# A MAMBA-BASED MODEL FOR AUTOMATIC CHORD RECOGNITION

**Chunyu Yuan**
CUNY Graduate Center
Computer Science

**Johanna Devaney**
Brooklyn College and CUNY Graduate Center
Data Analysis & Visualization and Music
johanna.devaney@brooklyn.cuny.edu

## ABSTRACT

In this work, we propose a new efficient solution, which is a Mamba-based model named BMACE (Bidirectional Mamba-based network, for Automatic Chord Estimation), which utilizes selective structured state-space models in a bidirectional Mamba layer to effectively model temporal dependencies. Our model achieves high prediction performance comparable to state-of-the-art models, with the advantage of requiring fewer parameters and lower computational resources.

## 1. INTRODUCTION AND BACKGROUND

Automatic chord recognition/estimation (ACE) has a long history in music information retrieval (MIR) research [1]. While the use of modern deep-learning techniques led to major improvements [2], even the recent state-of-art approaches still experience a performance ceiling [3–5]. Some challenges in ACE that have been previously identified are the large number of label permutations [3] and disagreements between expert annotators [6–9], which is particularly true for rare chords [8, 10]. While transformer-based models (e.g., [4] excel in capturing the necessary temporal dependencies for the ACE task, they also introduce significant computational overhead due to their quadratic complexity with respect to input length. The increased complexity of transformer architectures, combined with their high memory and processing requirements, limits their usability in low-latency environments, such as real-time music analysis systems or embedded devices. Such applications call for a careful balance between model accuracy and computational efficiency.

In this paper, we evaluate how much improvement and compactness can be achieved on the ACE task by updating the model architecture, specifically by adding selective structured state-space models in a bidirectional Mamba layer. Specifically, we introduce BMACE (Bidirectional Mamba-based Network for Automatic Chord Estimation), a novel Mamba-based model that incorporates selective structured state-space models within a bidirec-

tional Mamba layer to enhance the modeling of temporal dependencies. Notably, this model achieves performance comparable to its predecessors while utilizing fewer parameters, and lower computational costs.

## 2. BI-DIRECTIONAL MAMBA NETWORK

Inspired by the bidirectional Transformer, we propose a lightweight bidirectional Mamba-based network specifically designed for chord estimation/recognition: BMACE (Bidirectional Mamba-based network for Automatic Chord Estimation). The Mamba architecture was first introduced in late 2023 [11] and has been gaining rapid momentum since its release. It has been applied to some speech [12–14] and some MIR [15, 16] tasks, but not yet for ACE. Mamba distinguishes itself from other models by eschewing the usual attention and MLP blocks for a more streamlined approach. This results in a model that is not only lighter and faster but also uniquely capable of scaling linearly with sequence length, an achievement that sets it apart from its predecessors. Central to Mamba's design are its Selective-State-Spaces (SSM): these are recurrent models that selectively process information based on the current input, effectively filtering out irrelevant data to focus on what is most critical for efficient processing. Additionally, Mamba simplifies its architecture by replacing the complex attention and MLP blocks in Transformers with a single, unified SSM block, enhancing inference speed and reducing computational load. Mamba incorporates hardware-aware parallelism, using a specially designed parallel algorithm that optimizes recurrent operations for improved hardware efficiency, potentially boosting performance even further.

Figure 1 shows the three variants of our Mamba-based model that we experiment with. The first (MACE-V) is a vanilla Mamba model with two vertical Mamba layers. The second (MACE-H) is a Mamba model with two concatenated/horizontal models. The third is our proposed model, BMACE, which presents the structure of our bidirectional Mamba network. Bidirectional Mamba blocks and fully-connected layers are the main modules in the network. It processes a 10-second audio signal as a Constant Q Transform (CQT) feature. The model integrates a fully-connected layer into the input, which then proceeds to two Mamba blocks with opposite masking directions, represented as dotted boxes in Figure 1. The outputs from these blocks are concatenated and passed through a fully-
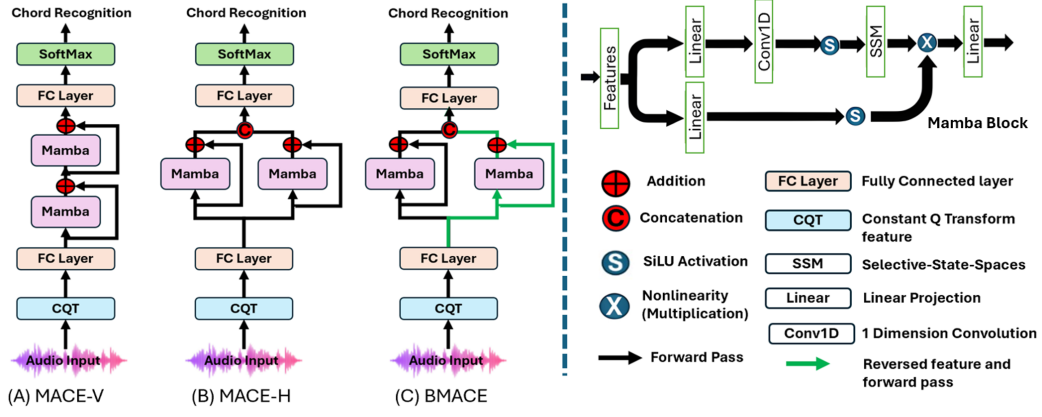
**Figure 1**. Architecture of Chord Recognition Models with Mamba Block. The diagram illustrates three variants, MACE-V (A), MACE-H (B), and BMACE (C). Each utilizes the Mamba block for improved feature processing. The Mamba block employs selective-state-spaces (SSM), SiLU activation, and 1D convolutions (Conv1D) for feature transformation. The figure highlights forward pass operations, feature reversal, addition, and concatenation mechanisms in the respective models, with fully connected (FC) layers leading to SoftMax output for chord recognition.

| Model | maj-min label type | | | | large vocabulary label type | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Root↑ | Maj-min↑ | GFlops↓ | Params↓ | Root↑ | Thirds↑ | Triads↑ | Sevenths↑ | Tetrads↑ | Maj-min↑ | MIREX↑ | GFlops↓ | Params↓ |
| CRNN [3] | 0.8185 | **0.7796** | 0.0957 | 435,609 | 0.8026 | 0.7459 | 0.6384 | 0.6426 | 0.5448 | **0.7544** | 0.7543 | 0.1038 | 472,874 |
| BTC [4] | 0.8202 | 0.7628 | 0.6282 | 2,910,361 | **0.8051** | **0.7524** | **0.6469** | 0.6506 | **0.5604** | 0.7531 | 0.7486 | 0.6362 | 2,929,066 |
| MACE-V | 0.7920 | 0.7309 | **0.0247** | **111,161** | 0.7739 | 0.715 | 0.6084 | 0.6137 | 0.5242 | 0.7166 | 0.7057 | **0.0328** | **129,866** |
| MACE-H | 0.7898 | 0.7347 | 0.0261 | 114,361 | 0.7833 | 0.7211 | 0.6188 | 0.6236 | 0.5314 | 0.7304 | 0.7238 | 0.0422 | 151,626 |
| BMACE | **0.8212** | 0.7678 | 0.0261 | 114,361 | 0.8043 | 0.7455 | 0.6426 | **0.6526** | 0.5571 | 0.7536 | **0.7595** | 0.0422 | 151,626 |

**Table 1**. Weighted Chord Symbol Recall (WCSR) scores for the performance of the CRNN [3], BTC [4], and our three Mamba variants (MACE-V, MACE-H, and BMACE) on the uspop2002 dataset.

connected layer to maintain the input's original dimensions. We added residual operation in the blocks and layers to increase the information entropy.

## 3. EXPERIMENT

### 3.1 Experiment Setting

Our models are implemented with Pytorch [17] framework. All experiments are conducted on the instance node at Lambda [1] that has a single NVIDIA RTX A6000 GPU (24 GB), 14vCPUs, 46 GiB RAM and 512 GiB SSD. Our model was trained and validated on the MARL annotations [2] of uspop2002 dataset [18]. Each 10-second audio signal was processed with a 5-second overlap between consecutive signals. The signals were sampled at 22,050 Hz and analyzed using a Constant Q Transform (CQT) that covered 6 octaves starting from C1, with 24 bins per octave and a hop size of 2048. The CQT features were then converted to log amplitude using the formula $S_{\log} = \ln(S + \epsilon)$, where S represents the CQT feature, and $\epsilon$ is an extremely small number. This was followed by the application of

global z-normalization, using the mean and variance derived from the training data.

We evaluate the three versions of our model described in Section 2 (MACE-V, MACE-H, and BMACE) against state of the art CRNN-based [3] and transformer-based [4] models on the 25-label maj-min and the 170-label large chord vocabularies.

### 3.2 Results

Table 1 presents the model validation results. BMACE performs slightly better on some label types than the CRNN and BTC models, though the difference is not likely to be statistically significant. However, there is a notable improvement over the non-bidirectional Mamba models (MACE-V and MACE-H). The most significant differences are in the size and processing requirements of the various models, as shown in Table 1. As previously observed, CRNNs are more efficient than transformer-based models, and the claim that Mamba networks are also more efficient holds true. All Mamba-based models use only 1/25th of the parameters of the transformer-based BTC model and are slightly less than 1/3 smaller than CRNN. This reduction in the number of parameters is reflected in the lower GFlops required to run the model.

---

[1] https://cloud.lambdalabs.com/instances
[2] https://github.com/tmc323/Chord-Annotations

## 4. ACKNOWLEDGMENTS

## 5. REFERENCES

[1] J. Pauwels, K. O'Hanlon, E. Gómez, and M. B. Sandler, "20 years of automatic chord recognition from audio," in *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, Delft, The Netherlands, 2019, pp. 54–63.

[2] F. Korzeniowski and G. Widmer, "A fully convolutional deep auditory model for musical chord recognition," in *2016 IEEE 26th International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE, 2016, pp. 1–6.

[3] B. McFee and J. P. Bello, "Structured training for large-vocabulary chord recognition," in *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, Suzhou, China, 2017, pp. 188–194.

[4] J. Park, K. Choi, S. Jeon, D. Kim, and J. Park, "A Bi-Directional Transformer for Musical Chord Recognition," in *Proceedings of the International Society for Music Information Retrieval Conference, ISMIR*, A. Flexer, G. Peeters, J. Urbano, and A. Volk, Eds., 2019, pp. 620–627. [Online]. Available: http://archives.ismir.net/ismir2019/paper/000075.pdf

[5] T. Ito and S. Arai, "Harmonic Representation for CNN-LSTM Automatic Chord Recognition," in *International Conference on Cybernetics and Intelligent System (ICORIS)*. IEEE, 2021, pp. 1–5.

[6] E. J. Humphrey and J. P. Bello, "Four timely insights on automatic chord estimation," in *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, Málaga, Spain, 2015, pp. 673–679.

[7] N. Condit-Schultz, Y. Ju, and I. Fujinaga, "A flexible approach to automated harmonic analysis: Multiple annotations of chorales by Bach and Prætorius." in *Proceedings of International Society for Music Information Retrieval Conference (ISMIR)*, 2018, pp. 66–73.

[8] H. V. Koops, W. B. de Haas, J. A. Burgoyne, J. Bransen, A. Kent-Muller, and A. Volk, "Annotator subjectivity in harmony annotations of popular music," *Journal of New Music Research*, vol. 48, no. 3, pp. 232–252, 2019.

[9] H. Koops, W. de Haas, J. Bransen, and A. Volk, "Chord label personalization through deep learning of integrated harmonic interval-based representations," in *Proceedings of the First International Conference on Deep Learning and Music*, 2017, pp. 19–25.

[10] Y. Ni, M. McVicar, R. Santos-Rodríguez, and T. D. Bie, "Understanding effects of subjectivity in measuring chord estimation accuracy," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 12, pp. 2607–2615, 2013.

[11] A. Gu and T. Dao, "Mamba: Linear-time sequence modeling with selective state spaces," *arXiv preprint arXiv:2312.00752*, 2023.

[12] X. Jiang, C. Han, and N. Mesgarani, "Dual-path mamba: Short and long-term bidirectional selective structured state space models for speech separation," *arXiv preprint arXiv:2403.18257*, 2024.

[13] K. Li and G. Chen, "Spmamba: State-space model is all you need in speech separation," *arXiv preprint arXiv:2404.02063*, 2024.

[14] C. Quan and X. Li, "Multichannel long-term streaming neural speech enhancement for static and moving speakers," *arXiv preprint arXiv:2403.07675*, 2024.

[15] J. Bai, Y. Fang, J. Wang, and X. Zhang, "A two-stage band-split mamba-2 network for music separation," *arXiv preprint arXiv:2409.06245*, 2024.

[16] J. Chen, T. Xie, X. Tang, J. Wang, W. Dong, and B. Shi, "Musicmamba: A dual-feature modeling approach for generating chinese traditional music with modal precision," *arXiv preprint arXiv:2409.02421*, 2024.

[17] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems*, vol. 32, 2019.

[18] A. Berenzweig, B. Logan, D. P. Ellis, and B. Whitman, "A large-scale evaluation of acoustic and subjective music-similarity measures," *Computer Music Journal*, pp. 63–76, 2004.