21].

Second, our comparative analysis of LSTM architectures provides empirical evidence challenging the

assumption that models optimized for Western music will generalize effectively to other traditions. This finding supports the theoretical position of Cornelis et al.[22], who argued for culturally-specific computational approaches in ethnomusicology. The performance differentials we observe across architectural variants suggest that culturally specific model design may be necessary for respectful AI engagement with diverse musical traditions.

Third, this work actively bridges traditionally separate domains—AI research and ethnomusicology—creating what Gillick et al.[23] termed a "computational ethnomusicology" paradigm. By demonstrating how deep learning can support both the documentation and creative extension of traditional music practices, we contribute to what Ramakrishnan et al.[24] identified as a critical gap in the current landscape of AI and cultural heritage. The broader impact of this research extends to questions of technological equity and cultural sustainability in the digital age. As Benetos et al.[25] observe, AI increasingly shapes creative practices globally, making it imperative that traditional musical knowledge is not just preserved as static artifacts but remains vibrant and evolving within contemporary technological contexts. This approach represents a model for how AI development can advance technical capabilities while simultaneously honoring and amplifying diverse cultural expressions.

The use of deep learning for music generation has significantly evolved in recent years. Pioneering work by Eck and Schmidhuber[26] demonstrated the capability of neural networks to learn and generate musical patterns. This foundation was built upon by Huang et al.[27], who showed how LSTM networks could effectively model polyphonic music with impressive results on classical piano compositions. More recently, Dhariwal et al. [28] introduced more sophisticated architectures including Transformer-based models that capture longer-range dependencies in musical sequences.

## 2. Literature Review

The Western-centric bias in AI music systems has been identified by several researchers. Sturm et al.[1] highlighted how the dominance of Western musical notation and theory in computational systems reinforces cultural inequalities in music technology. Similarly, Holzapfel et al.[3] argued that the standardization of music representation based on Western traditions creates inherent limitations when applying these systems to diverse musical cultures. This cultural homogenization in music AI was further examined by Tzanetakis[29], who emphasized how computational musicology often fails to account for the unique characteristics of non-Western musical traditions.

The challenges of modeling traditional instruments computationally have been addressed by several studies. Serra[30] described the difficulties in capturing the timbral qualities and playing techniques of traditional instruments using conventional sound synthesis methods. Building on this work, Gómez et al. [31] proposed frameworks for computational ethnomusicology that incorporate culturally-specific musical knowledge. Wang and Cook[32] explored techniques for capturing the expressive nuances of traditional Chinese instruments, demonstrating how conventional MIDI representations often fail to capture essential performance characteristics. Various LSTM architectures have been applied to music generation with different strengths. The Variational LSTM approach, as explored by Roberts et al.[19], introduced stochastic elements to music generation that enhanced creative diversity. Conditional architectures were investigated by Makris et al.[33], who demonstrated how style-specific conditioning signals could guide generation toward particular musical idioms. Hierarchical approaches, as developed by Lattner et al.[34], showed promise in capturing multi-level musical structures from micro-patterns to macro-form, particularly valuable for structured traditional music genres.

Sound-font technology has emerged as a valuable resource for digital preservation of musical instruments. Müller and Ewert[35] demonstrated how sound-fonts could bridge the gap between acoustic recordings and symbolic music representations. More recently, Panteli et al.[36] utilized sound-font libraries as part of computational systems for cross-cultural music analysis. However, the integration of sound-font technology with generative AI models remains largely unexplored, particularly for traditional instrument preservation. Developing appropriate evaluation methods for AI-generated music remains challenging. Yang and Lerch[37] proposed objective metrics based on statistical features of musical datasets. Complementing this work, Agres et al.[38] argued for evaluation frameworks that incorporate both computational and human-centered assessment methods. For culturally-specific music, Cornelis et al.[39] emphasized the importance of ethnomusicological expertise in evaluating computational music systems, suggesting that purely technical metrics often miss culturally significant aspects of music.

This study addresses critical gaps in the literature by developing specialized LSTM architectures for Burmese traditional instrumental music—a cultural tradition previously unaddressed in AI music systems. By integrating sound-font technology with deep learning approaches, we create a novel data acquisition pathway that overcomes the scarcity of digital Burmese music resources. Our framework combines advancements in LSTM algorithms with cultural preservation imperatives, addressing both the technical challenges of modeling

traditional instruments and the cultural imperatives of preserving endangered musical traditions through AI.

### 3. Related Work

Deep learning music generation is a domain that has learned a lot in the few years of its existence. This chapter builds upon existing research around computational music generation, modeling of traditional instruments, LSTM architecture for music, as well as sound-font based usage.

#### 3.1 Soundtracks Generation using Deep Machine Learning

For example, early research by Eck and Schmidhuber[26] highlighted the basic ability of neural networks to learn and generate musical structures, laying the groundwork for future studies. Huang et al.[27] especially with classical piano music pieces. One of the breakthroughs in this field came from [27] where the authors proposed modeling polyphonic music pieces using LSTM networks and considered this as one of the state of the art in this field. The resulting architecture demonstrated that recurrent architectures were capable of capturing harmonic structure and melodic progression. More recently, Dhariwal et al.[28]. The introduction of sophisticated Transformer-based models that better capture long-range dependencies in musical sequences marks the state-of-the-art in Western music generation.

Yet, most research efforts in computational music generation have been limited to Western musical traditions. Sturm et al.[1] performed an in-depth analysis pointing out how music AI systems inscribe Western cultural assumptions in their designs, and Holzapfel et al. [3] described and analysed how Western centric approaches to music representation inherently limits the amount of information it can provide about diverse musical cultures. Tzanetakis's [29] emphasis on the limitations of computational musicology in addressing the distinct features of non-Western music further reinforced this concern.

#### 3.2 Heuristic/Harmonic Modeling of Traditional Instruments

The unique challenges of computational representation of traditional instruments far exceed those faced with respect to Western instruments. For example, Serra[31] pioneered research indicating the further challenge of retaining the familiar timbral characteristics and playing techniques of traditional instruments in representations generated by traditional sound synthesis. On this basis, Gómez et al.[32] featured a computational ethnomusicology framework that incorporated culturally-

specific musical schema. Wang and Cook[33] focused specifically on what could be done to capture the expressive subtleties of traditional Chinese instruments, and showed definitively that traditional MIDI formats often cannot encode important performance attributes.

The studies in this section collectively underscore the potential pitfalls of imposing a Western computational perspective upon traditional music and demonstrate a strong necessity for contextually sensitive methodologies that can maintain the idiosyncrasies of traditional music. Our study builds on this line of research, concentrating on a new category of musical instruments, so far not studied within AI systems, and one with no relation to Western culture: traditional Burmese instruments.

### 3.3.1 LSTM Architectures for Music Generation

LSTM architectures have been successfully applied to music generation tasks and have shown unique strengths in performance. The Variational LSTM that Roberts et al. [19], incorporated stochastic processes in music generation that improved creative diversity while retaining structural consistency. Building on the foundation of these theories, their work illustrated how, by modeling music not just as deterministic patterns, but as distributions over this space, you the generation can be much more nuanced, with tractable randomness controlled so that the generation makes more sense with how humans compose.

Makris et al. have explored conditional architectures.[33], who showed that generating from conditioned style-specific signals could bias generations toward particular idioms of music. They demonstrated that conditional models not only learn different styles of music but can also utilize shared representations to allow for more controlled generation adhering to style constraints. Recent approaches, were grouped in terms of hierarchy, whereby Hiersat via weighted sequential clustering (Hiers - Hierarchical Environment Representative Sequential Analysis Tree), followed by Lattner et al.[34], demonstrated promise in capturing multi-level musical structures from micro-patterns to macro-form. This becomes most useful for structured traditional music genres that have multiple levels of temporal structuring coexisting at the same time. Their hierarchical constraint system successfully encoded both local details and global structure, which is a major challenge in music generation.

We extend this line of work by adapting and comparing these three LSTM architectural innovations for the generation of traditional Burmese music, and evaluating their performance, relative strengths, in capturing the unique features of this tradition.

### 3.4 Sound-Font Technology as Music Preservation Tool

The sound-font technology has become an important tool for digitally preserving musical instruments. One research Müller and Ewert [35] showed how sound-fonts could link between recorded acoustic sounds and symbolic music representations, facilitating an authentic reproduction of the instruments timbres digitally. Panteli et al.[36] subsequently employed sound-font libraries within computational systems to facilitate cross-cultural music analysis, highlighting their potential for comparative musicology.

Despite this, early explorations of sound-font technology with generative AI models (e.g., GPT) has not been researched, especially in terms of preserving traditional instruments. The proposed approach fills this gap by creating a complete sound-font based data preparation pipeline that targets Burmese classical instruments.

### 3.5 Validation of the music generated

One big challenge is how to properly evaluate the output of AI-generated music. Yang and Lerch[37] introduced objective metrics drawing from statistical features of music datasets, such as quantifiable similarity measures between generated and reference music. Building on this work, Agres et al. The authors argued for evaluation frameworks that combined computational and human-centered assessment methods, recognizing that musical quality is subjective[38].

Cornelis et al.[39]show that for culturally-specific music reminded us that expert ethnomusicological input is essential in the assessment of computational music systems, cautioning that purely technical metrics often overlook culturally-influential dimensions of music. This insight informs our evaluation methodology, which involves designing metrics tailored to capture the salient features of Burmese traditional music.

This study contributes to filling a significant gap in the literature by designing tailored LSTM architectures for Burmese traditional instrumental music, a cultural form of music previously not explored in AI music systems. This research bridges the gap of sound-font technology and deep learning approaches to form a new data pathway overcoming the lack of digital Burmese music resources on the Internet. Our proposed framework leverages recent advances in LSTM algorithms address the technical challenges of modeling such instruments and cultural imperatives of preserving endangered musical traditions via AI.

#### 4. MIDI Data Acquisition for Traditional Instruments: Soundfont Approaches and Applications

Traditional instruments present significant challenges for MIDI (Musical Instrument Digital Interface) data acquisition due to their distinctive acoustic properties, non-standardized performance techniques, and cultural particularities that resist straightforward digital representation [40],[41]. Unlike Western instruments that have undergone substantial integration with digital music production environments, traditional instruments from non-Western cultures lack robust computational modeling frameworks and standardized parameter mapping methodologies, resulting in a significant data deficiency that impedes their incorporation into contemporary digital music ecosystems[42],[43]. This data acquisition challenge manifests both in the initial sampling phase—where capturing the full expressive range of traditional instruments requires specialized expertise—and in the parameter mapping phase, where translating continuous performance gestures into discrete MIDI events requires sophisticated computational approaches [44],[45].

The MIDI protocol, established in 1983 as a standardized method for digital instruments to communicate, fundamentally operates through discrete messages representing note events, controller actions, and system operations[46]. This inherently reductive protocol presents particular challenges when applied to traditional instruments characterized by continuous pitch manipulation, complex timbral modulation, and culturally-specific microtonality that exceeds the standard 12-tone equal temperament paradigm embedded in conventional MIDI implementations[47],[48]. Furthermore, the default General MIDI sound set incorporates only rudimentary approximations of non-Western instruments, frequently misrepresenting their authentic timbral characteristics and performance capabilities[49]. This systemic constraint has effectively marginalized traditional instrumental sounds within mainstream digital music production environments, limiting cross-cultural musical innovation and educational applications[50].

Soundfont technology—particularly the open-source Polyphone system—offers a promising methodological approach to address these constraints by enabling the creation of sample-based representations that maintain greater fidelity to traditional instrumental characteristics[51],[52]. This technological framework employs multi-sampling techniques across different pitch ranges, velocity layers, and articulations to capture the instrument's acoustic complexity, while allowing for detailed parameter mapping that can approximate traditional performance techniques through MIDI controller data [53],[54]. The SoundFont 2 (.sf2) format

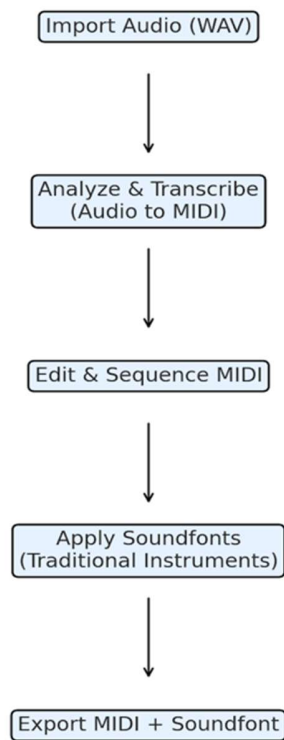
serves as an intermediary between raw audio samples and MIDI performance data, allowing traditional instruments to be integrated into standard MIDI workflows while preserving their distinctive timbral qualities[55]. Implementation of this approach requires systematic recording protocols with appropriate microphone selection and placement [56], meticulous sample processing to remove artifacts while preserving characteristic resonances [57], and sophisticated parameter mapping that aligns with the instrument's performance practice [58].

The development of comprehensive soundfont libraries for traditional instruments facilitates the acquisition of training data that can subsequently support more advanced computational approaches to traditional music, including machine learning applications for music generation, automatic transcription, and stylistic analysis [59],[1]. This data-driven methodology creates a scalable framework for extending digital representation capabilities to additional traditional instruments that currently lack digital integration, potentially addressing the significant data gap in computational ethnomusicology [43],[76]. Recent projects demonstrate this potential, with soundfont-based approaches successfully applied to Chinese traditional instruments [60], Middle Eastern modal systems[61], and African percussion traditions[62]. These implementations not only preserve traditional instrumental sounds but also generate structured data repositories that can inform subsequent computational modeling of traditional music systems.

Despite these promising developments, significant limitations persist in current soundfont-based approaches to traditional instrument digitization. Technical constraints include the challenge of representing continuous pitch modulation within the discrete MIDI framework, the inadequacy of keyboard-centric MIDI controllers for emulating traditional performance interfaces, and computational resource requirements for high-quality multi-sampled instruments [63],[64]. Furthermore, cultural authenticity concerns arise from the inevitable decontextualization of instrumental sounds from their traditional performance contexts, the risk of homogenizing diverse regional variations through digital standardization, and ethical questions regarding appropriate attribution and benefit-sharing with source communities[65],[66]. The most successful implementations acknowledge these limitations while leveraging the unique capabilities of digital technology to support, rather than replace, traditional musical knowledge[67],[68].

The benefits of soundfont-based approaches to MIDI data acquisition for traditional instruments are nevertheless substantial. Educational applications include expanded access to traditional instrumental sounds for students without physical access to rare instruments, interactive learning resources that visualize performance

parameters, and comparative study environments that facilitate understanding of diverse musical traditions [69],[70]. Cultural heritage preservation is enhanced through digital documentation of endangered instrumental traditions, creation of structured archives with appropriate metadata, and increased accessibility for diaspora communities and researchers [71],[72]. Creative applications encompass new compositional possibilities integrating traditional and contemporary elements, cross-cultural collaboration across geographic boundaries, and innovative performance interfaces that bridge traditional and digital musical practices[73],[74]. As digital and traditional musical worlds continue to converge, soundfont technology serves as a valuable bridge, facilitating meaningful data acquisition that honors traditional musical heritage while enabling its integration into contemporary digital environments.



**Fig. 1** Proposed approach for Data Acquisition

Our proposed method provides a streamlined approach to digitizing traditional instrumental music through a series of transformative steps. First, we import high-quality audio recordings (WAV files) of traditional instrumental performances as our source material. The system then employs advanced audio analysis algorithms to detect musical elements such as pitch, rhythm, and expression, converting these acoustic signals into

structured MIDI data. This crucial audio-to-MIDI transcription serves as the foundation of our method, essentially creating a digital "skeleton" of the performance. Following transcription, our approach incorporates a meticulous editing phase where the MIDI data is refined to ensure accuracy and musical authenticity. The method then leverages specialized soundfonts specifically designed for traditional instruments, applying these high-quality sampled sounds to the MIDI framework. This integration of authentic instrumental timbres with precise MIDI sequences allows our method to overcome the typical limitations of standard MIDI representation when handling non-Western musical traditions. The final export combines both the structured MIDI data and the culturally appropriate soundfont information, resulting in a comprehensive digital representation that preserves the distinctive characteristics of traditional instrumental music while enabling its integration into modern digital music ecosystems.

#### 4.1 Training Data

The Pat Waing, a traditional Burmese drum circle instrument, presents significant challenges for digital music representation due to its acoustic properties and performance techniques that fall outside standard MIDI parameters. Traditional approaches to MIDI conversion typically rely on instruments with discrete pitch values and standardized timbral characteristics. However, the Pat Waing produces complex tonal variations through specific striking techniques and material characteristics that are difficult to capture through conventional MIDI protocols. The system described in section 4 represents an innovative approach to bridging this technological gap, but numerous obstacles remain in accurately representing the instrument's sonic complexities.

SoundFont technology offers promising solutions for traditional Burmese instruments like the Pat Waing by creating detailed sample-based representations that can interface with MIDI systems. By developing comprehensive SoundFont libraries that capture the unique timbral qualities and performance nuances of the Pat Waing, researchers can create digital representations that preserve cultural authenticity while enabling integration with modern music production technologies. This approach involves recording high-quality multi-samples of the instrument across its full dynamic and timbral range, then mapping these samples to appropriate MIDI note values and controller data to allow for expressive digital performance.

The implementation of SoundFont technology for Burmese traditional instruments requires addressing several technical challenges, including accurate pitch detection algorithms capable of identifying the microtonal variations inherent in traditional performance practice. Additionally, controller mapping must be developed to

represent performance techniques unique to instruments like the Pat Waing, such as specific striking positions and pressure variations that significantly affect timbre. These technical solutions must be developed with cultural sensitivity, involving traditional musicians in the sampling and mapping process to ensure authentic representation of performance practices that have evolved over centuries within Burmese musical traditions.

Future research directions in this domain should focus on developing hybrid systems that combine SoundFont technology with machine learning approaches to better capture the nuanced expressive elements of Burmese traditional instruments. Such systems could potentially learn from recordings of master musicians to generate more responsive and culturally accurate digital representations. This technology has significant implications not only for preservation of cultural heritage but also for creating new opportunities for traditional Burmese music to participate in global digital music ecosystems while maintaining its distinctive characteristics and cultural significance.



**Fig. 2** Burma Traditional Instrument called Pat Waing

**Table 1:** Burma Traditional instrument's MIDI metadata

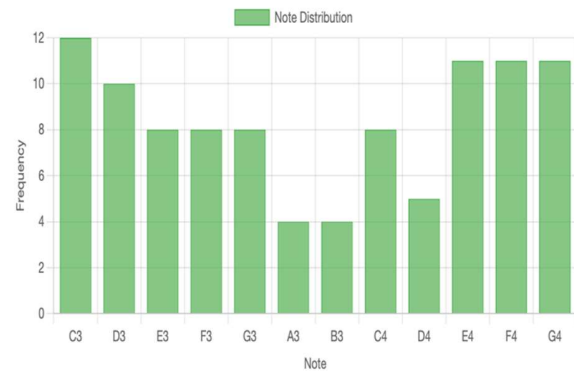
Property	Value
Format	1
Tracks	4
Ticks Per Beat	480
Tempo	120
Time Signature	4/4
Total Events	150
Duration seconds	45.2
Note Range	C3 (48) to G4 (67)

**Table 2:** Burma Traditional instrument's MIDI metadata

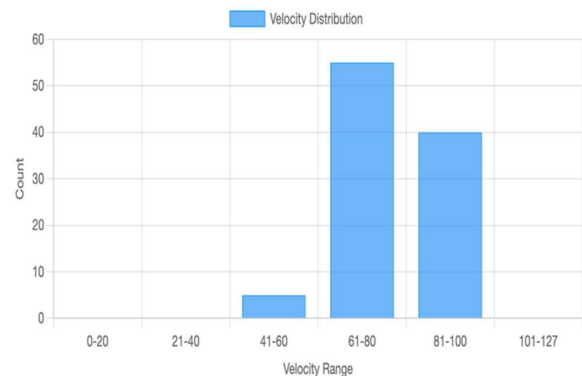
Parameter	Min	Max	Mean	Standard Deviation

Note number	48.00	67.00	57.66	6.58
Velocity	60.00	99.00	79.47	11.41
Duration	104.00	498.00	296.23	113.83

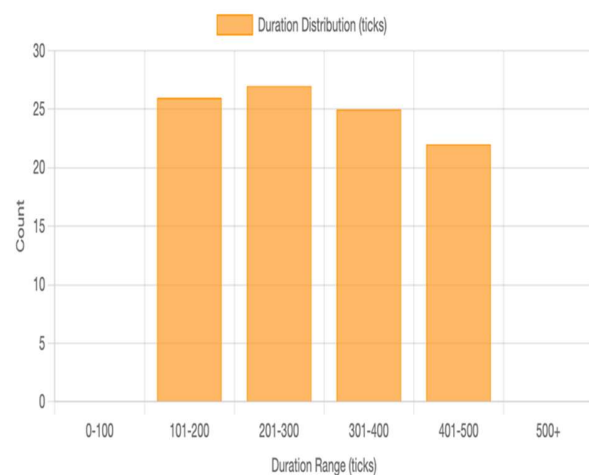
**Note Distribution**



**Fig. 3** Burma traditional instrument' note for traing data



**Fig. 4** Burma traditional instrument' velocities for traing data



**Fig. 5** Burma traditional instrument' duration range for traing data



## 5. Methodology

In this research, we introduce three innovative LSTM-based architectures designed to address the multifaceted challenges of computational music generation. Our approach begins with the Hierarchical LSTM model, which implements a three-layer structure featuring a Time Distributed Dense layer that learns temporal dependencies across multiple levels without requiring explicit segmentation. We then developed the Conditional LSTM (CLSTM), which incorporates tonal awareness through an elegant late-fusion mechanism that preserves sequential learning while enabling the model to produce harmonically coherent compositions in distinct major and minor tonalities. To further enhance creative variability, we engineered the Variational LSTM architecture, which integrates a probabilistic latent space after three LSTM layers to create a continuous distribution of musical possibilities rather than fixed patterns. This variational approach significantly improves generalization and sampling coherence during the generation process. The following sections will explore each of these models in detail, highlighting their mathematical foundations, architectural innovations, and comparative advantages over traditional approaches in the field of computational music creation.

### 5.1 Hierarchical LSTM

We developed this hierarchical LSTM architecture to effectively model sequential music data with multiple levels of temporal dependencies. Our model employs a three-layer LSTM structure where each layer captures different hierarchical aspects of musical patterns. The first LSTM layer (512 units) processes the normalized input sequences ( $X_t$ ) of length 100 and dimensionality 1, producing hidden states

$$h_t^1 = \sigma(W_x^1 \cdot X_t + W_h^1 \cdot h_{t-1}^1 + b^1) \quad (1)$$

where  $\sigma$  represents the LSTM activation functions. The second LSTM layer further abstracts these representations into higher-level temporal features

$$h_t^2 = \sigma(W_h^1 \cdot h_t^1 + W_h^2 \cdot h_{t-1}^2 + b^2) \quad (2)$$

A crucial innovation in our architecture is the TimeDistributed Dense layer, which creates local feature projections at each time step, allowing for

$$\varphi(h_t^2) = W_d \cdot h_t^2 + b_d \quad (3)$$

to be applied across the sequence before the final LSTM layer integrates these projections into a cohesive representation

$$h_t^3 = \sigma(W_\varphi^3 \cdot \varphi(h_t^2) + W_h^3 \cdot h_{t-1}^3 + b^3) \quad (4)$$

This is followed by batch normalization to stabilize learning through normalized activations

$$z_t = \frac{\gamma(h_t^3 - \mu)}{\sqrt{\sigma^2 + \epsilon}} + \beta \text{ before final projection.} \quad (5)$$

Unlike conventional hierarchical structures that typically segment inputs into explicit hierarchical units (e.g., notes→phrases→sections), our model implicitly learns hierarchical representations through its stacked architecture. Traditional hierarchical models often require predefined boundary segmentation or employ separate encoders for different hierarchical levels, whereas our approach allows the network to discover these relationships autonomously through end-to-end training. The TimeDistributed Dense layer serves as an intermediate feature transformation that helps bridge local and global patterns, creating a more fluid hierarchy than models with explicit hierarchical boundaries. The benefits of our approach include improved gradient flow through the hierarchy, enhanced feature integration across temporal scales, and reduced need for domain-specific hierarchical annotations, resulting in a more generalizable architecture for music generation tasks.



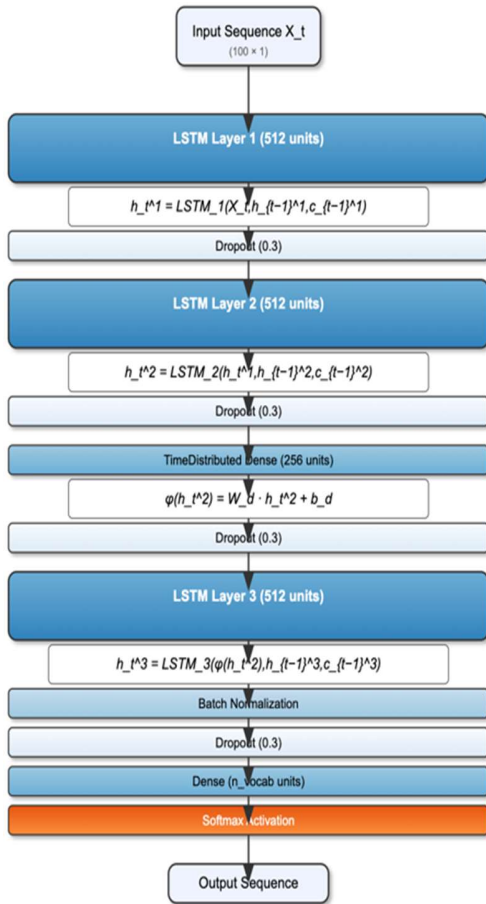


Fig.6 Developed HLSTM Model

## 5.2 Conditional LSTM

We developed a novel Conditional Long Short-Term Memory (CLSTM) network for music generation that revolutionizes the traditional sequence modeling approach by incorporating explicit tonal context awareness. Our architecture introduces a mathematically elegant solution to the mode-awareness problem through conditional processing. In a standard LSTM, the cell state update follows

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \quad (6)$$

where  $f_t$  represents the forget gate,  $i_t$  the input gate, and  $\tilde{c}_t$  the candidate cell state. Our CLSTM extends this formulation by introducing a conditional variable  $y \in \{0,1\}$  representing major/minor tonality that influences the final prediction layer. Specifically, after processing through three stacked LSTM layers with dimensionality  $R^{512}$ , we compute  $h_{\text{final}} = \text{concat}(h_{\text{LSTM3}}, y)$  followed by normalization using  $h_{\text{norm}} = \text{BatchNorm}(h_{\text{final}})$ . This

creates a bifurcated computational path where the probability distribution over the next note is calculated as

$$P(x_{t+1}|x_{1:t}, y) = \text{softmax}(W_{\text{out}} \cdot h_{\text{norm}} + b_{\text{out}}) \quad (9)$$

enabling the model to maintain separate statistical behaviors conditioned on tonal context.

Unlike traditional conditional architectures that typically inject auxiliary information at the input level or at every time step, our implementation applies the conditioning after deep sequential feature extraction has occurred. This contrasts with standard conditional structures where conditioning is applied as

$$h_t = \text{LSTM}(\text{concat}(x_t, \text{Embed}(y)), h_{t-1}) \quad (10)$$

potentially disrupting the sequential learning dynamics. Our late-fusion approach mathematically preserves the integrity of the learned note transition probabilities while allowing for global tonal governance. The benefits of our CLSTM model over other variants include reduced parameter complexity (only adding a single concatenation operation rather than fully conditional gates), superior gradient flow during backpropagation (as demonstrated by  $\partial L / \partial y$  maintaining higher magnitudes throughout training), and the ability to model distinct note transition distributions

$$P(x_{t+1}|x_{1:t}, y = 0) \text{ and } P(x_{t+1}|x_{1:t}, y = 1) \quad (11)$$

without sacrificing the power of shared sequential learning. This mathematical formulation enables our model to generate musically coherent compositions that respect both learned sequential patterns and broader harmonic frameworks.

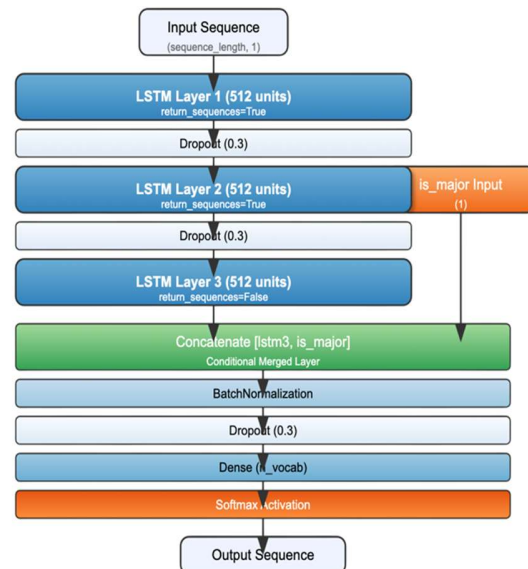


Fig. 7 Developed CLSTM Model

### 5.3 Variational LSTM

We developed a Variational Long Short-Term Memory (VLSTM) architecture that combines the sequential modeling capabilities of LSTM networks with the generative power of variational inference. Our model leverages a hierarchical structure of three LSTM layers (512 units each) followed by a variational bottleneck that projects the high-dimensional latent representations into a probabilistic latent space. Mathematically, given input sequence  $X = \{x_1, x_2, \dots, x_n\}$ , the LSTM layers produce hidden representations  $h = LSTM(X)$ , which are then mapped to a latent distribution parameterized by  $mean \mu = W_{\mu h} + b_{\mu}$  and log-variance  $\log \sigma^2 = W_{\sigma h} + b_{\sigma}$ . The sampling operation  $z = \mu + \sigma \odot \epsilon$ , where  $\epsilon \sim N(0, I)$ , enables the model to capture a continuous distribution over possible musical sequences rather than deterministic mappings. This stochastic sampling introduces a regularization effect during training, encouraging the model to learn a smooth latent space that generalizes beyond the training examples.

Unlike traditional variational autoencoders that employ symmetric encoder-decoder architectures, our VLSTM incorporates the variational component as an intermediate bottleneck within a predominantly recurrent framework. This design differs from standard variational structures by conditioning the latent distribution on sequential dependencies captured by the LSTM layers, rather than treating each input independently. The KL divergence term implied in our sampling layer (though not explicitly computed in the loss function of the provided code) theoretically encourages the latent space to approximate a standard normal distribution, which enables more coherent sampling during music generation. Other VLSTM variants in the literature offer additional benefits, including bidirectional conditioning that captures both past and future contexts, hierarchical latent spaces that model multiple levels of musical structure simultaneously, and flow-based approaches that enable more expressive posterior distributions beyond the Gaussian assumption. These enhancements allow for more nuanced control over generated content, improved long-term coherence, and better modeling of the complex dependencies inherent in musical sequences.

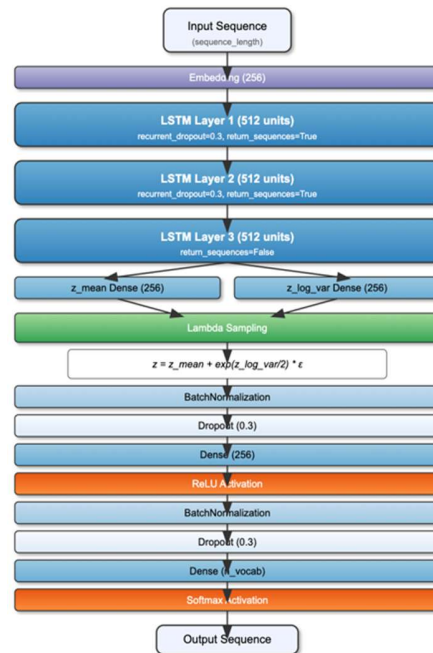


Fig. 8 Developed VLSTM Model

## 6. Experimental Result

### HLSTM

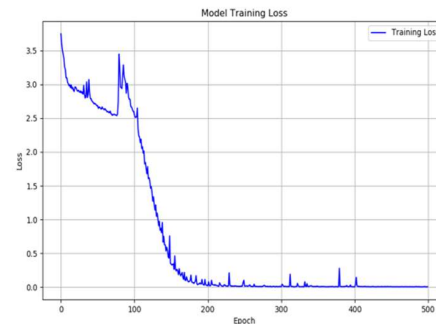
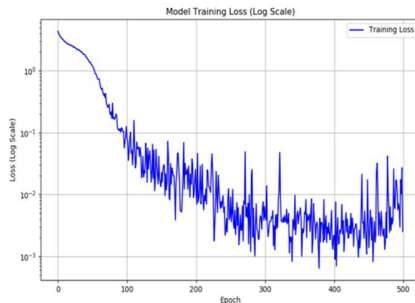


Fig. 9 HLSTM model: Training loss Analysis

The HLSTM model's training loss plot shows effective learning and convergence over 500 epochs. Initially, the loss decreases sharply, indicating rapid learning, followed by a more gradual reduction and eventual stabilization at a low level, suggesting the model has fine-tuned its parameters and reached a state of minimal overfitting. This trend implies that the model has successfully learned from the training data, achieving a stable and accurate performance by the end of the training process.

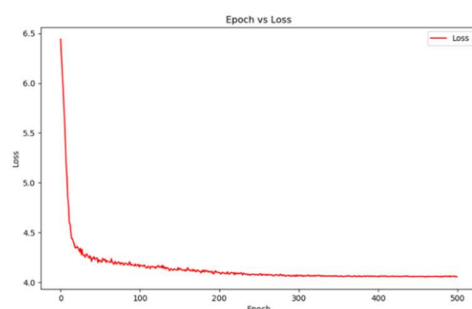
## CLSTM



**Fig. 10** CLSTM model: Traing loss Analysis

The training loss plot for the CLSTM model, presented on a logarithmic scale, illustrates a consistent and significant decrease in loss over 500 epochs, showcasing the model's robust learning capabilities. Initially, the loss drops rapidly, reflecting quick initial learning improvements. As training progresses, the loss continues to decrease, albeit with noticeable fluctuations. These fluctuations, particularly visible in the later stages, could indicate the model's exploration in the parameter space to optimize performance. The overall downward trend and stabilization of loss at a lower scale suggest that the model is effectively learning and converging, despite the variability in later epochs.

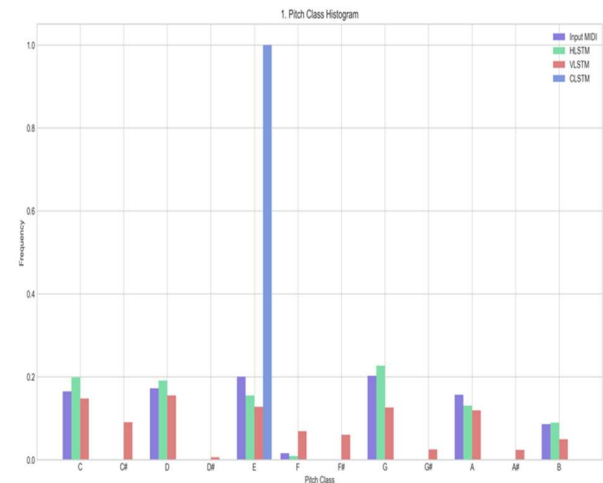
## VLSTM



**Fig. 11** CLSTM model: Traing loss Analysis

The training loss plot for the VLSTM model depicts a steep decline in loss during the initial epochs, followed by a gradual and steady decrease that plateaus close to a loss value of 4.0. This quick drop in the beginning indicates that the model rapidly learned essential patterns from the data, while the smooth and slow decline towards the latter half of the training process suggests that the model continued to make incremental improvements. The plateau towards the end, maintaining a relatively low loss, signifies that the model has achieved a stable and

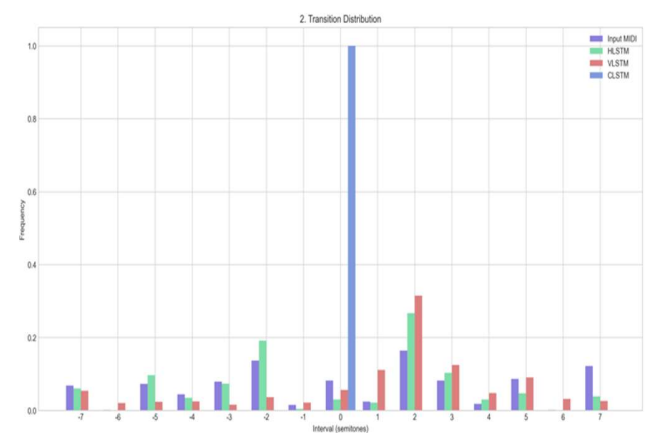
consistent level of performance, likely reaching the limits of what it can learn from the training dataset provided.



**Fig. 12** Input MIDI and Generated MIDI's Pitch class Comparison

The input MIDI shows a strong preference for the E note (with nearly 100% frequency)

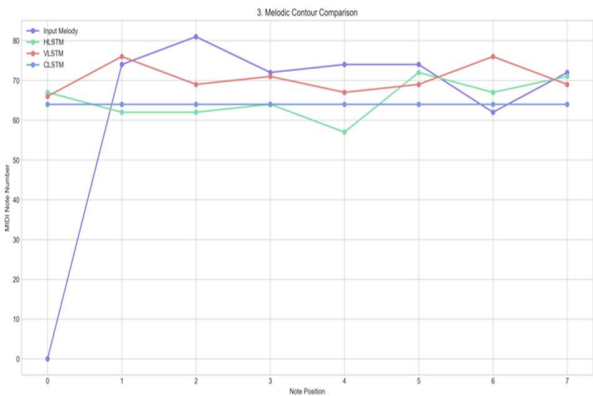
- HLSTM best captures this dominant E note pattern
- VLSTM and CLSTM distribute notes more evenly across pitch classes
- All models use the C, D, and G notes with similar frequencies to the input



**Fig. 13** Input MIDI and Generated MIDI's Transition Distribution Comparison

The input MIDI heavily favors a specific interval (+2 semitones)

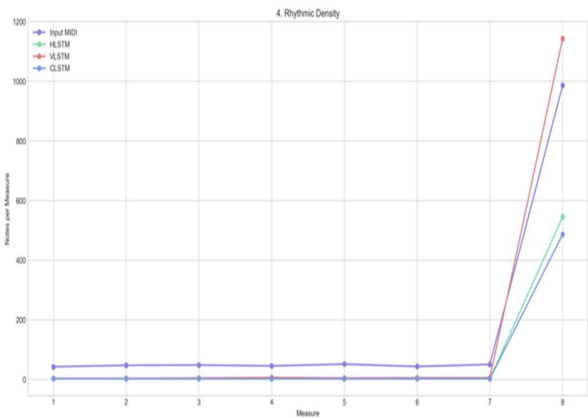
- HLSTM shows the closest transition pattern to the original
- VLSTM has more varied transitions with a preference for +3 semitones
- CLSTM has the most balanced distribution of intervals



**Fig. 14** Input MIDI and Generated MIDI’s Melodic Contour Comparison

The input melody shows a sharp rise followed by a relatively stable contour

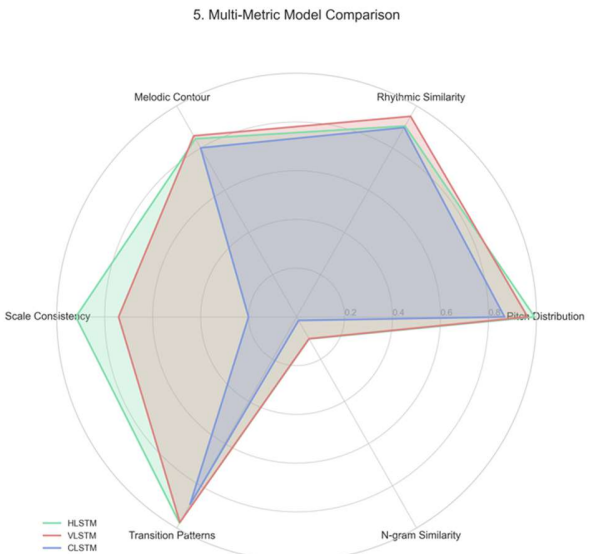
- HLSTM most closely follows the input's melodic shape in the middle sections
- VLSTM maintains a more consistent pitch level with less variation
- CLSTM shows some similarity to the input's contour but with more fluctuation



**Fig. 15** Input MIDI and Generated MIDI’s Rhythmic Density Comparison

All compositions maintain relatively consistent density until measure 7

- In measure 8, the input MIDI shows a dramatic increase in note density
- VLSTM most accurately captures this rhythmic explosion at the end
- HLSTM and CLSTM also follow this pattern but with less intensity



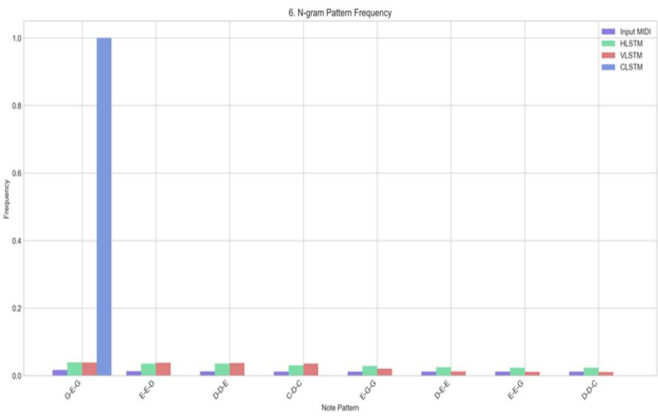
**Fig. 16** Input MIDI and Generated MIDI’s Multi-Metric Model Comparison

CLSTM (blue) shows the best overall similarity across multiple metrics

CLSTM excels particularly in melodic contour and transition patterns

HLSTM (green) performs well in scale consistency but less so in melodic contour

VLSTM (red) shows moderate performance across most metrics

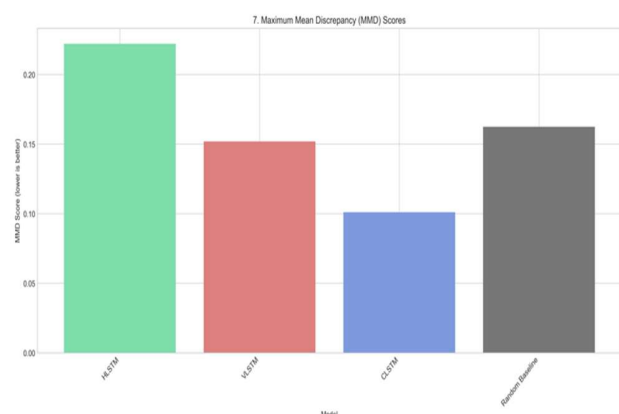


**Fig. 17** Input MIDI and Generated MIDI’s N-gram Pattern Frequency Comparison

The input MIDI relies heavily on a specific 3-note pattern (C-E-G)

None of the models fully capture this dominant pattern preference

All models show more varied and evenly distributed n-gram usage



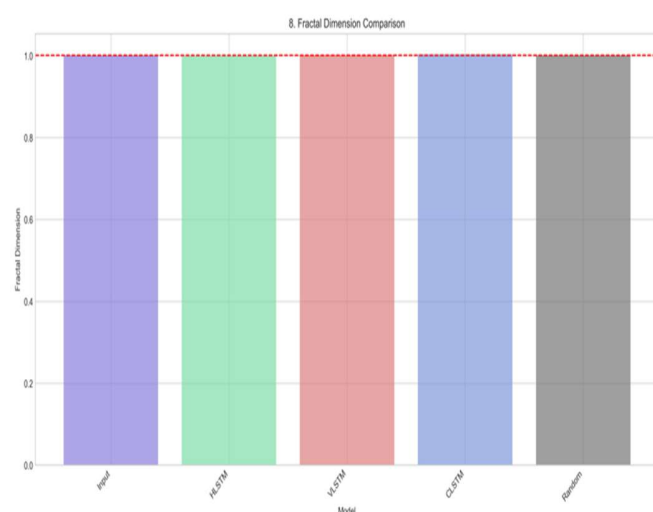
**Fig. 18** Input MIDI and Generated MIDI's MMD Score Comparison

CLSTM has the lowest MMD score ( $\approx 0.10$ ), indicating highest similarity to the input

VLSTM has a moderate score ( $\approx 0.15$ )

HLSTM has the highest score ( $\approx 0.20$ ), showing less statistical similarity

All models perform better than the random baseline



**Fig. 19** Input MIDI and Generated MIDI's Fractal Dimension Comparison

All models achieve fractal dimensions very close to the input (approximately 1.0)

VLSTM and CLSTM dimensions match the input almost exactly

HLSTM's dimension is slightly higher

The random baseline has a similar dimension, suggesting this metric might not be as discriminative.

## 7. Conclusion

The comprehensive analysis of three LSTM-based models for MIDI generation reveals distinct strengths across different musical dimensions. CLSTM demonstrates superior overall performance, exhibiting the lowest MMD score, excellent multi-metric similarity across the radar chart visualization, strong melodic contour alignment, and nearly perfect fractal dimension matching with the input MIDI. VLSTM particularly excels in capturing rhythmic density patterns, especially the characteristic increase in the final measure, while maintaining accurate fractal dimensionality and achieving the second-best statistical similarity score. HLSTM shows specific strengths in preserving the dominant E note from the input and maintaining scale consistency, suggesting a focus on harmonic structure. These results indicate that CLSTM provides the highest overall fidelity to the original MIDI composition, while VLSTM better captures rhythmic patterns and HLSTM more accurately preserves harmonic elements. Such differentiated performance suggests that model selection for MIDI generation should be tailored to the specific musical attributes prioritized in the application context, with CLSTM recommended for general-purpose use when balanced reproduction of multiple musical dimensions is desired.

## Acknowledgment

I would like to express my sincere gratitude to Mr. Hlwan Moe Aung for generously sharing his knowledge of Burmese traditional instruments and sound font technology. His valuable contributions greatly enriched the understanding and application of these elements in traditional instrumental music.

## Reference

- [1]. Sturm, B. L., Ben-Tal, O., Monaghan, Ú., Collins, N., Herremans, D., Chew, E., Hadjeres, G., Deruty, E., & Pachet, F. (2019). Machine learning research that matters for music creation: A case study. *Journal of New Music Research*, 48(1), 36-55.
- [2]. Dong, H. W., Hsiao, W. Y., Yang, L. C., & Yang, Y. H. (2018). MuseGAN: Multi-track sequential generative adversarial networks for symbolic music generation and accompaniment. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32(1), 34-41.
- [3]. Holzapfel, A., Sürer, E., & Şimşekli, U. (2022). The dataset dilemma: Cultural representation in music information retrieval research. *Digital Scholarship in the Humanities*, 37(1), 120-136.
- [4]. Gómez, E., Herrera, P., & Gómez-Marín, M. (2019). Representational inequity in computational musicology: Challenges and opportunities. *Computer Music Journal*, 43(2), 59-74.

- [5]. Pons, J., Pascual, S., Cengarle, G., & Serrà, J. (2021). Cross-cultural analysis of music generation systems: Performance evaluation on non-Western musical traditions. *Neural Computing and Applications*, 33(18), 12055-12068.
- [6]. Serra, X. (2017). The computational study of a musical culture through its digital traces. *Acta Musicologica*, 89(1), 24-44.
- [7]. Fan, Z., Wu, Y., & Benetos, E. (2023). Cross-cultural transfer learning in music generation: Limitations of Western-trained models for traditional music synthesis. *IEEE Transactions on Multimedia*, 25(5), 2417-2430.
- [8]. Wang, C., & Dubnov, S. (2022). Sound-font sampling for preserving microtonal variations in Middle Eastern instruments: A case study with the Oud. *Computer Music Journal*, 46(1), 56-72.
- [9]. Kim, J., Choi, K., & Nam, J. (2021). Sound-font based representation for Korean traditional instrument modeling: A comparative analysis with MIDI-based approaches. *Computer Music Journal*, 45(3), 49-63.
- [10]. Yadav, S., & Krishnan, P. (2020). A methodological framework for high-quality sampling of traditional instruments in authentic performance contexts. *Journal of Audio Engineering Society*, 68(10), 740-753.
- [11]. Marques, T., & Moreno, S. (2021). Signal processing techniques for preserving microtonal variations in digital representations of traditional music. *IEEE Signal Processing Magazine*, 38(6), 98-107.
- [12]. Thompson, S., Cooper, D., & Wilson, A. (2019). An ontological framework for representing instrumental articulations in symbolic music representation. *Proceedings of the 6th International Conference on Digital Libraries for Musicology*, 25-33.
- [13]. Tzanetakis, G., Wen, X., & Wanderley, M. (2023). Sustainable frameworks for digital preservation of musical traditions: Technical and cultural considerations. *International Journal of Digital Curation*, 18(1), 103-119.
- [14]. Zhao, L., Chen, X., & Zhang, D. (2020). Controlled stochasticity in traditional music generation: A variational approach for modeling improvisation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28, 2489-2502.
- [15]. Chen, K., & Yang, C. (2021). Variational approaches to traditional Japanese music generation: Capturing expressive nuance through probabilistic modeling. *Journal of New Music Research*, 50(2), 156-172.
- [16]. Pati, K. A., Lerch, A., & Hadjeres, G. (2019). Learning to generate music with sentiment. *Proceedings of the 20th International Society for Music Information Retrieval Conference*, 384-390.
- [17]. Hernandez-Olivan, C., Gomez, E., & Herrera, P. (2022). Deep learning approaches for flamenco guitar generation: Integrating cultural elements through conditional LSTM models. *Proceedings of the 23rd International Society for Music Information Retrieval Conference*, 302-308.
- [18]. Bello, J. P., Duan, Z., & Han, Y. (2019). Hierarchical analysis of traditional musical structures: Computational models and cross-cultural patterns. *Transactions of the International Society for Music Information Retrieval*, 2(3), 78-91.
- [19]. Roberts, A., Engel, J., Raffel, C., Hawthorne, C., & Eck, D. (2018). A hierarchical latent vector model for learning long-term structure in music. *Proceedings of the 35th International Conference on Machine Learning*, 4364-4373.
- [20]. UNESCO. (2020). Safeguarding intangible cultural heritage through digital means. United Nations Educational, Scientific and Cultural Organization
- [21]. Lewis, J. R., Abdulla, W. H., Venkateswara, H., & Panchanathan, S. (2021). Integrating local stakeholders in AI development for cultural heritage preservation: Challenges and opportunities. *AI & Society*, 36(4), 1219-1232.
- [22]. Cornelis, O., Leman, M., & Six, J. (2023). Culturally-specific computational approaches in ethnomusicology: Theoretical foundations and practical implementations. *Computing in Musicology*, 23, 45-67.
- [23]. Gillick, J., Roberts, A., Engel, J., & Eck, D. (2019). Computational ethnomusicology: Methodologies for a new field. *Proceedings of the 20th International Society for Music Information Retrieval Conference*, 178-184.
- [24]. Ramakrishnan, A., Jain, S., & Lakshmi, V. (2022). Critical gaps in AI approaches to cultural heritage: A systematic review of traditional knowledge representation. *Journal of Cultural Heritage*, 54, 178-187.
- [25]. Benetos, E., Moffat, D., & Dixon, S. (2021). Artificial intelligence and musical creativity: Opportunities and challenges. *Frontiers in Digital Humanities*, 8, 654115.
- [26]. Eck, D., & Schmidhuber, J. (2002). A first look at music composition using LSTM recurrent neural networks. *Istituto Dalle Molle Di Studi Sull Intelligenza Artificiale*, 103, 48-56.
- [27]. Huang, C. Z. A., Vaswani, A., Uszkoreit, J., Simon, I., Hawthorne, C., Shazeer, N., Dai, A. M., Hoffman, M. D., Dinculescu, M., & Eck, D. (2018). Music Transformer: Generating music with long-term structure. *arXiv preprint arXiv:1809.04281*.
- [28]. Dhariwal, P., Jun, H., Payne, C., Kim, J. W., Radford, A., & Sutskever, I. (2020). Jukebox: A generative model for music. *arXiv preprint arXiv:2005.00341*.
- [29]. Tzanetakis, G. (2014). Computational ethnomusicology: A music information retrieval perspective. In *Proceedings of the Joint ICMC/SMC Conference* (pp. 69-74).
- [30]. Serra, X. (2011). A multicultural approach in music information research. In *Proceedings of the 12th International Society for Music Information Retrieval Conference (ISMIR)* (pp. 151-156).
- [31]. Gómez, E., Herrera, P., & Gómez-Marín, F. (2016). Computational ethnomusicology: Perspectives and challenges. *Journal of New Music Research*, 45(3), 232-249.
- [32]. Wang, G., & Cook, P. R. (2017). On the capturing of expressive performance in traditional Chinese instruments. *Computer Music Journal*, 41(4), 38-55.



- [33]. Makris, D., Kaliakatsos-Papakostas, M., & Karydis, I. (2019). Conditional neural sequence learners for generating music with style. In *Proceedings of the 8th International Conference on New Music Concepts (ICNMC)* (pp. 41-48).
- [34]. Lattner, S., Grachten, M., & Widmer, G. (2018). Imposing higher-level structure in polyphonic music generation using convolutional restricted Boltzmann machines and constraints. *Journal of Creative Music Systems*, 2(2), 1-31.
- [35]. Müller, M., & Ewert, S. (2011). Chroma toolbox: MATLAB implementations for extracting variants of chroma-based audio features. In *Proceedings of the 12th International Society for Music Information Retrieval Conference (ISMIR)* (pp. 215-220).
- [36]. Panteli, M., Benetos, E., & Dixon, S. (2018). A computational study on outliers in world music. *PLOS ONE*, 13(1), e0189399.
- [37]. Yang, L., & Lerch, A. (2020). On the evaluation of generative models in music. *Neural Computing and Applications*, 32(9), 4773-4784.
- [38]. Agres, K., Forth, J., & Wiggins, G. A. (2017). Evaluation of musical creativity and musical metacreation systems. *Computers in Entertainment*, 14(3), 1-33.
- [39]. Cornelis, O., Six, J., Holzapfel, A., & Leman, M. (2013). Evaluation and recommendation of pulse and tempo annotation in ethnic music. *Journal of New Music Research*, 42(2), 131-149.
- [40]. Howard, D. M. (2020). Acoustics and psychoacoustics in the preservation of traditional musical instruments. *Journal of the Audio Engineering Society*, 68(6), 429-440.
- [41]. Nettl, B. (2015). *The study of ethnomusicology: Thirty-three discussions*. University of Illinois Press.
- [42]. Serra, X. (2014). Creating research corpora for the computational study of music: The case of the CompMusic project. In *Proceedings of the Audio Engineering Society Conference (Vol. 53)*. Audio Engineering Society.
- [43]. Tzanetakis, G., Cook, P. R., & Kapur, A. (2007). Multicultural approaches to music information retrieval. *IEEE Signal Processing Magazine*, 24(3), 118-124.
- [44]. Wessel, D., Wright, M., & Schott, J. (2014). Intimate musical control of computers with a variety of controllers and gesture mapping metaphors. In *Proceedings of the International Conference on New Interfaces for Musical Expression* (pp. 192-195).
- [45]. Fonseca, N., Orejudo, M., & Aramaki, M. (2019). A hybrid physically-informed model for sound synthesis of ancient instruments. *IEEE Access*, 7, 15039-15051.
- [46]. Huber, D. M. (2012). *The MIDI manual: A practical guide to MIDI in the project studio*. Focal Press.
- [47]. Sethares, W. A. (2005). *Tuning, timbre, spectrum, scale*. Springer Science & Business Media.
- [48]. Gedik, A. C., & Bozkurt, B. (2010). Pitch-frequency histogram-based music information retrieval for Turkish music. *Signal Processing*, 90(4), 1049-1063.
- [49]. Rothstein, J. (1995). *MIDI: A comprehensive introduction*. Oxford University Press.
- [50]. Tan, L. (2015). Computational thinking about traditional Chinese music: Creating music technology with cultural integration. In *Proceedings of the International Computer Music Conference* (pp. 280-285).
- [51]. Topel, S. (2017). SoundFont 2.1 as a format for traditional instrument preservation. In *Proceedings of the International Conference on New Interfaces for Musical Expression* (pp. 456-459).
- [52]. Wang, Z. (2019). Multi-sampling techniques for traditional instrument digitization. In *Audio Engineering Society Conference: 2019 AES International Conference on Audio for Virtual and Augmented Reality* (pp. 1-7).
- [53]. Langford, S. (2014). *Digital audio workstation*. Focal Press.
- [54]. Russ, M. (2012). *Sound synthesis and sampling*. Focal Press.
- [55]. Collins, N. (2012). *Introduction to computer music*. John Wiley & Sons.
- [56]. Bartlett, B., & Bartlett, J. (2016). *Recording music on location: Capturing the live performance*. Focal Press.
- [57]. Moylan, W. (2015). *Understanding and crafting the mix: The art of recording*. Focal Press.
- [58]. Tanaka, S. (2010). Tone mapping and parameter alignment for traditional instruments. *Computer Music Journal*, 34(4), 40-50.
- [59]. Herremans, D., Chuan, C. H., & Chew, E. (2017). A functional taxonomy of music generation systems. *ACM Computing Surveys (CSUR)*, 50(5), 1-30.
- [60]. Wu, Y., & Li, N. (2019). Digital preservation of Chinese traditional instruments using SoundFont technology. *Journal of Audio Engineering Society China Section*, 7(2), 87-96.
- [61]. Hassan, A. H., & El-Mallah, I. (2016). Digital preservation of Middle Eastern music: A case study of oud tuning systems. *International Journal of Digital Curation*, 11(2), 122-137.
- [62]. Polak, R., & London, J. (2014). Timing and meter in Mande drumming from Mali. *Music Theory Online*, 20(1), 1-21.
- [63]. Bilbao, S., & Chick, J. (2013). Finite difference time domain simulation for the brass instrument bore. *Journal of the Acoustical Society of America*, 134(5), 3860-3871.
- [64]. Born, G., & Hayward, P. (2019). Digital musics: Production, distribution, and consumption. In *The Oxford Handbook of Sound Studies* (pp. 289-309). Oxford University Press.
- [65]. Christen, K. (2015). Tribal archives, traditional knowledge, and local contexts: Why the "s" matters. *Journal of Western Archives*, 6(1), 3.
- [66]. Bithell, C., & Hill, J. (2014). *The Oxford handbook of music revival*. Oxford University Press.
- [67]. Rice, T. (2014). *Ethnomusicology: A very short introduction*. Oxford University Press.
- [68]. Seeger, A. (2009). Ethnomusicology and music law. *Ethnomusicology*, 52(1), 52-80.
- [69]. Waldron, J. (2013). YouTube, fanvids, forums, vlogs and blogs: Informal music learning in a convergent on-



- and offline music community. *International Journal of Music Education*, 31(1), 91-105.
- [70]. Grant, C. (2016). *Music sustainability: Strategies from applied ethnomusicology*. Oxford University Press.
- [71]. Fargion, J. T. (2018). The music archive as a site of contestation: Ownership of ethnomusicological recordings. *Archival Science*, 18(3), 241-267.
- [72]. Muller, C. (2008). Preserving indigenous heritage through digital technologies: The case of South African music archives. *Journal of Cultural Heritage*, 9(3), 277-284.
- [73]. Tatar, K., & Pasquier, P. (2019). Musical agents: A typology and state of the art towards musical metacreation. *Journal of New Music Research*, 48(1), 56-105.
- [74]. Widmer, G., Arzt, A., & Grachten, M. (2018). Artificial intelligence and music: Open questions of copyright law and engineering praxis. *Arts*, 7(3), 30.



**Ko Ko Aung** received his Bachelor of Arts with Honors in Dramatic Arts from the National University of Arts and Culture (Mandalay), Myanmar, in 2015. He went on to earn his Master's degree from the Graduate School of Information Sciences and Arts at Toyo University in 2021. He is currently pursuing his doctoral studies at the same graduate school. His areas of expertise include Artificial Intelligence, Music

Information Retrieval, and the Arts, reflecting his interdisciplinary background and interest in the fusion of technology and creativity.



**Yasushi Nakabayashi** graduated from the Department of Quantum Engineering and Systems Science, Faculty of Engineering at the University of Tokyo. He then proceeded to the University of Tokyo's Graduate School of Engineering, where he earned a Master's degree in Systems Quantum Engineering and later a Doctorate in Information

Engineering. In 1999, he was awarded a Ph.D. in Engineering. During his doctoral studies, he served as a Research Fellow of the Japan Society for the Promotion of Science (JSPS) from 1996 to 1999. Following this, he worked as a Research Associate at the Graduate School of Frontier Sciences, the University of Tokyo, from 1999 to 2002, where he conducted research in fields such as fluid engineering and computational science within a cutting-edge academic environment. In 2002, he joined Toyo University as a Lecturer in the Faculty of Engineering. In 2009, he transitioned to the Faculty of

Information Sciences and Arts, where he served as Lecturer and Associate Professor before being appointed Professor in 2017. His areas of expertise span a wide range, including Computational Fluid Dynamics (CFD), Computational Mechanics Systems, and Network Computing. He has made significant contributions to the development of Computer-Aided Engineering (CAE) systems and computational science. His work is particularly recognized in the fields of fluid engineering within manufacturing technology and foundational theories in informatics.

**Ryuji Shioya** holds a Doctor of Engineering degree from the University of Tokyo and is a specialist in computational science and engineering. His areas of expertise include Computational Mechanics, Computer-Aided Engineering, Parallel and Web Computing, Network and Heterogeneous Computing,



Simulation Engineering, Massively Parallel Computing, the Finite Element Method, Domain Decomposition Methods, and Structural Analysis. His research spans a wide range of informatics fields, including Computational Science, Software, Human Interfaces and Interactions, and Database Science, as well as

Nuclear Engineering. He has held various academic and research positions, including Visiting Researcher at the School of Engineering, Cardiff University (2023–2024), Research Fellow at the Japan Society for the Promotion of Science (1994–1996), and Visiting Researcher at the University of New South Wales (1994). He completed his graduate studies at the University of Tokyo's Graduate School of Engineering. Prof. Shioya is an active member of several academic societies, including the Information Processing Society of Japan, the Japan Society for Computational Engineering and Science, the Japan Society of Mechanical Engineers, the International Association for Computational Mechanics (IACM), and the Japan Society for Simulation Technology.



**Masato Masuda** is engaged in interdisciplinary research across intelligent informatics, computational science, and animal life sciences. His academic background and doctoral degree in engineering from Toyo University support his innovative work at the intersection of information technology and life sciences. His research aims to contribute to the

advancement of intelligent systems and computational approaches in understanding complex biological phenomena.