KAD: No More FAD! An Effective and Efficient Evaluation Metric for Audio Generation

Yoonjin Chung*1, Pilsun Eu*1, Junwon Lee2, Keunwoo Choi^{1,3}, Juhan Nam2, Ben Sangbae Chon¹
¹Gaudio Lab Inc., ²KAIST, ³Genentech

Abstract

Although being widely adopted for evaluating generated audio signals, the Fréchet Audio Distance (FAD) suffers from significant limitations, including reliance on Gaussian assumptions, sensitivity to sample size, and high computational complexity. As an alternative, we introduce the Kernel Audio Distance (KAD), a novel, distribution-free, unbiased, and computationally efficient metric based on Maximum Mean Discrepancy (MMD). Through analysis and empirical validation, we demonstrate KAD's advantages: (1) faster convergence with smaller sample sizes, enabling reliable evaluation with limited data; (2) lower computational cost, with scalable GPU acceleration; and (3) stronger alignment with human perceptual judgments. By leveraging advanced embeddings and characteristic kernels, KAD captures nuanced differences between real and generated audio. Open-sourced in the kadtk² toolkit, KAD provides an efficient, reliable, and perceptually aligned benchmark for evaluating generative audio models.

Index Terms: audio generation, audio quality evaluation, Fréchet audio distance, maximum mean discrepancy, kernel methods

1 Introduction

Distribution-free?		Bias-free?	Computation Cost	
FAD [1]	X	X	$\mathcal{O}(dN^2+d^3)$	
KAD (Ours)	✓	✓	$\mathcal{O}(dN^2)$	

Table 1: Comparison of FAD and KAD: KAD is distribution-free, unbiased, and computationally efficient, even for high embedding dimensions ($d \le 2048$) and large sample sizes ($N \le 10k$).

As the demand for neural audio generation continues to grow across various domains such as content creation and virtual environments, innovative models are emerging to address a wide range of tasks. These include generating audio from textual descriptions, visual inputs, temporal data, or other audio signals, underscoring the importance of models that can process diverse types of inputs. Consequently, the need for robust and reliable methods to evaluate the quality of these models is becoming increasingly critical.

The Fréchet Audio Distance (FAD) [1] is a widely used metric for evaluating the overall performance of audio generation models, measuring the dissimilarity between the statistical distributions of real and generated audio samples. FAD is considered a simple yet effective measure for objective evaluation, making it a popular choice for assessing generative audio models across various tasks.

^{*}Equal contribution

²https://github.com/YoonjinXD/kadtk

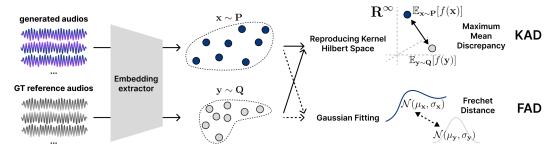


Figure 1: Comparison between KAD (Kernel Audio Distance) and FAD (Fréchet Audio Distance). KAD is a distribution-free metric that does not require any underlying assumptions for embedding distributions P and Q.

However, FAD has significant shortcomings that limit its effectiveness as a benchmark for generative audio models. First, it relies on the assumption that audio embeddings follow a Gaussian distribution, which often does not apply to real-world audio data with complex and diverse characteristics. Second, FAD suffers from an inherent bias in finite-sample estimation, which produces unreliable results particularly with smaller datasets. Finally, the computational cost of FAD is substantial, due to its high dependence on the embedding dimension size. These issues collectively challenge the practical use of FAD in evaluating modern generative audio models [2–5].

Motivated by these limitations of FAD, we propose the Kernel Audio Distance (KAD), a novel evaluation metric for audio generation. The proposed metric leverages the Maximum Mean Discrepancy (MMD), a non-parametric measure that makes no assumptions about the underlying distribution of sample embeddings such as the normality assumption in FAD. A comparison between KAD and FAD is illustrated in Figure 1. In summary, our contributions are:

- We propose KAD, a novel metric based on MMD for evaluating generative audio models.
- We provide empirical evidence demonstrating faster convergence with sample size, computational advantages, and stronger rank correlation with human evaluations of KAD.
- We provide guidelines to establish the practical applicability of KAD for selecting kernels, parameters, and embedding models to ensure consistent and reliable audio quality evaluation.

The implementation of KAD is provided as an open-source toolkit named kadtk (for more detail, refer to Appendix C).

2 Related Works and Preliminaries

2.1 Fréchet Audio Distance and its Limitations

The Fréchet Audio Distance (FAD) [1] measures the difference between two sets of audio samples within their data embedding space. Specifically, it is an estimation of the Fréchet distance between the underlying distributions of two given embedding sample sets. FAD is an adaptation of the Fréchet Inception Distance (FID)[6] – originally proposed for evaluating image generation models – to the audio domain. The embeddings are typically extracted using an audio encoder model pretrained on real-world data such as VGGish[7], ensuring that the embeddings capture representative features of the audio samples for reliable evaluation.

Given the ground-truth reference set embeddings $X = \{x_i\}_{i=1}^n$ and the target evaluation set $Y = \{y_j\}_{j=1}^m$, FAD is defined by:

$$FAD^{2}(X,Y) = \|\mu_{X} - \mu_{Y}\|_{2}^{2} + tr\left(\Sigma_{X} + \Sigma_{Y} - 2\sqrt{\Sigma_{X}\Sigma_{Y}}\right), \tag{1}$$

where X and Y are assumed to be sampled from multivariate Gaussian distributions, fully characterized by their means μ_X, μ_Y and covariances Σ_X, Σ_Y .

FAD is a conventional choice of metric for evaluating generative models in various domains, including text-to-audio [8–20] and vision-to-audio [21–36] tasks, and is considered one of the standards for

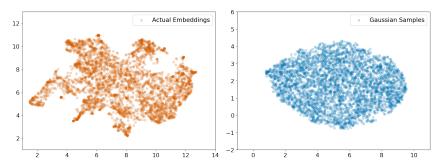


Figure 2: The left panel is the complex, non-Gaussian shape of the VGGish embeddings from the Clotho [38] train dataset, while the right panel depicts its assumed Gaussian distribution.

generative performance. Despite its popularity, FAD has three crucial limitations that undermine its efficacy and efficiency:

- 1. **Normality assumption**: The assumption that audio embeddings follow a Gaussian distribution often fails for real-world data. Such data are often asymmetrically distributed or unevenly clustered, which compromises the metric's ability to accurately measure similarity. Similar limitations have been observed in the computer vision domain, where the normality assumption was shown to be unsuitable for the Inception [37] embedding space used for FID [2]. Figure 2 shows how the actual distribution of VGGish embeddings from the Clotho dataset (reduced via UMAP) differs significantly from the samples drawn from a Gaussian distribution with the same mean and covariance. This deviation highlights the limitations of the normality assumption for representing real-world audio data, leading to potential inaccuracies and over- or under-estimations in FAD values.
- 2. Sample size bias: FAD is an inherently biased metric [5] which requires a large number of audio samples for a reliable result. Given the embedding sample size $N = \max(n, m) < \infty$, The bias in FAD from the finite-sample estimation of sample covariance increases as $\mathcal{O}(1/N)$ (for a detailed derivation, refer to Appendix B). Similar findings have also been reported for FID [3]. This necessitates the use of larger datasets for more accurate evaluations, which is particularly undesirable in the audio domain where high-quality data is relatively scarce compared to image datasets. Furthermore, this sample size bias creates the potential for manipulation: increasing N can artificially reduce bias, leading to better FAD scores and the appearance of improved performance. Naturally, this further undermines the credibility of FAD as a reliable evaluation metric.
- 3. **High computational cost**: The time complexity of FAD is given by $\mathcal{O}(dN^2+d^3)$, which scales poorly with the number of dimensions of the embedding space. This makes the calculations cumbersome when using audio embedding models that produce high-dimensional embeddings. Moreover, the calculation of square-roots of covariance matrices in Equation 1 is not easily parallelized, limiting the utilization of parallel computing.

2.2 Maximum Mean Discrepancy

To address the limitations of FAD, we adopt the Maximum Mean Discrepancy (MMD) [39]. Originally proposed for a statistical test to distinguish whether two samples come from the same distribution, MMD is capable of capturing differences not only in mean and variance but also in higher-order moments. It is also distribution-free, meaning that it does not assume that the samples belong to a specific family of distributions (e.g. Gaussian). This allows for a more comprehensive comparison of how two sets of audio samples differ in their embedding spaces.

The MMD between two distributions P and Q is defined as:

$$MMD(\mathcal{F}, P, Q) = \sup_{f \in \mathcal{F}} \Big(\mathbb{E}_{x \sim P}[f(x)] - \mathbb{E}_{y \sim Q}[f(y)] \Big), \tag{2}$$

where \mathcal{F} is a class of functions chosen to detect differences between P and Q.

When \mathcal{F} is chosen to be the Reproducing Kernel Hilbert Space (RKHS) induced by a kernel function $k(\cdot,\cdot)$, calculating the MMD corresponds to measuring the Euclidean distance between the mean embedding positions after mapping the data into a high-dimensional feature space. The high-dimensional mapping is necessary because it reveals the nonlinear differences between two distributions that may not be apparent in a lower-dimensional setting.

Rather than computing these high-dimensional representations explicitly, kernel operations can be used to calculate the distances in the RKHS directly in the original embedding space. This technique, often referred to as the "kernel trick," allows the metric to leverage high-dimensional – or even infinite-dimensional – representations that would otherwise be infeasible to compute. With this setup, the MMD can be computed entirely through pairwise comparisons of the samples:

$$MMD^{2}(P,Q) = \mathbf{E}_{x,x'}[k(x,x')] + \mathbf{E}_{y,y'}[k(y,y')] - 2\mathbf{E}_{x,y}[k(x,y)],$$
(3)

where x, x' are drawn from P and y, y' are drawn from Q.

For the finite samples in the reference set $X = \{x_i\}_{i=1}^n$ and the evaluation set $Y = \{y_j\}_{j=1}^m$, an unbiased estimator of 3 is:

$$\widehat{\text{MMD}}_{\text{unbiased}}^{2}(X,Y) = \frac{1}{n(n-1)} \sum_{i \neq j} k(x_{i}, x_{j}) + \frac{1}{m(m-1)} \sum_{i \neq j} k(y_{i}, y_{j}) - \frac{2}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} k(x_{i}, y_{j}).$$
(4)

One of the most widely used MMD-based metrics for evaluating generative models is the Kernel Inception Distance (KID) proposed by Bińkowski et al. [40], as the squared MMD value between Inception embeddings of images with a cubic polynomial kernel. KID was used as an evaluation metric in several audio-related works [41–44], particularly in music generation. For image quality, Jayasumana et al. [2] proposed the CLIP Maximum Mean Discrepancy (CMMD) for evaluation on CLIP embeddings [45] and a Gaussian kernel, with Novack et al. [46] following similar practices but on CLAP [47] embeddings. Many of these works acknowledge the potential advantages of MMD. However, to our knowledge, there has been no comprehensive study of its reliability in comparison to FAD or the effect of various design choices.

3 Kernel Audio Distance

In this section, we propose the *Kernel Audio Distance* (KAD), a reliable and computationally efficient metric for evaluating audio generation models.

The fundamental requirement for such a metric is its ability to capture perceptually meaningful differences between generated and reference audio. Provided that the embedding space sufficiently encodes these perceptually relevant features, a reliable metric must be capable of accurately comparing embedding distributions without imposing restrictive assumptions. To this end, we adopt the MMD, whose distribution-free nature eliminates the need for parametric assumptions and enables a comprehensive comparison between two embedding distributions.

We define KAD as follows:

$$KAD = \alpha \cdot \widehat{MMD}_{\text{unbiased}}^{2}, \tag{5}$$

where α is a resolution scaling factor introduced for convenient score comparison. We set $\alpha=100$ as the default.

3.1 The Strengths of KAD

Along with its robust theoretical foundation, KAD also provides key practical advantages over FAD:

Unbiased Nature: The KAD score is independent of the sample size, making it robust to smaller samples without employing bias-correction procedures. By contrast, the bias-correction for FAD (e.g., FAD_{∞} [5]) relies on linear fitting of results at multiple sample sizes. This independence makes KAD especially robust in data-scarce conditions, such as early-stage evaluations of generative models or when high-quality reference datasets are limited.

Overall Computational Efficiency: KAD has a time complexity of $\mathcal{O}(dN^2)$. In practice, this can be significantly faster than $\mathcal{O}(dN^2+d^3)$ for FAD, as the d^3 term can dominate with higher dimensionality. Therefore, KAD is more scalable for higher-dimensional embeddings typical of modern deep audio models.

Parallel Computation The pairwise operations in the computation of KAD (Eq. 4), $-\frac{2}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}k(x_i,y_j)$, can be performed in parallel, enabling substantial acceleration.

3.2 Kernel Function and Bandwidth Selection

KAD relies on evaluating pairwise relationships between embeddings through a kernel function. Many commonly used kernels such as Gaussian, Laplacian, and Matérn kernels are examples of a *characteristic* kernel, meaning that the MMD it induces i) fully distinguishes between two embedding distributions and ii) is zero if and only if the two distributions under test are identical [48]. This property is crucial for model evaluation as it ensures that the KAD metric captures meaningful differences in the embedding distributions of real and generated audio samples. As an example, a cubic polynomial kernel $(x^Ty+1)^2$ cannot differentiate between distributions with the same mean, variance and skewness, but different kurtosis [49].

For the KAD, we choose the Gaussian radial basis function (RBF) kernel:

$$k(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{y}\|^2}{2\sigma^2}\right),$$
 (6)

where σ is the bandwidth parameter. An implicit mapping $\phi(x)$ to the RKHS of the Gaussian RBF kernel is infinite-dimensional and is defined by the property $\langle \varphi(\mathbf{x}), \varphi(\mathbf{y}) \rangle = k(\mathbf{x}, \mathbf{y})$. This kernel has been extensively analyzed and validated in MMD applications [2, 50] for its smoothness, balanced sensitivity to both local and global variations, and well-studied performance across diverse datasets and modalities.

For the value of the bandwidth parameter σ , we follow the commonly adopted median distance heuristic [39], setting σ as the median pairwise distance between the embeddings within the reference set. This heuristic provides a stable baseline with minimal tuning, ensuring the kernel is neither too flat (insensitive to dissimilarities) nor too peaked (over-sensitive to noise). While exploring adaptive or data-driven kernel selection is beyond the scope of this work, our initial experiments indicate that the median heuristic is sufficiently effective. More details are discussed in Appendix A.

4 Experiments

In this section, we present our empirical findings on KAD and comparison with FAD across three key perspectives: (1) Alignment with Human Perception, (2) Convergence with Sample Size and (3) Computation Cost.

4.1 Experiment 1: Reliability of KAD in Perceptual Alignment

A reliable evaluation metric for generative audio should be closely aligned with the human perception of audio quality. While FAD relies on a normality assumption that may not accurately capture the multimodal nature of real-world audio embeddings, KAD takes a distribution-free approach. This flexibility allows it to handle complex acoustic feature representations and potentially align more closely with how humans perceive audio quality.

To validate the perceptual alignment of KAD in comparison with FAD, we use data from the DCASE 2023 Challenge Task 7 submissions [51–68] for Foley sound generation. This dataset provides human rating scores on audio quality for 9 different audio generation models, making it a reliable benchmark for correlating objective metrics with subjective judgments.

We compute both KAD and FAD using embedding from several well-known models, including VGGish [7], PANNs [69], CLAP [47, 70], and PaSST [71], all of which are trained on environmental sounds. These embedding models are widely used for the calculation of FAD scores for text-to-audio [8–20] and vision-to-audio generation[21–36]. Since music-focused models can differ substantially in their learned representations, we also include MERT [72] and CLAP-laion-music [47]

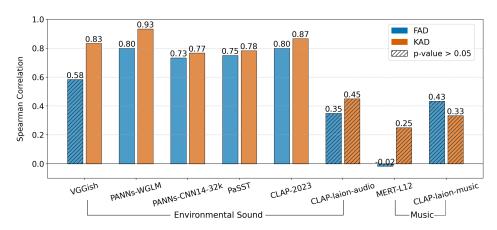


Figure 3: Spearman correlations between metric scores and human perceptual ratings for different embedding models. Since lower scores imply better results for both metrics, correlation values are negative. Correlation values are multiplied by -1 for the convenience of visualization. KAD (orange) consistently achieves higher alignment than FAD (blue).

for completeness. We then measure the Spearman rank correlation between each metric's scores and the average human evaluation scores, as well as the p-value. Correlations with p>0.05 are shaded in 3 to indicate a lack of statistical significance.

As shown in Figure 3, KAD exhibits a Spearman correlation of up to -0.93, notably outperforming FAD whose strongest correlation is -0.80. This suggests that KAD is more effective for differentiating the perceptual nuances captured within the audio data embeddings from a wide range of common audio representations. In contrast, embeddings trained on music data (MERT and CLAP-laion-music) show weaker alignment, consistent with previous findings[4].

Among the tested embedding models, PANNs-WGLM(WaveGram-LogMel) achieves the strongest correlation with human judgments, aligning with prior research that highlighted its suitability for FAD-based evaluations [4]. Based on this observation, we select PANNs-WGLM as the primary embedding model in subsequent experiments to further investigate the performance of KAD.

4.2 Experiment 2: Convergence with Sample Size

To compare how KAD and FAD converge as the evaluation set size N increases, we use the eval split of the Clotho 2 dataset [38] with 1045 samples as the reference set, and samples generated using AudioLDM [9] as the evaluation set. The evaluation samples were generated by conditioning on text captions from the dev split of the Clotho 2 dataset. The number of generated samples starts at N=100 and gradually increases up to N=3839 (the total size of the Clotho 2 dev split). We compute both KAD and FAD under these varying N values to observe their biases and convergence rates.

Figure 4 displays how KAD and FAD evolve as N increases, normalizing each metric by its extrapolated value at $N=\infty$. At small N, FAD shows a distinct positive bias, deviating substantially from its stable value. This deviation decreases roughly by half whenever the sample size doubles, indicating that a large N is needed for FAD to become reliable.

By contrast, KAD remains close to its asymptotic value even at relatively small N, reflecting its unbiased nature. While KAD does exhibit a relatively larger standard deviation (the shaded region) for smaller N, this uncertainty band narrows quickly. Notably, even when accounting for the standard deviation, the range of error for KAD is bounded by the magnitude of bias for FAD, up to the largest sample size tested (N=3839). These results show that KAD can serve as a more stable evaluation metric, especially when the availability of generated audio samples is limited.

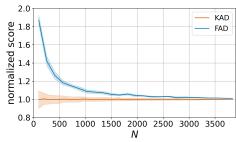


Figure 4: Normalized FAD and KAD scores against increasing embedding sample size. Scores are normalized by their respective extrapolated values at $N=\infty$. The shaded regions indicate standard deviations.

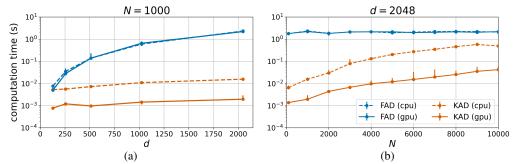


Figure 5: Comparison of FAD and KAD wall-clock computation times. (a) N=1000 with varying d. (b) d=2048 with varying N. Solid lines indicate CPU usage and dotted lines indicate GPU usage. Error bars mark the 5th to 95th percentile of 200 trials.

4.3 Experiment 3: Computation Cost Comparison

To assess the computational efficiency of KAD relative to FAD, we measure their wall-clock times on both CPU and GPU across varying embedding dimensions d and sample sizes N. We use PANNs-WGLM [69], VGGish [7], and CLAP [70] – encompassing dimension sizes from d=128 (VGGish) to d=2048 (PANNs-WGLM). The sample sizes range up to 10k to cover typical open-source audio-text datasets like Clotho [38] and AudioCaps [73].

For the measurements, AMD EPYC 7413 CPU (24 cores) and an Nvidia RTX 3090 GPU were used, and the code is implemented on PyTorch for both KAD and FAD calculations. For FAD, we refactored the Microsoft FAD toolkit [5] for consistency in CPU/GPU usage, thereby ensuring the comparability of runtime measurements. All values were calculated in single-precision floating points.

Figure 5 shows that FAD's computation time increases dramatically with dimension size d, whereas KAD remains relatively stable. This stark difference aligns with the theoretical d^3 scaling of FAD, in contrast to KAD's weaker dependence on d. FAD exhibits significant computational overhead at high dimensions even for small sample sizes. Figure 5a highlights how FAD's wall-clock time (blue lines) escalates with d, while KAD (orange lines) remains nearly flat. At d=2048, the runtime gap can reach three orders of magnitude. Figure 5b further confirms that the main bottleneck for FAD is dimension size, rather than the number of samples. This behavior indicates that FAD is less practical when evaluating embeddings with large d or on resource-limited systems.

Furthermore, KAD benefits considerably from GPU acceleration (dotted vs. solid orange lines), achieving more than an order of magnitude of speedup. Table 2 quantifies these observations, showing consistent performance advantages of KAD over FAD under both CPU and GPU conditions.

\overline{d}	N	CPU		GPU	
		KAD (ours)	FAD	KAD (ours)	FAD
128	100	2.8 ± 0.06	5.7 ± 0.03	0.6 ± 0.03	5.4 ± 0.02
	5000	102.8 ± 1.17	6.7 ± 0.09	4.1 ± 0.06	7.3 ± 0.12
	10000	424.2 ± 4.00	6.9 ± 0.19	12.8 ± 0.10	7.9 ± 0.08
512	100	2.8 ± 0.07	130.2 ± 0.65	1.3 ± 0.01	107.7 ± 0.29
	5000	132.0 ± 1.37	155.5 ± 2.72	5.4 ± 0.12	128.5 ± 1.70
	10000	461.5 ± 3.16	154.6 ± 2.65	17.3 ± 0.20	134.2 ± 1.83
2048	100	6.8 ± 0.12	1776.2 ± 14.5	1.4 ± 0.03	1829.1 ± 21.3
	5000	204.6 ± 1.98	1921.1 ± 14.1	13.0 ± 0.41	2136.5 ± 21.8
	10000	497.9 ± 4.35	2074.9 ± 20.5	46.3 ± 2.41	2174.4 ± 21.6

Table 2: Mean wall-clock times with 95% confidence intervals over 200 trials, in milliseconds. KAD on GPU consistently runs faster than FAD for N=100 and 5000, as well as for N=10000 at higher dimensions.

5 Conclusion

In this paper, we addressed key limitations of the Fréchet Audio Distance (FAD) for evaluating generative audio models and proposed the Kernel Audio Distance (KAD) as a more robust alternative. Built on the Maximum Mean Discrepancy (MMD), KAD avoids making statistical assumptions about the embedding distributions, provides unbiased results for all sample sizes, and offers a computational complexity that scale more efficiently, particularly at higher dimensionalities.

We define KAD as the MMD between reference and evaluation audio embedding sets using a Gaussian RBF kernel with the median-distance bandwidth heuristic. To validate its effectiveness, we compare both KAD and FAD against human evaluation data, observe their convergence behaviors with increasing sample sizes, and measure their CPU and GPU runtimes across a range of dimensionalities and sample sizes.

Our findings show that KAD aligns more strongly with human judgments than FAD across various common audio embedding models, with especially high correlation with PANNs-WGLM. Moreover, its score remains consistent regardless of sample size, making it practical for resource-constrained or early-stage model evaluations, and its computational overhead is up to orders of magnitude lower for higher dimensional (~2024) embeddings compared to FAD due to the reduction of the complexity from $\mathcal{O}(dN^2+d^3)$ to $\mathcal{O}(dN^2)$ and its amenability to parallel computation.

These advantages position KAD as an efficient, comprehensive, and scalable tool for benchmarking generative audio models. By more accurately capturing human-perceived audio quality, KAD can support the development of more reliable evaluation practices in the field. The accompanying open-source toolkit is provided to encourage widespread adoption, experimentation, and ongoing improvements to the development and assessment of generative audio models.

References

- [1] Kevin Kilgour, Mauricio Zuluaga, Dominik Roblek, and Matthew Sharifi. Fréchet audio distance: A metric for evaluating music enhancement algorithms. *arXiv preprint arXiv:1812.08466*, 2018.
- [2] Sadeep Jayasumana, Srikumar Ramalingam, Andreas Veit, Daniel Glasner, Ayan Chakrabarti, and Sanjiv Kumar. Rethinking fid: Towards a better evaluation metric for image generation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 9307–9315, 2024.
- [3] Min Jin Chong and David Forsyth. Effectively unbiased fid and inception score and where to find them. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6070–6079, 2020.
- [4] Modan Tailleur, Junwon Lee, Mathieu Lagrange, Keunwoo Choi, Laurie M Heller, Keisuke Imoto, and Yuki Okamoto. Correlation of fréchet audio distance with human perception of environmental audio is embedding dependant. *arXiv preprint arXiv:2403.17508*, 2024.

- [5] Azalea Gui, Hannes Gamper, Sebastian Braun, and Dimitra Emmanouilidou. Adapting frechet audio distance for generative music evaluation. In ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 1331–1335. IEEE, 2024.
- [6] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems, 30, 2017.
- [7] Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis, Jort F Gemmeke, Aren Jansen, R Channing Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, et al. Cnn architectures for large-scale audio classification. In 2017 ieee international conference on acoustics, speech and signal processing (icassp), pages 131–135. IEEE, 2017.
- [8] Chris Donahue, Julian McAuley, and Miller Puckette. Adversarial audio synthesis. arXiv preprint arXiv:1802.04208, 2018.
- [9] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. Audioldm: text-to-audio generation with latent diffusion models. In *Proceedings of the 40th International Conference on Machine Learning*, pages 21450–21474, 2023.
- [10] Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin Liu, Xiang Yin, and Zhou Zhao. Make-an-audio: Text-to-audio generation with prompt-enhanced diffusion models. In *International Conference on Machine Learning*, pages 13916–13932. PMLR, 2023.
- [11] Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. Text-to-audio generation using instruction guided latent diffusion model. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 3590–3598, 2023.
- [12] Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. arXiv preprint arXiv:2209.15352, 2022.
- [13] Yoonjin Chung, Junwon Lee, and Juhan Nam. T-foley: A controllable waveform-domain diffusion model for temporal-event-guided foley sound synthesis. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6820–6824. IEEE, 2024.
- [14] Haohe Liu, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Qiao Tian, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D Plumbley. Audioldm 2: Learning holistic audio generation with self-supervised pretraining. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 2024.
- [15] Yatong Bai, Trung Dang, Dung Tran, Kazuhito Koishida, and Somayeh Sojoudi. Accelerating diffusion-based text-to-audio generation with consistency distillation. *arXiv preprint arXiv:2309.10740*, 2023.
- [16] Navonil Majumder, Chia-Yu Hung, Deepanway Ghosal, Wei-Ning Hsu, Rada Mihalcea, and Soujanya Poria. Tango 2: Aligning diffusion-based text-to-audio generative models through direct preference optimization. In ACM Multimedia 2024, 2024.
- [17] Zhifeng Kong, Sang-gil Lee, Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, Rafael Valle, Soujanya Poria, and Bryan Catanzaro. Improving text-to-audio models with synthetic captions. arXiv preprint arXiv:2406.15487, 2024.
- [18] Jinlong Xue, Yayue Deng, Yingming Gao, and Ya Li. Auffusion: Leveraging the power of diffusion and large language models for text-to-audio generation. arXiv preprint arXiv:2401.01044, 2024.
- [19] Zach Evans, Julian D Parker, CJ Carr, Zack Zukowski, Josiah Taylor, and Jordi Pons. Stable audio open. arXiv preprint arXiv:2407.14358, 2024.
- [20] Jiarui Hai, Yong Xu, Hao Zhang, Chenxing Li, Helin Wang, Mounya Elhilali, and Dong Yu. Ezaudio: Enhancing text-to-audio generation with efficient diffusion transformer. arXiv preprint arXiv:2409.10819, 2024.
- [21] Marco Comunità, Riccardo F Gramaccioni, Emilian Postolache, Emanuele Rodolà, Danilo Comminiello, and Joshua D Reiss. Syncfusion: Multimodal onset-synchronized video-to-audio foley synthesis. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 936–940. IEEE, 2024.
- [22] Junwon Lee, Jaekwon Im, Dabin Kim, and Juhan Nam. Video-foley: Two-stage video-to-sound generation via temporal event condition for foley sound. arXiv preprint arXiv:2408.11915, 2024.

- [23] Yujin Jeong, Yunji Kim, Sanghyuk Chun, and Jiyoung Lee. Read, watch and scream! sound generation from text and video. arXiv preprint arXiv:2407.05551, 2024.
- [24] Yongqi Wang, Wenxiang Guo, Rongjie Huang, Jiawei Huang, Zehan Wang, Fuming You, Ruiqi Li, and Zhou Zhao. Frieren: Efficient video-to-audio generation with rectified flow matching. arXiv preprint arXiv:2406.00320, 2024.
- [25] Santiago Pascual, Chunghsin Yeh, Ioannis Tsiamas, and Joan Serrà. Masked generative video-to-audio transformers with enhanced synchronicity. In *European Conference on Computer Vision*, pages 247–264. Springer, 2025.
- [26] Ilpo Viertola, Vladimir Iashin, and Esa Rahtu. Temporally aligned audio for video with autoregression. arXiv preprint arXiv:2409.13689, 2024.
- [27] Xiulong Liu, Kun Su, and Eli Shlizerman. Tell what you hear from what you see video to audio generation through text. In The Thirty-eighth Annual Conference on Neural Information Processing Systems, 2024.
- [28] Ziyang Chen, Prem Seetharaman, Bryan Russell, Oriol Nieto, David Bourgin, Andrew Owens, and Justin Salamon. Video-guided foley sound generation with multimodal controls. arXiv preprint arXiv:2411.17698, 2024.
- [29] Wei Guo, Heng Wang, Weidong Cai, and Jianbo Ma. Gotta hear them all: Sound source aware vision to audio generation. *arXiv preprint arXiv:2411.15447*, 2024.
- [30] Saksham Singh Kushwaha and Yapeng Tian. Vintage: Joint video and text conditioning for holistic audio generation. *arXiv preprint arXiv:2412.10768*, 2024.
- [31] Ho Kei Cheng, Masato Ishii, Akio Hayakawa, Takashi Shibuya, Alexander Schwing, and Yuki Mitsufuji. Taming multimodal joint training for high-quality video-to-audio synthesis. arXiv preprint arXiv:2412.15322, 2024.
- [32] Hao-Wen Dong, Xiaoyu Liu, Jordi Pons, Gautam Bhattacharya, Santiago Pascual, Joan Serrà, Taylor Berg-Kirkpatrick, and Julian McAuley. Clipsonic: Text-to-audio synthesis with unlabeled videos and pretrained language-vision models. In 2023 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA), pages 1–5. IEEE, 2023.
- [33] Heng Wang, Jianbo Ma, Santiago Pascual, Richard Cartwright, and Weidong Cai. V2a-mapper: A lightweight solution for vision-to-audio generation by connecting foundation models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 15492–15501, 2024.
- [34] Roy Sheffer and Yossi Adi. I hear your true colors: Image guided audio generation. In ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 1–5. IEEE, 2023
- [35] Xinhao Mei, Varun Nagaraja, Gael Le Lan, Zhaoheng Ni, Ernie Chang, Yangyang Shi, and Vikas Chandra. Foleygen: Visually-guided audio generation. In 2024 IEEE 34th International Workshop on Machine Learning for Signal Processing (MLSP), pages 1–6. IEEE, 2024.
- [36] Xihua Wang, Yuyue Wang, Yihan Wu, Ruihua Song, Xu Tan, Zehua Chen, Hongteng Xu, and Guodong Sui. Tiva: Time-aligned video-to-audio generation. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 573–582, 2024.
- [37] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and* pattern recognition, pages 2818–2826, 2016.
- [38] Konstantinos Drossos, Samuel Lipping, and Tuomas Virtanen. Clotho: An audio captioning dataset. In ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 736–740. IEEE, 2020.
- [39] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. The Journal of Machine Learning Research, 13(1):723–773, 2012.
- [40] Mikołaj Bińkowski, Danica J Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. arXiv preprint arXiv:1801.01401, 2018.
- [41] Javier Nistal, Stefan Lattner, and Gaël Richard. Comparing representations for audio synthesis using generative adversarial networks. In 2020 28th European Signal Processing Conference (EUSIPCO), pages 161–165. IEEE, 2021.

- [42] Javier Nistal, Marco Pasini, Cyran Aouameur, Maarten Grachten, and Stefan Lattner. Diff-a-riff: Musical accompaniment co-creation via latent diffusion models. *arXiv* preprint arXiv:2406.08384, 2024.
- [43] Maarten Grachten. Measuring audio prompt adherence with distribution-based embedding distances. arXiv preprint arXiv:2404.00775, 2024.
- [44] Jiatong Shi, Hye-jin Shim, Jinchuan Tian, Siddhant Arora, Haibin Wu, Darius Petermann, Jia Qi Yip, You Zhang, Yuxun Tang, Wangyou Zhang, et al. Versa: A versatile evaluation toolkit for speech, audio, and music. arXiv preprint arXiv:2412.17667, 2024.
- [45] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [46] Zachary Novack, Ge Zhu, Jonah Casebeer, Julian McAuley, Taylor Berg-Kirkpatrick, and Nicholas J Bryan. Presto! distilling steps and layers for accelerating music generation. arXiv preprint arXiv:2410.05167, 2024.
- [47] Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation. In 2023 IEEE International Conference on Acoustics, Speech and Signal Processing, pages 1–5. IEEE, 2023.
- [48] Bharath K Sriperumbudur, Kenji Fukumizu, and Gert RG Lanckriet. Universality, characteristic kernels and rkhs embedding of measures. *Journal of Machine Learning Research*, 12(7), 2011.
- [49] Bharath K Sriperumbudur, Arthur Gretton, Kenji Fukumizu, Bernhard Schölkopf, and Gert RG Lanckriet. Hilbert space embeddings and metrics on probability measures. *The Journal of Machine Learning Research*, 11:1517–1561, 2010.
- [50] Raif M Rustamov. Closed-form expressions for maximum mean discrepancy with applications to wasserstein auto-encoders. arXiv preprint arXiv:1901.03227, 2019.
- [51] Keunwoo Choi, Jaekwon Im, Laurie Heller, Brian Mcfee, Keisuke Imoto, Yuki Okamoto, Mathieu Lagrange, and Shinnosuke Takamichi. Foley sound synthesis at the dcase 2023 challenge. In 2023 Workshop on Detection and Classification of Acoustic Scenes and Events (DCASE 2023), 2023.
- [52] Siwei Huang, Jisheng Bai, Yafei Jia, and Jianfeng Chen. Jless submission to dcase2023 task7: Foley sound synthesis using non-autoagressive generative model. Technical report, School of Marine Science and Technology, Northwestern Polytechnical University, Xi'an, China, LianFeng Acoustic Technologies Co., Ltd. Xi'an, China, 2023.
- [53] Won-Gook Choi and Joon-Hyuk Chang. Hyu submission for the dcase 2023 task 7: Diffusion probabilistic model with adversarial training for foley sound synthesis. Technical report, Department of Electronic Engineering, Hanyang University, Seoul, Republic of Korea, Department of Electronic Engineering, Hanyang University, Seoul, Republic of Korea, 2023.
- [54] Minsung Kang, Sangshin Oh, Hyeongi Moon, Kyungyun Lee, and Ben Sangbae Chon. Fall-e: Gaudio foley synthesis system. Technical report, Gaudio Lab, Inc., Seoul, South Kore, 2023.
- [55] Chae-Woon Bang, Nam Kyun Kim, and Chanjun Chun. High-quality foley sound synthesis using monte carlo dropout. Technical report, Chosun University, Gwangju, South Korea, Korea Automotive Technology Institute, Gwanjgu, South Korea, 2023.
- [56] Yoonjin Chung, Junwon Lee, and Juhan Nam. Foley sound synthesis in waveform domain with diffusion model. Technical report, Graduate School of AI, KAIST, Graduate School of Culture Technology, KAIST, 2023.
- [57] Shitong Fan, Qiaoxi Zhu, Feiyang Xiao, Haiyan Lan, Wenwu Wang, and Jian Guan1. Foley sound synthesis with audioldm for dcase2023 task 7. Technical report, Group of Intelligent Signal Processing (GISP), College of Computer Science and Technology, Harbin Engineering University, Harbin, China, Centre for Audio, Acoustic and Vibration (CAAV), University of Technology Sydney, Ultimo, Australia, Centre for Vision, Speech and Signal Processing (CVSSP), University of Surrey, Guildford, UK, 2023.
- [58] Hae Chun Chung, Yuna Lee, and Jae Hoon Jung. Foley sound synthesis based on gan using contrastive learning without label information. Technical report, KT Corporation, Republic of Korea, 2023.

- [59] Purnima Kamath, Tasnim Nishat Islam, Chitralekha Gupta, Lonce Wyse, and Suranga Nanayakkara. Dcase task-7: Stylegan2-based foley sound synthesis. Technical report, National University of Singapore, Singapore and Bangladesh University of Engineering and Technology, Bangladesh and Universitat Pompeu Fabra, Barcelona, Spain, 2023.
- [60] Kyungsu Kim, Jinwoo Lee, Hayoon Kim, and Kyogu Lee. Conditional foley sound synthesis with limited data: Two-stage data augmentation approach with stylegan2-ada. Technical report, Seoul National University Department of Intelligence and Information, 2023.
- [61] Junhyeok Lee1, Hyeonuk Nam, and Yong-Hwa Park. Vifs: An end-to-end variational inference for foley sound synthesis. Technical report, maum.ai Inc., Republic of Korea and Korea Advanced Institute of Science and Technology, Republic of Korea, 2023.
- [62] Ashwin Pillay, Sage Betko, Ari Liloia, Hao Chen, and Ankit Shah. Dcase task 7: Foley sound synthesis. Technical report, Carnegie Mellon University, Pittsburgh, USA, 2023.
- [63] Timo Wendner, Patricia Hu, Tara Jadidi, and Alexander Neuhauser. Audio diffusion for foley sound synthesis. Technical report, Johannes Kepler University, Linz, Austria, 2023.
- [64] Zeyu Xie, Xuenan Xu, Baihan Li, Mengyue Wu, and Kai Yu. The x-lance system for dcase2023 challenge task 7: Foley sound synthesis track b. Technical report, MoE Key Lab of Artificial Intelligence X-LANCE Lab, Department of Computer Science and Engineering AI Institute, Shanghai Jiao Tong University, Shanghai, China, 2023.
- [65] Yi Yuan, Haohe Liu, Xubo Liu, Xiyuan Kang, Mark D. Plumbley, and Wenwu Wang. Latent diffusion model based foley sound generation system for dcase challenge 2023 task 7. Technical report, University of Surrey, Guildford, United Kingdom, 2023.
- [66] Anbin Qi. Auto-bit for dcase2023 task7 technical reports: Assemble system of bitdiffusion and pixelsnail. Technical report, School Information and Electronics, Beijing Institute of Technology, Beijing, China, 2023.
- [67] Haojie Zhang, Kun Qian, Lin Shen, Lujundong Li, Kele Xu, and Bin Hu. From noise to sound: Audio synthesis via diffusion models. Technical report, Key Laboratory of Brain Health Intelligent Evaluation and Intervention, Ministry of Education (Beijing Institute of Technology), P. R. China, School of Medical Technology, Beijing Institute of Technology, P. R. China, National University of Defense Technology, P. R. China, 2023.
- [68] Robin Scheibler, Takuya Hasumi, Yusuke Fujita, Tatsuya Komatsu, Ryuichi Yamamoto, and Kentaro Tachibana. Class-conditioned latent diffusion model for dcase 2023 foley sound synthesis challenge. Technical report, LINE Corporation, Tokyo, Japan, 2023.
- [69] Qiuqiang Kong, Yin Cao, Talha Iqbal, Yuxuan Wang, Zihao Yin, Wenwu Wang, and Mark D. Plumbley. Panns: Large-scale pretrained audio neural networks for audio pattern recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:2880–2894, 2020.
- [70] Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. Clap learning audio concepts from natural language supervision. In 2023 IEEE International Conference on Acoustics, Speech and Signal Processing, pages 1–5. IEEE, 2023.
- [71] Khaled Koutini, Jan Schlüter, Hamid Eghbal-zadeh, and Gerhard Widmer. Efficient training of audio transformers with patchout. In *Interspeech 2022, 23rd Annual Conference of the International Speech Communication Association, Incheon, Korea, 18-22 September 2022*, pages 2753–2757. ISCA, 2022.
- [72] Yizhi Li, Ruibin Yuan, Ge Zhang, Yinghao Ma, Xingran Chen, Hanzhi Yin, Chenghao Xiao, Chenghua Lin, Anton Ragni, Emmanouil Benetos, et al. Mert: Acoustic music understanding model with large-scale self-supervised training. *arXiv preprint arXiv:2306.00107*, 2023.
- [73] Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating captions for audios in the wild. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 119–132, 2019.
- [74] Iver Jordal. Audiomentations. https://github.com/iver56/audiomentations?tab=readme-ov-file, August 2021. Accessed: 2025-02-03.
- [75] Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. High fidelity neural audio compression. Transactions on Machine Learning Research, 2023.

- [76] Aurora Linh Cramer, Ho-Hsiang Wu, Justin Salamon, and Juan Pablo Bello. Look, listen, and learn more: Design choices for deep audio embeddings. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3852–3856. IEEE, 2019.
- [77] Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar. High-fidelity audio compression with improved ryagan. Advances in Neural Information Processing Systems, 36, 2024.
- [78] Pranay Manocha, Zeyu Jin, Richard Zhang, and Adam Finkelstein. Cdpam: Contrastive learning for perceptual audio similarity. In ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 196–200. IEEE, 2021.
- [79] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. Advances in neural information processing systems, 33:12449–12460, 2020.
- [80] Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. IEEE/ACM transactions on audio, speech, and language processing, 29:3451–3460, 2021.
- [81] Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, et al. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. IEEE Journal of Selected Topics in Signal Processing, 16(6):1505–1518, 2022.
- [82] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR, 2023.
- [83] Junwon Lee, Modan Tailleur, Laurie M Heller, Keunwoo Choi, Mathieu Lagrange, Brian McFee, Keisuke Imoto, and Yuki Okamoto. Challenge on sound scene synthesis: Evaluating text-to-audio generation. In Audio Imagination: NeurIPS 2024 Workshop AI-Driven Speech, Music, and Sound Generation, 2024.
- [84] Mathieu Lagrange, Junwon Lee, Modan Tailleur, Laurie M. Heller, Keunwoo Choi, Brian McFee, Keisuke Imoto, and Yuki Okamoto. Sound scene synthesis at the dcase 2024 challenge. arXiv preprint arXiv:2501.08587, 2025.

A Median Pair-wise Distance Heuristic for the Kernel Bandwidth

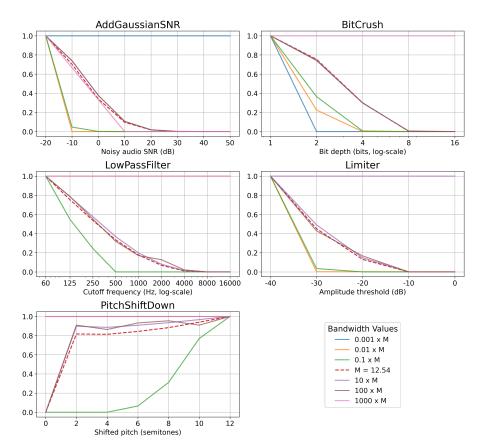


Figure 6: Effect of audio degradations to the MMD values, where the maximum value is normalized to 1 for each bandwidth result. MMD values that are smaller than $\varepsilon=10^{-12}$ are clipped to ε before normalization.

To assess the effectiveness of the median pairwise distance heuristic for setting the bandwidth parameter σ , we compared the MMD values calculated at different bandwidth settings as the quality of the evaluation data is artificially degraded. If the heuristic to works as intended, the MMD at the median distance bandwidth should be well-correlated with the audio quality for the perceptually relevant ranges of degradation.

We compared the MMD scores between the clean audio files from the Clotho dataset's *eval* split, and the same files under various audio signal transformations including Gaussian noise injection, bit depth reduction, low-pass filtering, dynamic range reduction (limiter), and pitch shifts. The degradations were applied using the *audiomentations* python library [74].

For each transformation, MMD scores were calculated using the median pairwise distance as the bandwidth, as well as bandwidths scaled by factors from $0.001 \times$ up to $1000 \times$ relative to the median. The scores were normalized such that the maximum MMD score for each bandwidth result is 1, allowing us to observe trends in the scores independently of their absolute values.

The results are shown in Figure 6. A reliable metric is expected to show a monotonic increase in score as the severity of the degradation increases, because higher KAD scores should correlate with worse audio quality. When the median pairwise distance is used as the bandwidth, the MMD scores consistently respect the expected monotonic trend across all degradations tested. This suggests that the median heuristic provides stable and meaningful results. Bandwidths scaled by $10\times$ and $100\times$ the median also produced reliable results, maintaining the monotonic relationship. However, smaller bandwidths $(0.001\times$ to $0.1\times$ resulted in scores that drop too quickly, diminishing the discriminative

power of the metric. Larger bandwidths ($1000 \times$ the median) tend to flatten the scores, reducing sensitivity to dissimilarities.

These initial experiments suggest that the median bandwidth heuristic is a robust and effective choice for calculating MMD scores. While certain scaled bandwidths may also be viable, the median heuristic avoids nonsensical results, ensuring stable and interpretable evaluations without additional tuning. Therefore, we recommend its use for kernel bandwidth selection in this context.

B Bias of FAD

In order to analyze the bias of FAD, we revisit Equation 1:

$$FAD^{2}(X,Y) = \|\mu_{X} - \mu_{Y}\|_{2}^{2} + tr\left(\Sigma_{X} + \Sigma_{Y} - 2\sqrt{\Sigma_{X}\Sigma_{Y}}\right).$$

For finite samples, μ_X , μ_Y and Σ_X , Σ_Y are replaced by their sample estimates $\hat{\mu}_X$, $\hat{\mu}_Y$ and $\hat{\Sigma}_X$, $\hat{\Sigma}_Y$. This introduces bias due to finite sample effects.

The first term, $\|\mu_X - \mu_Y\|^2$, is unbiased since the sample means $\hat{\mu}_X$ and $\hat{\mu}_Y$ are unbiased estimators of the true means μ_X and μ_Y . Thus,

$$\mathbb{E}[\|\hat{\mu}_X - \hat{\mu}_Y\|^2] = \|\mu_X - \mu_Y\|^2.$$

The second term, tr $(\Sigma_X + \Sigma_Y - 2\sqrt{\Sigma_X \Sigma_Y})$, is affected by finite sample sizes. The sample covariance matrix $\hat{\Sigma}$ is a biased estimator of the true covariance Σ :

$$\mathbb{E}[\hat{\Sigma}] = \frac{N-1}{N} \cdot \Sigma,$$

where N is the sample size. Consequently, the expected value of the trace of the covariance matrices is:

$$\mathbb{E}[\operatorname{tr}(\hat{\Sigma}_X + \hat{\Sigma}_Y)] = \frac{N-1}{N} \cdot \operatorname{tr}(\Sigma_X + \Sigma_Y).$$

In the second term, $2(\Sigma_X \Sigma_Y)^{1/2}$ is nonlinear in the covariances, and its exact bias is difficult to compute. However, to first-order approximation, the bias of FAD due to finite samples can be expressed as:

$$\mathrm{Bias}_{\mathrm{FAD}} \approx \tfrac{1}{N} \cdot \left(\mathrm{tr}(\Sigma_X) + \mathrm{tr}(\Sigma_Y) \right) + \mathcal{O}\left(\tfrac{1}{N^2} \right).$$

Here, the primary contribution to the bias arises from the underestimation of the covariance matrices, which scales inversely with the sample size N.

C kadtk: KAD Toolkit Release

We release a toolkit named kadtk that can calculate KAD scores from input audios. Given input audio directories for the ground-truth reference set and target evaluation set, the toolkit calculates the score and saves it to the given output file path in CSV format. The toolkit is written in Python and supports Pytorch and Tensorflow environments (refer to our Readme document for further details). It supports numerous models for embedding extraction: CLAP [47, 70], Encodec [75], MERT [72], VGGish [7], PANNs [69], OpenL3 [76], PaSST [71], DAC [77], CDPAM [78], Wav2vec2.0 [79], HuBERT [80], WavLM [81], and Whisper [82]. This covers a wide range of the audio domain including general environment sounds (sound effects, foley sounds, etc.), music, and speech. We also support FAD calculation for comparison.

The kadtk is released under the MIT License, allowing unrestricted use, modification, and distribution with proper attribution. The full license text is included in the repository. Some codes were brought from the FAD toolkit (fadtk) [5, 83, 84]. We sincerely thank the authors for sharing the code as an open source. Note that fadtk was also licensed under the MIT License.

For the computational cost comparison in Section 4.3, we used torch.linalg.eigval and torch.sqrt to compute the covariance matrix of FAD, ensuring a fair comparison with KAD

in a fully parallelized GPU setting. However, this approach is prone to accuracy issues, often leading to discrepancies with fadtk, which relies on scipy.linalg.sqrtm (available only for CPU computation). To ensure high FAD score accuracy at the cost of increased runtime, we used scipy.linalg.sqrtm for the perceptual alignment experiment (Section 4.1) and in the released toolkit.