

# BEYOND THE SAMPLE ALONE: FAMILIARITY AND CONTEXT IN HUMAN EVALUATION OF GENERATED MUSIC

Anonymous Authors  
Anonymous Affiliation  
anonymous@ksmi.kr

## ABSTRACT

Human evaluation of generated music is often modeled as a direct function of sample-level representations, yet it remains unclear how much of subjective judgment can actually be explained by the generated sample alone. To investigate how listeners evaluate generated music, we conducted a continuation-based listening test using 200 piano continuations generated from shared primers across ten classical composer styles, rated by 190 listeners on six rating dimensions and Familiarity. Sample-only predictions captured general rating trends at the aggregate level, but performed poorly for individual listeners and failed to generalize to unseen composers. Variance decomposition further revealed that the unexplained portion of ratings was largely attributable to individual listener differences. We therefore examined Familiarity as a key interpretable listener-side factor. While Familiarity showed a strong style-level baseline across composer contexts, its explanatory role for ratings was driven primarily by continuation-level familiarity, and it also modulated how relative musical change was evaluated. These results suggest that evaluation of generated music is layered: beyond signal-level properties, human judgments reflect structured listener-side responses that are not fully captured by coarse similarity or sample-only representations.

## 1. INTRODUCTION

Evaluating generated music remains a central challenge in Music Information Retrieval (MIR). Recent work has largely framed this problem as one of predicting human judgments from sample-level representations, whether through learned embeddings, per-sample metrics, or benchmark datasets with listener ratings [1–3]. This framing is useful for comparing systems at scale, but it also raises a basic question of how much of human evaluation can actually be explained by the generated sample alone.

This question is particularly important for music, where judgments are often context-dependent and listener-dependent. Prior studies in MIR have shown that human annotations can exhibit structured disagreement rather than pure noise [4–6]. Related research in music cognition further suggests that listener-dependent factors such as familiarity, expectation, and processing fluency can shape aes-

thetic response [7, 8]. Among these, familiarity is especially relevant because it provides an interpretable link between stylistic context and listener response. Yet familiarity has rarely been modeled in a structured way in generated music evaluation. We therefore treat familiarity as one listener-side explanatory variable and ask whether it helps account for part of the variation that sample-only representations leave unexplained.

In this paper, we address this gap through a continuation-based evaluation design. Instead of comparing unrelated clips, we construct a dataset around shared composer-style primers. This makes it possible to examine evaluation at three levels: across stylistic contexts, across continuations within the same context, and across listeners rating the same continuation. More specifically, we use this setting to connect three questions that are often studied separately: what sample-only prediction captures, what structured variation remains unexplained, and whether part of that variation can be interpreted through familiarity and context-relative change.

## 2. EXPERIMENTAL SETUP

We conducted a continuation-based listening experiment using piano continuations generated from shared primer excerpts. Ten four-bar piano primers were selected from MAESTRO dataset [9] to provide distinct stylistic baselines drawn from classical composers while avoiding highly recognizable excerpts. For each primer, 20 continuations were generated with a pre-trained MusicTransformer [10], yielding a final set of 200 continuations. This controlled shared-context design allowed us to distinguish broad style-level effects from continuation-level variation while minimizing inter-model variability and broader genre-level preference confounds.

Subjective ratings were collected through an online listening study with 190 participants. Each participant rated 20 samples in a balanced design such that two continuations were heard from each primer context. After listening, participants rated each sample on 7-point Likert scales for Overall, Melodiousness, Creativity, Rhythmicity, Coherence, Naturalness, and Familiarity. No information was provided about the use of AI generation.

For signal-based modeling, we extracted MERT [11] audio embeddings from the full sequence, primer, and continuation segments. These were projected into a shared low-dimensional space via PCA, yielding full-sequence

88 and relative-change (continuation minus primer) compo-  
 89 nents for use in subsequent models. Analyses combined  
 90 ridge regression for sample-only prediction with mixed-  
 91 effects models for variance decomposition, familiarity de-  
 92 composition, and interaction testing.

### 93 3. RESULTS

94 Human ratings were multidimensional, with a consistent  
 95 evaluative structure across rating dimensions. Overall was  
 96 broadly aligned with Melodiousness, Coherence, and Nat-  
 97 uralness, while Creativity and Familiarity were more dis-  
 98 tinct from the main rating axis. This motivated treating  
 99 Familiarity not simply as another rating dimension, but  
 100 as a listener-side variable potentially explaining variation  
 101 across ratings.

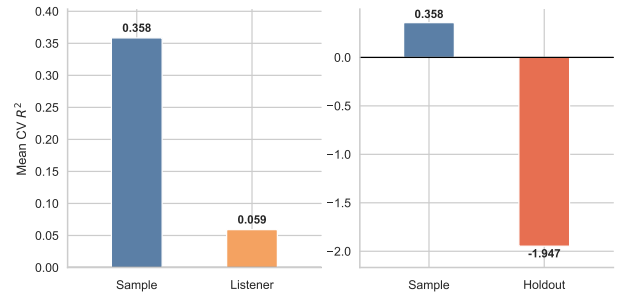
102 We first tested whether ratings could be predicted from  
 103 the continuation sample alone. As shown in Figure 1,  
 104 sample-only models captured part of the average rating  
 105 tendency: predictions aligned moderately with aggregated  
 106 sample ratings, but performance dropped substantially  
 107 when the same predictions were compared with individual  
 108 listener ratings. Generalization also deteriorated sharply  
 109 under composer-holdout evaluation. Thus, sample-level  
 110 representations recovered some broad ranking of continua-  
 111 tions, but not the fuller structure of ratings across listeners  
 112 and stylistic contexts.

113 To characterize what remained unexplained, we fit null  
 114 mixed-effects models that decomposed variance across  
 115 composer context, continuation identity within composer,  
 116 listener identity, and residual variance. Variance decom- 145  
 117 position showed that listener-level variance accounted for 146  
 118 a substantial share of total variance across the six rating 147  
 119 dimensions. In Overall ratings, for example, listener iden- 148  
 120 tity explained a larger share of variance than either com- 149  
 121 poser context or continuation identity alone. Importantly, 150  
 122 this should not be interpreted as evidence for Familiarity 151  
 123 alone: listener-side variance may reflect multiple factors, 152  
 124 including familiarity, preference, musical background, or 153  
 125 response style. However, it does show that the unexplained  
 126 portion of evaluation is structured rather than purely noisy. 154

127 We next examined how Familiarity itself was patterned  
 128 across contexts. Familiarity showed a strong style-level 155  
 129 baseline, and composer contexts that received higher Over- 156  
 130 all ratings also tended to be rated as more familiar. How- 157  
 131 ever, it was not reducible to composer context alone. 158  
 132 Within the same stylistic context, continuations still varied 159  
 133 substantially in perceived Familiarity, and the Familiarity- 160  
 134 Overall relation differed across composers. 161

135 We then asked whether this variation in Familiarity 162  
 136 helped explain the other ratings. As summarized in Ta- 163  
 137 ble 1, adding raw Familiarity improved model fit sub- 164  
 138 stantially over the sample-only baseline, and decomposed 165  
 139 models showed that continuation-level familiarity pro- 166  
 140 vided the clearest contribution across most rating dimen- 167  
 141 sions. Thus, raw Familiarity improved prediction, but its 168  
 142 contribution became more interpretable when decomposed 169  
 143 into style-level and continuation-level components. 170

144 Finally, Familiarity also shaped how relative musical 171



**Figure 1.** Performance of the sample-only prediction. Left: mean cross-validated  $R^2$  at the sample and listener levels. Right: mean  $R^2$  for sample-level CV and composer holdout.

Outcome	Sample	+Raw	+Decomp	Key component
Overall	0.09	0.26	0.27	continuation
Melodiousness	0.09	0.29	0.31	continuation
Creativity	0.03	0.03	0.07	listener
Rhythmicity	0.07	0.21	0.22	continuation
Coherence	0.11	0.31	0.31	continuation
Naturalness	0.07	0.24	0.24	continuation

**Table 1.** Marginal  $R^2$  of sample-only, raw-Familiarity, and decomposed-Familiarity models across the six rating dimensions. The final column indicates the familiarity component showing the clearest unique contribution in nested component comparisons.

change was evaluated. Interaction analyses revealed that the relationship between relative-change axes and ratings depended on continuation-level familiarity. In some cases, the same change dimension tended to be evaluated more positively when familiarity was high than when it was low; in others, familiarity attenuated the strength of the change effect. These patterns suggest that Familiarity influences not only rating magnitude but also how listeners interpret the continuation relative to its preceding context.

### 4. DISCUSSION AND CONCLUSION

Taken together, these findings suggest that generated music evaluation is layered. Sample-only representations capture part of the broad rating pattern, but a substantial portion of rating variation remains structured across contexts, continuations, and especially listeners. Familiarity provides one interpretable listener-side factor that explains part of this remaining variation.

Importantly, we do not argue that Familiarity exhaustively explains human evaluation. Rather, our results show that listener-side structure matters and that continuation-based designs provide a useful way to study it. For MIR, this implies that improving evaluation is not only a matter of learning better sample-level representations, but also of accounting for how listeners interpret musical events relative to stylistic context. This provides a useful first step toward context-sensitive and user-aware evaluation frameworks for generated music.

## 5. References

- 173 [1] C. Liu, H. Wang, J. Zhao, S. Zhao, H. Bu, X. Xu,  
174 J. Zhou, H. Sun, and Y. Qin, “Musiceval: A generative  
175 music dataset with expert ratings for automatic text-to-  
176 music evaluation,” in *ICASSP 2025-2025 IEEE Inter-  
177 national Conference on Acoustics, Speech and Signal  
178 Processing (ICASSP)*. IEEE, 2025, pp. 1–5.
- 179 [2] G. Ma, Y. Xia, J. Yao, H. Xue, H. Liu, S. Wang,  
180 H. Liu, and L. Xie, “The icassp 2026 automatic  
181 song aesthetics evaluation challenge,” *arXiv preprint  
182 arXiv:2601.07237*, 2026.
- 183 [3] A. Lerch, C. Arthur, N. Bryan-Kinns, C. Ford, Q. Sun,  
184 and A. Vinay, “Survey on the evaluation of generative  
185 models in music,” *ACM Computing Surveys*, vol. 58,  
186 no. 4, pp. 1–36, 2025.
- 187 [4] N. L. Masclef, A. Vaglio, and M. Moussallam, “User-  
188 centered evaluation of lyrics-to-audio alignment.” in  
189 *ISMIR*, 2021, pp. 420–427.
- 190 [5] A. Flexer, T. Lallai, and K. Rašl, “On evaluation of  
191 inter-and intra-rater agreement in music recommenda-  
192 tion,” *Transactions of the International Society for Mu-  
193 sic Information Retrieval*, vol. 4, no. 1, 2021.
- 194 [6] J. S. Gómez-Cañón, E. Cano, P. Herrera, and  
195 E. Gómez, “Joyful for you and tender for us: The in-  
196 fluence of individual characteristics and language on  
197 emotion labeling and classification,” in *Proceedings of  
198 the 21st International Society for Music Information  
199 Retrieval Conference*, 2020, pp. 853–860.
- 200 [7] A. Korkor, V. Noreika, C. Di Bernardi Luft, and  
201 M. Pearce, “Relationships between surprise, liking,  
202 and error perception in musical listening.” *Psychology  
203 of Aesthetics, Creativity, and the Arts*, 2025.
- 204 [8] E. Schubert, D. J. Hargreaves, and A. C. North, “A  
205 dynamically minimalist cognitive explanation of mu-  
206 sical preference: is familiarity everything?” *Frontiers  
207 in psychology*, vol. 5, p. 38, 2014.
- 208 [9] C. Hawthorne, A. Stasyuk, A. Roberts, I. Si-  
209 mon, C.-Z. A. Huang, S. Dieleman, E. Elsen,  
210 J. Engel, and D. Eck, “Enabling factorized piano  
211 music modeling and generation with the MAE-  
212 STRO dataset,” in *International Conference on  
213 Learning Representations*, 2019. [Online]. Available:  
214 <https://openreview.net/forum?id=r11YRjC9F7>
- 215 [10] H. C.-Z. Anna, V. Ashish, U. Jakob, S. Ian, H. Curtis,  
216 S. Noam, D. Monica, E. Douglas *et al.*, “Music trans-  
217 former: Generating music with long-term structure,”  
218 *arXiv preprint*, 2018.
- 219 [11] Y. Li, R. Yuan, G. Zhang, Y. Ma, X. Chen, H. Yin,  
220 C. Xiao, C. Lin, A. Ragni, E. Benetos *et al.*,  
221 “Mert: Acoustic music understanding model with  
222 large-scale self-supervised training,” *arXiv preprint  
223 arXiv:2306.00107*, 2023.