# MAPSS: Manifold-based Assessment of Perceptual Source Separation

**Anonymous authors**Paper under double-blind review

000

001

002003004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

025

026

027

028

029

031

032

034

037

040

041

042

043

044

046

047

048

051

052

#### **ABSTRACT**

Objective assessment of audio source-separation systems still mismatches subjective human perception, especially when interference from competing talkers and distortion of the target signal interact. We introduce Perceptual Separation (PS) and Perceptual Match (PM), a complementary pair of measures that, by design, isolate these leakage and distortion factors. Our intrusive approach generates a set of fundamental distortions, e.g., clipping, notch filter, and pitch shift from each reference waveform signal in the mixture. Distortions, references, and system outputs from all sources are independently encoded by a pre-trained self-supervised model, then aggregated and embedded with a manifold learning technique called diffusion maps, which aligns Euclidean distances on the manifold with dissimilarities of the encoded waveform representations. On this manifold, PM captures the self-distortion of a source by measuring distances from its output to its reference and associated distortions, while PS captures leakage by also accounting for distances from the output to non-attributed references and distortions. Both measures are differentiable and operate at a resolution as high as 75 frames per second, allowing granular optimization and analysis. We further derive, for both measures, frame-level deterministic error radius and non-asymptotic, high-probability confidence intervals. Experiments on English, Spanish, and music mixtures show that, against 14 widely used measures, the PS and PM are almost always placed first or second in linear and rank correlations with subjective human mean-opinion scores.

#### 1 Introduction

Reliable perceptual evaluation is critical for source-separation progress, yet gold-standard listening tests are costly and slow (ITU-T., 1996; 2003; 2018). Thus, research relies on objective metrics that blur two distinct failures, interference from competing talkers and target distortion. Disentangling these modes can better align with listener perception and accelerate trustworthy development.

Existing measures such as the signal-to-distortion ratio (SDR), signal-to-interference ratio (SIR), signal-to-artifacts ratio (SAR) (Vincent et al., 2006), scale-invariant SDR (SI-SDR) (Le Roux et al., 2019) and alike usually compute ratios between source to various disturbances in the waveform domain, offering low complexity and widespread adoption. However, even jointly, they are defined with an intrinsic ambiguity to whether an error stems for leakage or self-distortion. Classical intrusive perceptual and intelligibility metrics like the PESQ (Rix et al., 2001), STOI (Taal et al., 2011) and ESTOI (Jensen & Taal, 2016) map an entire utterance to a scaled mean-opinion score (MOS) using hand-crafted auditory features. Designed preliminary for speech enhancement, they perform well for corrupted noisy-reverberant speech utterances but may not account for leakage, while also lacking to provide access to their inherent granular processing. Learned black-box metrics such as the DNSMOS family (Reddy et al., 2022) that are trained end-to-end to predict crowd-sourced MOS, as well as SpeechBERTscore (Saeki et al., 2024) and Sheet-SSQA (Huang et al., 2025), have shown promising results on various speech tasks, but do not offer confidence in their decisions. Spectral-distance metrics are interpretable but tend to mask where degradations occur, e.g., the popular Mel-Cepstral Distortion (MCD) (Fukada et al., 1992) collapses the spectral envelope into a global value. Even when taking into account a broader set of metrics, as available in recently developed speech quality assessment toolkits (Shi et al., 2024), no existing family of measures can simultaneously disentangle leakage from distortion, offer granular analysis, and provide error estimates for their decisions.

We introduce the Perceptual Separation (PS) and Perceptual Match (PM), the first measures for source separation that functionally disentangle leakage and self-distortion. Inspired by auditory theory (Gabrielsson & Sjögren, 1979; Jekosch, 2004; Wilson & Fazenda, 2014; Bannister et al., 2024), we apply a set of fundamental distortions to every reference waveform, intended to create a wide cover of perceptual auditory field around the reference. These distortions range from mildly-intrusive short-tailed reverberations to highly degrading hard clipping. A pretrained self-supervised model, e.g., wav2vec 2.0 (Baevski et al., 2020), is used to independently encode the waveforms of references, distortions, and system outputs across all sources, in a resolution as high as 75 framesper-second. These representations are aggregated and projected via a manifold learning technique called diffusion maps (Coifman & Lafon, 2006) onto a low-dimensional manifold. A key property of diffusion maps aligns Euclidean distances between points on the manifold with dissimilarities between their encoded representations. On the manifold, PM quantifies self-distortion by measuring how far an output lies from its attributed reference and the distortions, whereas PS quantifies leakage by comparing these distances with the output proximity to non-attributed references and distortions.

Evaluations on the SEBASS database (Kastner & Herre, 2022) with mixtures of English, Spanish, and music, show that compared to 14 widely used measures, PS and PM almost always achieve first-or second-place rankings in both linear and rank correlations with human scores, with the exception of Spanish rank correlations, where they remain within the top third. We derive granular theoretical deterministic error radius and high-probability confidence intervals (CIs) for both measures, enabling frame-level guarantees on the reliability of the measures. In almost all scenarios, the worst-case error radius would not lower the PS and PM rankings. In addition, the normalized mutual information (NMI) (Danon et al., 2005) between the PS and PM values shows that they are highly complementary.

\*Two lines about appendices and disclaimers on LLM usage\*

#### 2 PROBLEM FORMULATION

Notational remark. Column vectors and matrices are written in bold and other symbols in non-bold.

Consider a source separation system performing inference on an audio mixture (Vincent et al., 2018). In a time frame f that consists of L samples, let  $N_f \geq 2$  denote the number of active sources and  $S_f$  their index set. The observed mixture  $\mathbf{z}_f \in \mathbb{R}^L$  is modeled as:

$$\mathbf{z}_f = \sum_{i \in \mathcal{S}_f} \mathbf{y}_{i,f} + \mathbf{v}_f,\tag{1}$$

For  $i \in \mathcal{S}_f$ , we denote  $\mathbf{y}_{i,f} \in \mathbb{R}^L$  the reference signal of the *i*-th source in frame f, potentially including interference inherent to its original conditions.  $\mathbf{v}_f$  represents system and environmental interference, assumed statistically independent of the sources. The estimation of  $\mathbf{y}_{i,f}$  is denoted  $\hat{\mathbf{y}}_{i,f}$ .

Given source indices  $i, j \in S_f$  in time frame f, our goal is to introduce these two measures:

- The perceptual separation (PS) measure quantifies how well  $\hat{\mathbf{y}}_{i,f}$  is perceptually separated from all interfering sources  $\{\mathbf{y}_{j,f}\}_{j\neq i}$ .
- The perceptual match (PM) measure quantifies how closely the estimated source  $\hat{\mathbf{y}}_{i,f}$  perceptually aligns with its reference  $\mathbf{y}_{i,f}$ .

# 3 DIFFUSION MAPS: THEORETICAL FOUNDATIONS

**Notational remark.** Sections are denoted by  $\S$ . Symbols are carried over from  $\S 2$ , except for indices i, j that are repurposed, and the subscript f that is dropped since we concern a fixed time frame.

Diffusion maps is a manifold learning method that represents high-dimensional data in a low-dimensional space by capturing geometric and structural relationships (Coifman & Lafon, 2006). Consider the set  $\mathcal{X} = \{\mathbf{x}_i\}_{i=1}^N$  with  $\mathbf{x}_i \in \mathbb{R}^M$  for all i, e.g., feature vectors from wav2vec 2.0 (Baevski et al., 2020). An affinity matrix  $\mathbf{K} \in \mathbb{R}^{N \times N}$  is calculated between the high-dimensional vectors:

$$\mathbf{K}_{i,j} = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma_{\mathbf{K}}^2}\right),\tag{2}$$

where  $i, j \in \{1, ..., N\}$  and  $\forall i, j : 0 \le K_{i,j} \le 1$ , and  $\sigma_{\mathbf{K}}^2 = \text{median}\left\{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2 \mid i \ne j\right\}$ . To account for non-uniform sampling density of points, an  $\alpha$ -normalization replaces  $\mathbf{K}$  by  $\mathbf{K}^{(\alpha)}$ :

$$\mathbf{K}_{i,j}^{(\alpha)} = \frac{\mathbf{K}_{i,j}}{\left(v_i v_j\right)^{\alpha}},\tag{3}$$

where  $\alpha \in [0,1]$  and  $v_i = \sum_{j=1}^N \mathbf{K}_{i,j}$ . Then, we define the diagonal degree-matrix  $\mathbf{D}^{(\alpha)}$ , given by  $\mathbf{D}^{(\alpha)} = \operatorname{diag}\left(v_0^{(\alpha)},\ldots,v_{N-1}^{(\alpha)}\right) \in \mathbb{R}^{N\times N}$ , where  $v_i^{(\alpha)} = \sum_{j=1}^N \mathbf{K}_{i,j}^{(\alpha)}$  and  $\forall i: v_i^{(\alpha)} > 0$  by construction. We assume  $\alpha$  is fixed and for readability we neglect the  $\alpha$  notation from now on.

The probability transition operator P on K is defined with (3) as:

$$\mathbf{P} = \mathbf{D}^{-1} \mathbf{K} \in \mathbb{R}^{N \times N}. \tag{4}$$

Note  $\mathbf{P}$  is row-stochastic, so  $\forall i, j: \mathbf{P}_{ij} \geq 0, \ \sum_{j=1}^{N} \mathbf{P}_{ij} = 1$ . Spectral decomposition on  $\mathbf{P}$  reveals a trivial right eigenvector  $\mathbf{u}_0 = \mathbf{1} \in \mathbb{R}^N$  with eigenvalue  $\lambda_0 = 1$ . Remaining eigenvectors  $\{\mathbf{u}_\ell\}_{\ell=1}^{N-1}$  are associated with eigenvalues  $\{\lambda_\ell\}_{\ell=1}^{N-1}$  and ordered as  $1 > \lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_{N-1} > 0$ , so that:

$$\mathbf{P}\mathbf{u}_{\ell} = \lambda_{\ell}\mathbf{u}_{\ell}.\tag{5}$$

Denoting  $\mathbf{u}_j(i)$  the *i*-th element of the *j*-th eigenvector, then the embedding of  $\mathbf{x}_i$  onto manifold  $\mathcal{M}$  can be expressed with the eigenfunctions in (5), by the embedding operation  $\Psi_t : \mathbb{R}^M \to \mathbb{R}^{N-1}$ :

$$\mathbf{\Psi}_t(\mathbf{x}_i) = \left(\lambda_1^t \mathbf{u}_1(i), \, \lambda_2^t \mathbf{u}_2(i), \, \dots, \, \lambda_{N-1}^t \mathbf{u}_{N-1}(i)\right)^T. \tag{6}$$

where t is the number of Markov chain steps, controlling the diffusion scale of the embedding. The eigenvalues in  $\{\mathbf{u}_\ell\}_{\ell=1}^{N-1}$  are orthonormal under the stationary measure  $\boldsymbol{\pi} = [\pi_1, \pi_2, \dots, \pi_N]^T$ :

$$\pi_i = \frac{\mathbf{D}_{ii}}{\sum_{j=1}^{N} \mathbf{D}_{jj}}, \quad \pi_i \in (0, 1).$$
 (7)

Let  $D_t(i, j)$  be the diffusion distance at time step t between two points  $\mathbf{x}_i$  and  $\mathbf{x}_j$ :

$$D_t^2(i,j) = \sum_{m=1}^N \frac{\left(\mathbf{P}_{im}^t - \mathbf{P}_{jm}^t\right)^2}{\pi_m}$$
 (8)

where  $\mathbf{P}_{im}^t$  (4) denotes the probability of transitioning from node i to node m in t time steps. Intuitively, the diffusion distance measures the similarity between the probability distributions of random walks starting from nodes i and j. A key strength of diffusion maps is the equivalence (6):

$$D_t^2(i,j) = \left\| \mathbf{\Psi}_t(\mathbf{x}_i) - \mathbf{\Psi}_t(\mathbf{x}_j) \right\|_2^2, \tag{9}$$

which is fundamental to our approach, as it ensures that the Euclidean distances between every two points on the manifold, which we measure in  $\S 4.2$  and  $\S 4.3$ , align with dissimilarities between their matching high-dimensional points, represented by the diffusion distance (8). The embedding in (6) is truncated to its first d coordinates and discards the rest. This reduces noise sensitivity and retains the most meaningful geometric structures (Nadler et al., 2006). The mapping  $\Psi_t^{(d)}: \mathbb{R}^M \to \mathbb{R}^d$  gives:

$$\mathbf{\Psi}_t^{(d)}(\mathbf{x}_i) = \left(\lambda_1^t \mathbf{u}_1(i), \, \lambda_2^t \mathbf{u}_2(i), \, \dots, \, \lambda_d^t \mathbf{u}_d(i)\right)^T. \tag{10}$$

Consider  $\tau \in [0,1]$  as the minimal normalized retained sum of the eigenvalues, then d is given by:

$$d = \min \left\{ k \in \{1, \dots, N\} : \frac{\sum_{\ell=1}^{k} \lambda_{\ell}}{\sum_{\ell=1}^{N} \lambda_{\ell}} \ge \tau \right\}.$$
 (11)

#### 4 The Perceptual Separation and Perceptual Match Measures

# 4.1 Constructing Perceptual Clusters on the Manifold

The waveform reference signal of the i-th source,  $\mathbf{y}_i$ , undergoes  $N_p$  perceptual distortions, e.g., noise gating in different thresholds, vibrato in various rates, and a comb filter with several delay-gain pairs. Typically,  $N_p \in [60, 70]$ . We define the i-th distortion set  $\mathcal{D}_i$  as:

$$\mathcal{D}_i = \left\{ \hat{\mathbf{y}}_i, \mathbf{y}_i, \mathbf{y}_{i,1}, \dots, \mathbf{y}_{i,N_p} \right\}, \quad \forall p \in \{1, \dots, N_p\} : \mathbf{y}_{i,p} \in \mathbb{R}^L,$$
(12)

 with L from (1). Each waveform in  $\mathcal{D}_i$  is independently encoded via a pre-trained self-supervised model, e.g., wav2vec 2.0 (Baevski et al., 2020). Let  $\Phi: \mathbb{R}^L \to \mathbb{R}^M$  be this encoding operator, with M from §3, so  $\mathbf{x}_{i,p} = \Phi\left(\mathbf{y}_{i,p}\right)$ ,  $\mathbf{x}_i = \Phi\left(\mathbf{y}_i\right)$ ,  $\hat{\mathbf{x}}_i = \Phi\left(\hat{\mathbf{y}}_i\right)$ . Applying (12) across all  $N_f$  sources results in the high-dimensional set of representations:

$$\mathcal{X} = \left\{ \hat{\mathbf{x}}_i, \mathbf{x}_i, \mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,N_p} \mid i = 1, \dots, N_f \right\}, \tag{13}$$

with  $|\mathcal{X}| = N_f(N_p + 2) := N$ . We define the *i*-th perceptual cluster  $\mathcal{C}_i^{(d)}$  on manifold  $\mathcal{M}^{(d)}$  (10):

$$C_i^{(d)} = \left\{ \boldsymbol{\Psi}_t^{(d)}(\mathbf{x}_i), \boldsymbol{\Psi}_t^{(d)}(\mathbf{x}_{i,p}) \mid p = 1, \dots, N_p \right\}. \tag{14}$$

where we exclude the embedding of the system output  $\Psi_t^{(d)}(\hat{\mathbf{x}}_i) \in \mathbb{R}^d$  (10) from  $\mathcal{C}_i^{(d)}$ , since this embedding will be measured against the cluster statistics to produce the PS and PM measures. Including  $\Psi_t^{(d)}(\hat{\mathbf{x}}_i)$  in the cluster would create a circular dependency that will bias the PS and PM. These distortions were hand-crafted to create a wide perceptual auditory coverage relative to the reference, e.g., by considering mildly-intrusive additive colored noise with signal-to-noise-ratios (SNRs) of 15 dB on one hand, and severely degrading heavy-tailed reverberations on the other hand.

Given  $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{X}$ , the property in (9) guarantees that as the Euclidean distance between  $\mathbf{\Psi}_t^{(d)}(\mathbf{x}_i)$  and  $\mathbf{\Psi}_t^{(d)}(\mathbf{x}_j)$  lowers, so does the diffusion distance between  $\mathbf{x}_i$  and  $\mathbf{x}_j$ . In §4.2 and §4.3, we define our PS and PM measures using Euclidean distances, based on our hypothesis that this diffusion distance also aligns with the perceptual alignment between the corresponding waveforms,  $\mathbf{y}_i$  and  $\mathbf{y}_j$ . In §7, we explore if this perceptual-geometric hypothesis is valid by comparing our measures with human perception.

## 4.2 THE PERCEPTUAL SEPARATION (PS) MEASURE

For readability, we denote the elements of the clusters in (14) as  $\psi$  both here and in §4.3.

For source i, we aim to quantify the perceptual separation of  $\hat{\mathbf{y}}_i$  from its non-attributed references  $\{\mathbf{y}_j\}_{i\neq j}$  with the Mahalanobis distance (Evans et al., 2021). The empirical centroid and unbiased covariance matrix of the cluster  $\mathcal{C}_i^{(d)}$  are:

$$\hat{\boldsymbol{\mu}}_{j}^{(d)} = \frac{1}{\left|\mathcal{C}_{j}^{(d)}\right|} \sum_{\boldsymbol{\psi} \in \mathcal{C}_{i}^{(d)}} \boldsymbol{\psi}, \quad \widehat{\boldsymbol{\Sigma}}_{j}^{(d)} = \frac{1}{\left|\mathcal{C}_{j}^{(d)}\right| - 1} \sum_{\boldsymbol{\psi} \in \mathcal{C}_{i}^{(d)}} \left(\boldsymbol{\psi} - \hat{\boldsymbol{\mu}}_{j}^{(d)}\right) \left(\boldsymbol{\psi} - \hat{\boldsymbol{\mu}}_{j}^{(d)}\right)^{T}, \quad (15)$$

where  $\hat{\mu}_j^{(d)} \in \mathbb{R}^d$ ,  $\hat{\Sigma}_j^{(d)} \in \mathbb{R}^{d \times d}$ . The squared Mahalanobis distance from the embedding of the *i*-th output  $\Psi_t^{(d)}(\hat{\mathbf{x}}_i)$  to  $\mathcal{C}_i^{(d)}$  is given by:

$$d_{M}^{2}\left(\boldsymbol{\Psi}_{t}^{(d)}(\hat{\mathbf{x}}_{i});\hat{\boldsymbol{\mu}}_{j}^{(d)},\hat{\boldsymbol{\Sigma}}_{j}^{(d)}\right) = \left(\boldsymbol{\Psi}_{t}^{(d)}(\hat{\mathbf{x}}_{i}) - \hat{\boldsymbol{\mu}}_{j}^{(d)}\right)^{T}\left(\hat{\boldsymbol{\Sigma}}_{j}^{(d)} + \epsilon I^{(d)}\right)^{-1}\left(\boldsymbol{\Psi}_{t}^{(d)}(\hat{\mathbf{x}}_{i}) - \hat{\boldsymbol{\mu}}_{j}^{(d)}\right), (16)$$

where we use for regularization  $\epsilon = 10^{-6}$  with the *d*-dimensional identity matrix  $I^{(d)}$ . We define the measured Mahalanobis distance from  $\Psi_t^{(d)}(\hat{\mathbf{x}}_i)$  to its attributed and closest non-attributed clusters as:

$$\hat{A}_{i}^{(d)} = d_{M} \left( \mathbf{\Psi}_{t}^{(d)}(\hat{\mathbf{x}}_{i}); \hat{\boldsymbol{\mu}}_{i}^{(d)}, \hat{\boldsymbol{\Sigma}}_{i}^{(d)} \right), \quad \hat{B}_{i}^{(d)} = d_{M} \left( \mathbf{\Psi}_{t}^{(d)}(\hat{\mathbf{x}}_{i}); \hat{\boldsymbol{\mu}}_{j^{*}}^{(d)}, \hat{\boldsymbol{\Sigma}}_{j^{*}}^{(d)} \right), \tag{17}$$

with  $j^* = \arg\min_{j \in \{1,...,N_f\}, j \neq i} d_M\left(\Psi_t^{(d)}\left(\hat{\mathbf{x}}_i\right); \boldsymbol{\mu}_j^{(d)}, \boldsymbol{\Sigma}_j^{(d)}\right)$ . Notice that (17) resembles the source permutation minimization processing in source separation evaluations (Le Roux et al., 2019).

The measured PS score for  $\hat{\mathbf{y}}_i$  in the truncated dimension d is:

$$\widehat{PS}_{i}^{(d)} = 1 - \frac{\hat{A}_{i}^{(d)}}{\hat{A}_{i}^{(d)} + \hat{B}_{i}^{(d)}}, \quad \widehat{PS}_{i}^{(d)} \in [0, 1].$$
(18)

where by design  $\hat{A}_i^{(d)} + \hat{B}_i^{(d)} > 0$  and a higher score is better. Functionally, when  $\hat{A}_i^{(d)} \ll \hat{B}_i^{(d)}$  then the i-th output perceptually resembles its cluster members significantly more than competing cluster members and  $\widehat{\text{PS}}_i^{(d)}$  approaches 1.  $\hat{B}_i^{(d)} \ll \hat{A}_i^{(d)}$  indicates the opposite, and  $\widehat{\text{PS}}_i^{(d)}$  drops towards 0.

#### 4.3 THE PERCEPTUAL MATCH (PM) MEASURE

The PM measure aims to quantify how perceptually aligned the estimated output  $\hat{\mathbf{y}}_i$  is with its reference  $\mathbf{y}_i$ . Let  $\tilde{\mathcal{C}}_i^{(d)} = \mathcal{C}_i^{(d)} \setminus \Psi_t^{(d)}(\mathbf{x}_i)$  denote the reference-free *i*-th cluster. Unlike Equation (15), we compute the unbiased empirical covariance matrix of  $\tilde{\mathcal{C}}_i^{(d)}$  relative to its reference embedding:

$$\widehat{\widetilde{\Sigma}}_{i}^{(d)} = \frac{1}{\left|\widetilde{\mathcal{C}}_{i}^{(d)}\right| - 1} \sum_{\boldsymbol{\psi} \in \widetilde{\mathcal{C}}_{i}^{(d)}} \left(\boldsymbol{\psi} - \boldsymbol{\Psi}_{t}^{(d)}(\mathbf{x}_{i})\right) \left(\boldsymbol{\psi} - \boldsymbol{\Psi}_{t}^{(d)}(\mathbf{x}_{i})\right)^{T}.$$
(19)

Then, for  $p \in \{1, \dots, N_p\}$ , the squared Mahalanobis distance from the p-th distortion to its attributed reference in the i-th cluster, is given by  $d_M^2\left(\boldsymbol{\Psi}_t^{(d)}(\mathbf{x}_{i,p});\boldsymbol{\Psi}_t^{(d)}(\mathbf{x}_i),\widehat{\widetilde{\boldsymbol{\Sigma}}}_i^{(d)}\right)$ , following the definition in Equation (16). Let us define the set of distances:

$$\hat{\mathcal{G}}_{i}^{(d)} = \left\{ d_{M}^{2} \left( \mathbf{\Psi}_{t}^{(d)}(\mathbf{x}_{i,p}); \mathbf{\Psi}_{t}^{(d)}(\mathbf{x}_{i}), \widehat{\widetilde{\boldsymbol{\Sigma}}}_{i}^{(d)} \right) \mid p = 1, \dots, N_{p} \right\}.$$
 (20)

Empirically, we validated that nearly always these distances are well-approximated by a Gamma distribution, using Kolmogorov-Smirnov goodness-of-fit tests (Smirnov, 1948; Kolmogorov, 1986). The sample mean and unbiased variance of  $\hat{\mathcal{G}}_i^{(d)}$  are estimated by:

$$\hat{\mu}_{\mathcal{G}_{i}^{(d)}} = \frac{1}{\left|\hat{\mathcal{G}}_{i}^{(d)}\right|} \sum_{q \in \hat{\mathcal{G}}_{i}^{(d)}} g, \quad \hat{\sigma}_{\mathcal{G}_{i}^{(d)}}^{2} = \frac{1}{\left|\hat{\mathcal{G}}_{i}^{(d)}\right| - 1} \sum_{q \in \hat{\mathcal{G}}_{i}^{(d)}} \left(g - \hat{\mu}_{\mathcal{G}_{i}^{(d)}}\right)^{2}, \tag{21}$$

and can be moment-matched with a Gamma distribution, assuming  $\hat{\mu}_{\mathcal{G}_i^{(d)}}, \hat{\sigma}_{\mathcal{G}_i^{(d)}}^2 > 0$ , with parameters:

$$\hat{k}_{i}^{(d)} = \frac{\hat{\mu}_{\mathcal{G}_{i}^{(d)}}^{2}}{\hat{\sigma}_{\mathcal{G}_{i}^{(d)}}^{2}}, \quad \hat{\theta}_{i}^{(d)} = \frac{\hat{\sigma}_{\mathcal{G}_{i}^{(d)}}^{2}}{\hat{\mu}_{\mathcal{G}_{i}^{(d)}}}.$$
(22)

Similarly, the squared Mahalanobis distance from the output embedding to its attributed reference is  $\hat{a}_i^{(d)} = d_M^2 \left( \Psi_t^{(d)}(\hat{\mathbf{x}}_i); \Psi_t^{(d)}(\mathbf{x}_i), \widehat{\widetilde{\boldsymbol{\Sigma}}}_i^{(d)} \right)$ . Consider  $Q(k,x) = \Gamma(k,x)/\Gamma(k)$  as the regularized upper incomplete Gamma function (NIST, 2024). Then, the PM score for  $\hat{\mathbf{y}}_i$  in dimension d is:

$$\widehat{PM}_{i}^{(d)} = Q\left(\hat{k}_{i}^{(d)}, \frac{\hat{a}_{i}^{(d)}}{\hat{\theta}_{i}^{(d)}}\right), \quad \widehat{PM}_{i}^{(d)} \in [0, 1],$$
(23)

where  $\hat{k}_i^{(d)}, \hat{\theta}_i^{(d)}$  are well-defined by design for  $N_p \geq 1$  and a higher score is better. If the output  $\hat{a}_i^{(d)}$  lies well within the bulk of its distortion cluster, the Gamma-tail probability is near 1, which may indicate a strong perceptual match. As  $\hat{a}_i^{(d)}$  drifts away, the score decays smoothly toward zero, reflecting degradation. When distortions are tightly concentrated and  $\hat{k}_i^{(d)}$  or  $\hat{\theta}_i^{(d)}$  lower, even small mismatches in  $\hat{a}_i^{(d)}$  lower PM sharply. As  $\hat{k}_i^{(d)}$  and  $\hat{\theta}_i^{(d)}$  grow, the PM tolerates larger  $\hat{a}_i^{(d)}$  deviations.

# 5 ERROR GUARANTEES FOR THE PS AND PM MEASURES

Standing notation and assumptions. We fix frame f with generally  $N_f \geq 2$  (1). For this proof, consider the specific case of  $N_f = 2$  with indices  $i, j \in S_f$ . Consider source index j and set m = N - 1 - d as the dimension of the omitted space in the diffusion maps process, so the embedding notations in the retained, omitted, and complete N - 1-dimensional spaces are respectively  $\Psi_t^{(d)}(\mathbf{x}_j), \Psi_t^{(m)}(\mathbf{x}_j), \Psi_t(\mathbf{x}_j)$  (10). Similarly, clusters  $\mathcal{C}_j^{(d)}, \mathcal{C}_j^{(m)}$ , and  $\mathcal{C}_j$  are formed as in §4.1, with means and covariances pairs  $\left(\boldsymbol{\mu}_j^{(d)}, \boldsymbol{\Sigma}_j^{(d)}\right), \left(\boldsymbol{\mu}_j^{(m)}, \boldsymbol{\Sigma}_j^{(m)}\right)$ , and  $(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$ , and cross-covariance  $C_j \in \mathbb{R}^{d \times m}$ . We have empirically prevented ill-conditioning, as all matrix inversions are Tikhonov-regularized (Tikhonov & Arsenin, 1977) with  $\epsilon I$ , where  $\epsilon = 10^{-6}$  and I the identity matrix,

with context-dependent dimension. When we quantify sampling uncertainty, we use dependence-adjusted effective sample sizes via Bartlett method (Bartlett, 1946), and sub-Gaussian tails for quadratic forms via the dependent Hanson–Wright inequalities (Adamczak, 2015; Vershynin, 2024).

Schur decomposition of full versus truncated Mahalanobis distances. For radius error calculations, from (24) to (36), we assume access to clusters statistics. For the output embedding of source i against cluster j, define  $\Delta_{i,j}^{(d)} = \Psi_t^{(d)}(\hat{\mathbf{x}}_i) - \boldsymbol{\mu}_j^{(d)}$  and  $\Delta_{i,j}^{(m)} = \Psi_t^{(m)}(\hat{\mathbf{x}}_i) - \boldsymbol{\mu}_j^{(m)}$ . The full cluster statistics aggregate as:

$$\boldsymbol{\mu}_{j} = \begin{bmatrix} \boldsymbol{\mu}_{j}^{(d)} \\ \boldsymbol{\mu}_{j}^{(m)} \end{bmatrix}, \quad \boldsymbol{\Delta}_{i,j} = \begin{bmatrix} \boldsymbol{\Delta}_{i,j}^{(d)} \\ \boldsymbol{\Delta}_{i,j}^{(m)} \end{bmatrix}, \quad \boldsymbol{\Sigma}_{j} = \begin{bmatrix} \boldsymbol{\Sigma}_{j}^{(d)} & \boldsymbol{C}_{j} \\ \boldsymbol{C}_{j}^{T} & \boldsymbol{\Sigma}_{j}^{(m)} \end{bmatrix}. \tag{24}$$

Block inversion to (16) via the Schur complement (Horn & Johnson, 2013) yields:

$$d_M^2(\boldsymbol{\Psi}_t(\hat{\mathbf{x}}_i); \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) = \underbrace{\left(\boldsymbol{\Delta}_{i,j}^{(d)}\right)^T \left(\boldsymbol{\Sigma}_j^{(d)} + \epsilon I^{(d)}\right)^{-1} \boldsymbol{\Delta}_{i,j}^{(d)}}_{:=a} + \underbrace{\boldsymbol{r}_{i,j}^T \boldsymbol{S}_j^{-1} \boldsymbol{r}_{i,j}}_{:=b}, \tag{25}$$

$$\mathbf{r}_{i,j} = \mathbf{\Delta}_{i,j}^{(m)} - \mathbf{C}_{j}^{T} \left( \mathbf{\Sigma}_{j}^{(d)} + \epsilon I^{(d)} \right)^{-1} \mathbf{\Delta}_{i,j}^{(d)}, \quad \mathbf{S}_{j} = \mathbf{\Sigma}_{j}^{(m)} - \mathbf{C}_{j}^{T} \left( \mathbf{\Sigma}_{j}^{(d)} + \epsilon I^{(d)} \right)^{-1} \mathbf{C}_{j}.$$
 (26)

Since  $\forall a, b \geq 0 : |\sqrt{a+b} - \sqrt{a}| \leq \sqrt{b}$  (Rudin, 1976, Ch. 5), we bound truncation error to (25):

$$|\delta_{i,j}| := \left| d_M(\mathbf{\Psi}_t(\hat{\mathbf{x}}_i); \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j) - d_M^2 \left( \mathbf{\Psi}_t^{(d)}(\hat{\mathbf{x}}_i); \boldsymbol{\mu}_j^{(d)}, \boldsymbol{\Sigma}_j^{(d)} \right) \right| \le \sqrt{\boldsymbol{r}_{i,j}^T \boldsymbol{S}_j^{-1} \boldsymbol{r}_{i,j}}. \tag{27}$$

**PS radius.** Let  $A_i, B_i$  be the full-space versions of  $A_i^{(d)}, B_i^{(d)}$  in (17). Set  $|\delta_{i,i}| := |A_i - A_i^{(d)}|$  and  $|\delta_{i,j^*}| := |B_i - B_i^{(d)}|$ , with  $j^*$  as in (17). We empirically confirmed that truncation introduces only mild changes, i.e.,  $|\delta_{i,i}|, |\delta_{i,j^*}| \ll A_i^{(d)} + B_i^{(d)}$ . Thus, a first-order Taylor expansion of  $\mathrm{PS}_i$ , the full-space version of  $\mathrm{PS}_i^{(d)}$ , around  $\left(A_i^{(d)}, B_i^{(d)}\right)$ , is valid. Ultimately, we drop quadratic components that were found negligible, and use  $|\delta_{i,i}|$  and  $|\delta_{i,j^*}|$  inside the Taylor expansion, to yield:

$$\left| \text{PS}_i - \text{PS}_i^{(d)} \right| \le \frac{B_i^{(d)} \left| \delta_{i,i} \right| + A_i^{(d)} \left| \delta_{i,j^*} \right|}{\left( A_i^{(d)} + B_i^{(d)} \right)^2}.$$
 (28)

Combining with (27), the deterministic PS error radius is:

$$\left| \operatorname{PS}_{i} - \operatorname{PS}_{i}^{(d)} \right| \leq \frac{B_{i}^{(d)} \sqrt{\boldsymbol{r}_{i,i}^{T} \boldsymbol{S}_{i}^{-1} \boldsymbol{r}_{i,i}} + A_{i}^{(d)} \sqrt{\boldsymbol{r}_{i,j^{*}}^{T} \boldsymbol{S}_{j^{*}}^{-1} \boldsymbol{r}_{i,j^{*}}}}{\left(A_{i}^{(d)} + B_{i}^{(d)}\right)^{2}}.$$
 (29)

We notice that large residual spread  $\Sigma_i^{(m)}$  or cross-block coupling  $C_i$  inflate (29) through  $S_i^{-1}$ .

**PM radius.** For source i and every distortion index  $p \in \{1, \dots, N_p\}$ , we center cluster coordinates at the reference  $\mathbf{x}_i$ , so  $\boldsymbol{\Delta}_{i,p}^{(d)} = \boldsymbol{\Psi}_t^{(d)}(\mathbf{x}_{i,p}) - \boldsymbol{\Psi}_t^{(d)}(\mathbf{x}_i)$  and  $\boldsymbol{\Delta}_{i,p}^{(m)} = \boldsymbol{\Psi}_t^{(m)}(\mathbf{x}_{i,p}) - \boldsymbol{\Psi}_t^{(m)}(\mathbf{x}_i)$ . By repeating (24) for the reference-free clusters in the d, m, and N-1 spaces, we obtain:

$$\widetilde{\Sigma}_{i} = \begin{bmatrix} \widetilde{\Sigma}_{i}^{(d)} & \widetilde{C}_{i} \\ \widetilde{C}_{i}^{T} & \widetilde{\Sigma}_{i}^{(m)} \end{bmatrix}, \tag{30}$$

where  $\widetilde{\Sigma}_i^{(d)}$  is defined in (19). Exactly as in (25)-(26), we use Schur complement:

$$d_M^2\left(\boldsymbol{\Psi}_t(\mathbf{x}_{i,p});\boldsymbol{\Psi}_t(\mathbf{x}_i),\widetilde{\boldsymbol{\Sigma}}_i\right) = \left(\boldsymbol{\Delta}_{i,p}^{(d)}\right)^T \left(\widetilde{\boldsymbol{\Sigma}}_i^{(d)} + \epsilon I\right)^{-1} \boldsymbol{\Delta}_{i,p}^{(d)} + \boldsymbol{r}_{i,p}^T \boldsymbol{S}_i^{-1} \boldsymbol{r}_{i,p}, \tag{31}$$

$$\boldsymbol{r}_{i,p} = \boldsymbol{\Delta}_{i,p}^{(m)} - \widetilde{\boldsymbol{C}}_{i}^{T} \left( \widetilde{\boldsymbol{\Sigma}}_{i}^{(d)} + \epsilon \boldsymbol{I} \right)^{-1} \boldsymbol{\Delta}_{i,p}^{(d)}, \quad \boldsymbol{S}_{i} = \widetilde{\boldsymbol{\Sigma}}_{i}^{(m)} - \widetilde{\boldsymbol{C}}_{i}^{T} \left( \widetilde{\boldsymbol{\Sigma}}_{i}^{(d)} + \epsilon \boldsymbol{I} \right)^{-1} \widetilde{\boldsymbol{C}}_{i}.$$
(32)

Let  $\mathcal{G}_i$  be the set of the squared distances in (31) over p and  $\mathcal{G}_i^{(d)}$  its d-dimensional analogue (20). Define per-sample truncation gaps  $\delta_{\mathcal{G}_i,p} := \boldsymbol{r}_{i,p}^T \boldsymbol{S}_i^{-1} \boldsymbol{r}_{i,p} \geq 0$  and  $\delta_{\max} = \max_p \delta_{\mathcal{G}_i,p}$ . Employing elementary algebra and Cauchy–Schwarz inequality (Vershynin, 2024), we obtain the relations (21):

$$\left| \mu_{\mathcal{G}_i} - \mu_{\mathcal{G}_i^{(d)}} \right| = \frac{1}{N_p} \sum_{p=1}^{N_p} \delta_{\mathcal{G}_i, p}, \quad \left| \sigma_{\mathcal{G}_i}^2 - \sigma_{\mathcal{G}_i^{(d)}}^2 \right| \le \frac{N_p}{N_p - 1} \left( 2\delta_{\max} \left( \sigma_{\mathcal{G}_i} + \sigma_{\mathcal{G}_i^{(d)}} \right) + \delta_{\max}^2 \right). \tag{33}$$

Again, simple algebra bounds the Gamma-matching parameters (22), with constants  $C_1, C_2 > 0$ :

$$\left| k_i - k_i^{(d)} \right| \le C_1 \, \delta_{\max} \frac{N_p}{N_p - 1} \frac{\mu_{\mathcal{G}_i} + \mu_{\mathcal{G}_i^{(d)}}}{\sigma_{\mathcal{G}_i^{(d)}}^2}, \quad \left| \theta_i - \theta_i^{(d)} \right| \le C_2 \, \delta_{\max} \frac{N_p}{N_p - 1} \frac{\sigma_{\mathcal{G}_i}^2 + \sigma_{\mathcal{G}_i^{(d)}}^2}{\mu_{\mathcal{G}_i^{(d)}}^2}. \quad (34)$$

Recalling the distance of the output from its cluster, denoted  $a_i^{(d)}$  (23), we can define using (32):

$$d_M^2\left(\boldsymbol{\Psi}_t(\hat{\mathbf{x}}_i);\boldsymbol{\Psi}_t(\mathbf{x}_i),\tilde{\boldsymbol{\Sigma}}_i\right) - d_M^2\left(\boldsymbol{\Psi}_t^{(d)}(\hat{\mathbf{x}}_i);\boldsymbol{\Psi}_t^{(d)}(\mathbf{x}_i),\tilde{\boldsymbol{\Sigma}}_i^{(d)}\right) = \boldsymbol{r}_{i,a}^T\boldsymbol{S}_i^{-1}\boldsymbol{r}_{i,a} := \delta_{\mathcal{G}_i,a}.$$
(35)

In the full space,  $\mathrm{PM}_i = Q(k_i, a_i/\theta_i)$  (23). Its derivatives with respect to its variables are standard and bounded on compact sets (NIST, 2024). Let the truncation ellipsoid  $\mathcal{B}_i$  be such set, so that  $\mathcal{B}_i = \left\{ (k_i', \theta_i', a_i') : \left| k_i' - k_i^{(d)} \right| \le \delta_{\mathcal{G}_i,k}, \left| \theta_i' - \theta_i^{(d)} \right| \le \delta_{\mathcal{G}_i,\theta}, \left| a_i' - a_i^{(d)} \right| \le \delta_{\mathcal{G}_i,a} \right\}$ , with  $\delta_{\mathcal{G}_i,k}, \delta_{\mathcal{G}_i,\theta}$  denoting the bounds in (34). Since Q increases in k and decreases in  $k = a/\theta$  for k, k > 0 (NIST, 2024), the maximum deviation over  $\mathcal{B}_i$  occurs at a corner, and the radius can be obtained by:

$$\left| \text{PM}_i - \text{PM}_i^{(d)} \right| \le \max_{(k_c, \theta_c, a_c) \in \partial \mathcal{B}_i} \left| Q(k_c, a_c/\theta_c) - Q(k_i^{(d)}, a_i^{(d)}/\theta_i^{(d)}) \right|.$$
 (36)

**Dependence-adjusted sample size.** For any cluster  $C_j^{(d)}$  with  $n_j$  dependent points, we use  $n_{j,\text{eff}} := n_j \left(1 + 2\sum_{\ell=1}^{L_j} \hat{\rho}_{j,\ell}\right)^{-1}$  with  $L_j = \min\{\ell : |\hat{\rho}_{j,\ell}| < z_{0.975}/\sqrt{n_j - \ell}\}$  (Bartlett, 1946).

**PS tail bound.** We now resort to the retained d-dimensional space, and estimate cluster statistics due the finite number of  $n_{j,\text{eff}}$  samples. Let  $\widehat{\mu}_{j}^{(d)}$  and  $\widehat{\Sigma}_{j}^{(d)}$  be the empirical cluster statistics. Vector and matrix-Bernstein and dependent Hanson–Wright theories (Vershynin, 2024, Props. 2.8.1, 4.7.1), (Adamczak, 2015, Thm. 2.5) give for  $\delta_{j,\mu}^{\text{PS}}$ ,  $\delta_{j,\Sigma}^{\text{PS}} \in (0,1/2)$ , with least probabilities  $1 - \delta_{j,\mu}^{\text{PS}}$ ,  $1 - \delta_{j,\Sigma}^{\text{PS}}$ .

$$\left\| \boldsymbol{\mu}_{j}^{(d)} - \widehat{\boldsymbol{\mu}}_{j}^{(d)} \right\|_{2} \le \sqrt{2 \lambda_{\max} \left(\widehat{\boldsymbol{\Sigma}}_{j}^{(d)}\right) \ln(2/\delta_{j,\mu}^{\mathrm{PS}})/n_{j,\mathrm{eff}}} =: \Delta_{j,\mu}, \tag{37}$$

$$\left\| \mathbf{\Sigma}_{j}^{(d)} - \widehat{\mathbf{\Sigma}}_{j}^{(d)} \right\|_{2} \le C \lambda_{\max} \left( \widehat{\mathbf{\Sigma}}_{j}^{(d)} \right) r_{j} + \ln(2/\delta_{j,\Sigma}^{\text{PS}}) / n_{j,\text{eff}} =: \Delta_{j,\Sigma}, \tag{38}$$

with  $r_j = \operatorname{trace}(\widehat{\Sigma}_j^{(d)})/\lambda_{\max}\left(\widehat{\Sigma}_j^{(d)}\right)$ . A first-order perturbation of  $A_i^{(d)}, B_i^{(d)}$  (17) yields the bounds:

$$\varepsilon^{\mathrm{PS}}\left(\widehat{A}_{i}^{(d)}\right) \leq 2\sqrt{\widehat{A}_{i}^{(d)}}\Delta_{i,\mu}\sqrt{\lambda_{\mathrm{max}}\left(\widehat{\boldsymbol{\Sigma}}_{i}^{(d)}\right)/\widetilde{\lambda}_{\mathrm{min}}\left(\widehat{\boldsymbol{\Sigma}}_{i}^{(d)}\right)} + \widehat{A}_{i}^{(d)}\Delta_{i,\Sigma}/\lambda_{\mathrm{max}}\left(\widehat{\boldsymbol{\Sigma}}_{i}^{(d)}\right),\tag{39}$$

$$\varepsilon^{\text{PS}}\left(\widehat{B}_{i}^{(d)}\right) \leq 2\sqrt{\widehat{B}_{i}^{(d)}}\Delta_{j^{*},\boldsymbol{\mu}}\sqrt{\lambda_{\text{max}}\left(\widehat{\boldsymbol{\Sigma}}_{j^{*}}^{(d)}\right)/\widetilde{\lambda}_{\text{min}}\left(\widehat{\boldsymbol{\Sigma}}_{j^{*}}^{(d)}\right)} + \widehat{B}_{i}^{(d)}\Delta_{j^{*},\boldsymbol{\Sigma}}/\lambda_{\text{max}}\left(\widehat{\boldsymbol{\Sigma}}_{j^{*}}^{(d)}\right). \tag{40}$$

With  $L_i^{\text{PS}}$  being the Euclidean gradient norm of  $\text{PS}_i^{(d)}$  at  $(A_i^{(d)}, B_i^{(d)})$ , for  $\delta_i^{\text{PS}} = \delta_{i,\mu}^{\text{PS}} + \delta_{i,\Sigma}^{\text{PS}} \in (0,1)$ 

$$\left|\widehat{\mathrm{PS}}_{i}^{(d)} - \mathrm{PS}_{i}^{(d)}\right| \leq L_{i}^{\mathrm{PS}} \sqrt{\varepsilon^{\mathrm{PS}} \left(\widehat{A}_{i}^{(d)}) + \varepsilon^{\mathrm{PS}} (\widehat{B}_{i}^{(d)})} \quad \text{w.p. } \geq 1 - \delta_{i}^{\mathrm{PS}}. \tag{41}$$

**PM tail bound.** Let  $R_i = \max_{g \in \mathcal{G}_i} g$  and choose confidence levels  $\delta_{i,\mu}^{\mathrm{PM}}, \delta_{i,\sigma}^{\mathrm{PM}}, \delta_{i,a}^{\mathrm{PM}} \in (0,1/3)$ . Concentration bounds for the output quadratic form via Hanson–Wright (Vershynin, 2024, Props. 2.8.1, 4.7.1) yields

$$\left| \mu_{\mathcal{G}_{i}^{(d)}} - \widehat{\mu}_{\mathcal{G}_{i}^{(d)}} \right| \le \sqrt{\frac{2\widehat{\sigma}_{\mathcal{G}_{i}^{(d)}}^{2} \ln(2/\delta_{i,\mu}^{\text{PM}})}{N_{p}} + \frac{3R_{i} \ln(2/\delta_{i,\mu}^{\text{PM}})}{N_{p}}},$$
 (42)

$$\left|\sigma_{\mathcal{G}_{i}^{(d)}} - \widehat{\sigma}_{\mathcal{G}_{i}^{(d)}}\right| \leq \sqrt{\frac{2R_{i}^{2}\ln(2/\delta_{i,\sigma}^{\mathrm{PM}})}{N_{p}}} + \frac{3R_{i}^{2}\ln(2/\delta_{i,\sigma}^{\mathrm{PM}})}{N_{p}}, \left|a_{i}^{(d)} - \widehat{a}_{i}^{(d)}\right| \leq R_{i}\sqrt{\frac{\ln(2/\delta_{i,a}^{\mathrm{PM}})}{N_{p}}}.$$
(43)

Just like in (39) and (40), these are mapped to bounds on the parameters  $(k,\theta,a)$ , yielding  $\Delta_{i,k}, \Delta_{i,\theta}, \Delta_{i,a}$ . Consider  $\delta_i^{\mathrm{PM}} = \delta_{i,\mu}^{\mathrm{PM}} + \delta_{i,\sigma}^{\mathrm{PM}} + \delta_{i,a}^{\mathrm{PM}} \in (0,1)$ , then for the local box  $\mathcal{B}_i^{\mathrm{loc}} = \{ |k - k_i^{(d)}| \leq \Delta_{i,k}, |\theta - \theta_i^{(d)}| \leq \Delta_{i,\theta}, |a - a_i^{(d)}| \leq \Delta_{i,a} \}$ :

$$\left|\widehat{\mathrm{PM}}_{i}^{(d)} - \mathrm{PM}_{i}^{(d)}\right| \leq \max_{(k_{c}, \theta_{c}, a_{c}) \in \partial \mathcal{B}_{i}^{\mathrm{loc}}} \left| Q(k_{c}, a_{c}/\theta_{c}) - Q(\widehat{k}_{i}^{(d)}, \widehat{a}_{i}^{(d)}/\widehat{\theta}_{i}^{(d)}) \right| \quad \text{w.p. } \geq 1 - \delta_{i}^{\mathrm{PM}}. \tag{44}$$

Table 1: SRCC and PCC of the PS and PM measures (underlined), their waveform counterparts, and 14 comparative measures, across scenarios. The top-3 results in every column are in bold.

	English		Spanish		Music (Drums)		Music (No Drums)	
Measure	SRCC	PCC	SRCC	PCC	SRCC	PCC	SRCC	PCC
PS	84.12%	83.74%	82.33%	85.01%	72.87%	77.38%	87.23%	87.81%
PM	<b>84.69</b> %	<del>86.36</del> %	83.41%	<b>85.30</b> %	<b>75.18</b> %	<b>69.88</b> %	<b>88.12</b> %	<b>85.26</b> %
PS (waveform)	$\overline{73.42}\%$	$\overline{71.04}\%$	74.69%	$\overline{75.05}\%$	51.75%	$\overline{61.83}\%$	$\overline{78.88}\%$	$\overline{78.95}\%$
PM (waveform)	69.30%	66.62%	68.27%	67.35%	49.52%	51.77%	74.37%	75.51%
STOI	80.85%	78.40%	78.79%	82.56%	67.29%	71.27%	75.64%	78.13%
eSTOI	82.14%	82.28%	79.20%	82.68%	54.68%	57.35%	70.06%	74.45%
PESQ	85.56%	84.05%	86.06%	84.98%	61.60%	53.87%	61.26%	60.24%
SI-SDR	78.11%	76.96%	84.07%	81.38%	42.08%	56.98%	70.42%	71.96%
SDR	77.72%	73.13%	84.29%	76.07%	44.78%	54.33%	74.51%	75.35%
SIR	51.28%	56.20%	45.67%	55.19%	18.64%	35.76%	51.00%	55.12%
SAR	75.54%	72.98%	78.21%	73.29%	36.98%	40.81%	66.15%	68.96%
CI-SDR	78.66%	77.41%	84.32%	81.48%	45.02%	55.42%	74.25%	75.11%
DNSMOS-OVRL	63.70%	67.77%	35.34%	43.57%	21.79%	34.27%	13.81%	19.47%
MCD	43.05%	33.86%	45.90%	37.97%	30.27%	42.23%	33.49%	32.19%
SpeechBERTscore	68.58%	67.44%	69.55%	70.48%	52.33%	59.71%	75.60%	81.13%
Sheet-SSQA	41.17%	51.38%	61.06%	73.01%	39.40%	29.03%	14.19%	5.17%
UTMOS	55.53%	55.43%	52.22%	55.75%	-9.24%	-8.25%	12.59%	7.72%
NISQA	60.78%	67.62%	63.37%	66.58%	27.27%	41.73%	42.33%	48.07%

# 6 EXPERIMENTAL SETUP

#### 6.1 Database

 We use the Subjective Evaluation of Blind Audio Source Separation (SEBASS) database (Kastner & Herre, 2022), a public collection of expertly curated listening tests that aggregates 11,000 ratings for more than 900 separated signals across five evaluation campaigns. SEBASS covers speech mixtures of 4 male or 4 female speakers, each consisting of English and Spanish pairs. As realistic conversations are monolingual, we separate each mixture into English and Spanish speakers pairs. Also included are music mixtures with drums and without drums, each with 3 sources. Namely,  $N_f \in \{2,3\}$  (1). The split between drum and no-drum mixtures is crucial, as percussion transients create perceptual and algorithmic masking distinct from harmonic content. Each mixture was processed by 32 source separation systems, ranging from classic approaches to deep-learning models. Outputs with 10 s duration, sampled at 16 kHz, were judged by 15 certified raters under the MUSHRA standard (Schoeffler et al., 2018), which grades output quality between 0 and 100 relative to a reference.

#### 6.2 Pre-processing and Performance Evaluation

SEBASS provides MOS values at the utterance level. Since our PM and PS measures operate at much finer temporal resolutions, with frame sizes of L=400 for speech and L=324 samples for music (1), aggregation from the frame-level to the utterance-level is required to enable comparison with human MOS. PM values are aggregated using a simple average, while PS values are aggregated with a perceptually-weighted scheme inspired by PESQ. For performance evaluation, we correlate the aggregated PM and PS values with the utterance-level MOS values using the Pearson product-moment correlation coefficient (PCC) (Benesty et al., 2009) and the Spearman rank-order correlation coefficient (SRCC) (Sedgwick, 2014). We set  $\alpha=1$  (3) to eliminate density-dependent bias from the embedding, and t=1 (6) to keep the diffusion operator focused on local structures. The retained dimension d is in [20,40] (10), using  $\tau=0.99$  (11), as done on (Fjellström & Nyström, 2022).

#### 7 EXPERIMENTAL RESULTS

Results are from zero-shot SEBASS inference, without training or data-driven parameter tuning.

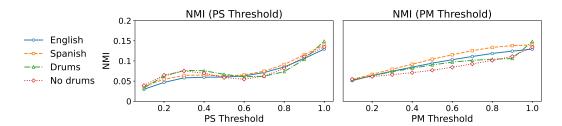


Figure 1: NMI between the PS and PM measures across their thresholded values.

Table 2: Deterministic error radius and probabilistic 95% CI of the SRCC and PCC across scenarios.

	English		Spanish		Music (Drums)		Music (No Drums)	
Measure	SRCC	PCC	SRCC	PCC	SRCC	PCC	SRCC	PCC
PS radius	0.16%	0.21%	0.10%	0.14%	0.40%	0.72%	0.14%	0.11%
PS CIs (95%)	30.03%	10.29%	26.39%	8.85%	28.71%	12.21%	12.69%	4.11%
PM radius	0.11%	0.99%	0.18%	1.23%	0.29%	1.39%	0.02%	1.04%
PM CIs (95%)	7.23%	3.83%	8.98%	4.28%	6.25%	4.15%	4.75%	1.77%

Table 1 benchmarks the proposed PS and PM measures against 14 widely-used metrics for audio quality and also versus its waveform-only version, denoted PS (waveform) and PM (waveform). For speech, we used a wav2vec 2.0-based (Baevski et al., 2020) model with features of dimension M=1024 (§3) and 24 transformer layers, and for music we use the MERT model (Li et al., 2023) with M=768 and 12 transformer layers. Since earlier layers are more stable, we pick layer 2 for speech, layer 1 for drums, and layer 3 for no-drums music. PS and PM consistently achieve top PCC values, aside from minor advantages by PESQ and STOI. For SRCC, our measures dominate in music, but trail PESQ in English and SDR-based metrics in Spanish. These results position the PS and PM as valid measures for leakage and self-distortion for source separation systems. Encoding proves essential, as waveform-only variants perform worse. Finally, PS and PM outperform SpeechBERTScore, showing the benefit of diffusion maps over cosine similarity.

We examine the complementary relationship between the PS and PM using NMI (Danon et al., 2005), which captures statistical dependence beyond linear effects, with lower NMI indicating less shared information. Each measure is normalized per utterance to [0,1]. For thresholds  $\{0.1,0.2,\ldots,1\}$ , we retain frames with PS below the threshold and compute the NMI between aligned PS-PM pairs. The procedure is repeated with thresholding on the PM. Figure 1 shows the NMI decreases toward zero as thresholds tighten, suggesting that the PS and PM become more complementary when separation quality is poor. At the loosest thresholds, NMI rises up to 0.15, yet full redundancy corresponds to 1. These results support reporting both PS and PM, as each captures failure modes missed by the other.

Frame-level error bounds of the measures were derived in  $\S5$ . Table 2 presents their propagation into PCC and SRCC error bounds. The error radius never exceeds 1.39%, a bias that rarely affects the performance ranking in Table 1. The 95% CIs highlight that the PS carries higher statistical uncertainty, whereas the PM is statistically more robust. This positions the PS as a complementary diagnostic, capturing perceptual leakage that the PM misses, at the cost of greater variability.

#### 8 Conclusions

We introduced the PS and PM, frame-level measures that showed competitive correlations with human MOS for source separation evaluation by operating on diffusion map embeddings of self-supervised audio representations. We derived a deterministic truncation bias and non-asymptotic CIs for both measures, making scores interpretable under quantified uncertainty. Looking forward, PS and PM can serve as diagnostic tools to localize whether errors stem from target distortion or cross-talk, while their differentiability enables use as loss terms or curriculum triggers to balance fidelity and separation under confidence monitoring. Finally, their uncertainty bounds offer a principled layer for benchmarking, supporting fairer hyper-parameter sweeps and reporting standards.

**Ethics Statement.** This work does not involve human subjects or personally identifiable information. We use only existing datasets under their respective licenses and terms of use, and we do not redistribute any data where licenses restrict sharing. Our study complies with the ICLR Code of Ethics. We assessed foreseeable risks and did not identify specific, material harms arising from the methods or results presented here. We will release implementation details and scripts to support responsible reuse and verification.

**Reproducibility Statement.** We provide complete code as an anonymous supplementary material in a separate .zip file. It contains the complete inference pipeline, including the frame-level calculation of the PS and PM measures and their determinstic and probabilistic error bounds.

#### REFERENCES

- Radosław Adamczak. A note on the Hanson–Wright inequality for random vectors with dependencies. *Electronic Journal of Probability*, 20:1–13, 2015. doi: 10.1214/EJP.v20-3412.
- Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. In *Proc. Advances in Neural Information Processing Systems*, volume 33, pp. 12449–12460, 2020.
- Scott Bannister, Alinka E. Greasley, Trevor J. Cox, Michael A. Akeroyd, Jon Barker, Bruno Fazenda, Jennifer Firth, Simone N. Graetzer, Gerardo Roa Dabike, Rebecca R. Vos, et al. Muddy, muddled, or muffled? understanding the perception of audio quality in music by hearing aid users. *Frontiers in Psychology*, 15:1310176, 2024.
- Maurice S. Bartlett. On the theoretical specification and sampling properties of autocorrelated time-series. *Supplement to the Journal of the Royal Statistical Society*, 8(1):27–41, 1946. doi: 10.2307/2983611. URL https://www.jstor.org/stable/2983611.
- Jacob Benesty, Jingdong Chen, Yiteng Huang, and Israel Cohen. Pearson correlation coefficient. In *Noise reduction in speech processing*, pp. 1–4. Springer, 2009.
- Jean Carletta, Simone Ashby, Sebastien Bourban, Mike Flynn, Mael Guillemot, Thomas Hain, Jaroslav Kadlec, Vasilis Karaiskos, Wessel Kraaij, Melissa Kronenthal, et al. The AMI meeting corpus: A pre-announcement. In *Proc. International workshop on machine learning for multimodal interaction*, pp. 28–39. Springer, 2005.
- Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, Jian Wu, Long Zhou, Shuo Ren, Yanmin Qian, Yao Qian, Jian Wu, Michael Zeng, and Furu Wei. 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. doi: 10.1109/JSTSP.2022.3187672.
- Ronald R. Coifman and Stéphane Lafon. Diffusion maps. *Applied and computational harmonic analysis*, 21(1):5–30, 2006.
- Alexis Conneau, Alexei Baevski, Ronan Collobert, Abdelrahman Mohamed, and Michael Auli. Unsupervised cross-lingual representation learning for speech recognition. *arXiv preprint arXiv:2006.13979*, 2020.
- Leon Danon, Albert Díaz-Guilera, Jordi Duch, and Alex Arenas. Comparing community structure identification. *Journal of Statistical Mechanics: Theory and Experiment*, (09), 2005. doi: 10.1088/1742-5468/2005/09/P09008.
- Bradley Efron and Robert J. Tibshirani. *An Introduction to the Bootstrap*, volume 57 of *Monographs on Statistics and Applied Probability*. Chapman & Hall/CRC, New York, 1994. ISBN 9780412042317.
- European Broadcasting Union (EBU). EBU recommendation R128: Loudness normalisation and permitted maximum level of audio signals. Technical recommendation, European Broadcasting Union, Geneva, Switzerland, 2011. URL https://tech.ebu.ch/docs/r/r128.pdf.
- Luke Evans, Maria K. Cameron, and Pratyush Tiwary. Computing committors in collective variables via mahalanobis diffusion maps. *arXiv* preprint arXiv:2108.08979, Aug 2021.
- Carmina Fjellström and Kaj Nyström. Deep learning, stochastic gradient descent and diffusion maps. *Journal of Computational Mathematics and Data Science*, 4(1), 2022.
- Toshiaki Fukada, Keiichi Tokuda, Takao Kobayashi, and Satoshi Imai. An adaptive algorithm for mel-cepstral analysis of speech. In *Proc. ICASSP*, volume 92, pp. 137–140, 1992.
- Alf Gabrielsson and Håkan Sjögren. Perceived sound quality of sound-reproducing systems. *The Journal of the Acoustical Society of America*, 65(4):1019–1033, 1979.

- Matthias Hein, Jean-Yves Audibert, and Ulrike von Luxburg. From graphs to manifolds weak and strong pointwise consistency of graph laplacians. In *Proc. COLT*, volume 3559 of *Lecture Notes in Computer Science*, pp. 470–485. Springer, 2005. URL https://doi.org/10.1007/11503415\_32.
  - Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, Cambridge, UK, 2nd edition, 2013. ISBN 9781107033413.
  - 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. doi: 10.1109/TASLP.2021.3122291.
  - Wen-Chin Huang, Erica Cooper, and Tomoki Toda. SHEET: A multi-purpose open-source speech human evaluation estimation toolkit. *arXiv preprint arXiv:2505.15061*, 2025.
  - Yu-Wen Hung, Szu-Wei Fu Cheng, Yu Tsao, and Hsin-Min Wang. Boosting self-supervised embeddings for speech quality assessment. In *Proc. Interspeech*, pp. 2288–2292, 2022.
  - ITU-T. ITU-T P.800: Methods for subjective determination of transmission quality. Technical report, Inter. Telecommunication Union, 1996. URL https://www.itu.int/rec/T-REC-P.800/en.
  - ITU-T. ITU-T P.835: Subjective test methodology for evaluating speech communication systems that include noise suppression algorithm. Technical report, Inter. Telecommunication Union, 2003. URL https://www.itu.int/rec/T-REC-P.835/en.
  - ITU-T. ITU-T P.862.2: Wideband extension to Recommendation P.862 for the assessment of wideband telephone networks and speech codecs. Recommendation, International Telecommunication Union, Telecommunication Standardization Sector (ITU-T), Geneva, Switzerland, November 2007. URL https://handle.itu.int/11.1002/1000/9275.
  - ITU-T. ITU-T P.808: Subjective evaluation of speech quality with a crowdsourcing approach. Technical report, Inter. Telecommunication Union, 2018. URL https://www.itu.int/rec/T-REC-P.808/en.
  - Ute Jekosch. Basic concepts and terms of "quality", reconsidered in the context of product-sound quality. *Acta Acustica united with Acustica*, 90(6):999–1006, 2004.
  - Jesper Jensen and Cees H. Taal. An algorithm for predicting the intelligibility of speech masked by modulated noise maskers. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 24(11):2009–2022, 2016. doi: 10.1109/TASLP.2016.2585878.
  - Thorsten Kastner and Jürgen Herre. The SEBASS-DB: A consolidated public data base of listening test results for perceptual evaluation of BSS quality measures. In *Proc. International Workshop on Acoustic Signal Enhancement (IWAENC)*, pp. 1–5. IEEE, 2022.
  - Maurice Kendall and Jean D. Gibbons. *Rank Correlation Methods*. Edward Arnold, London, 5 edition, 1990. ISBN 9780340506326.
  - Vahid Khanagha, Dimitrios Koutsaidis, Kshitij Kalgaonkar, and Sridha Srinivasan. Interference Aware Training Target for DNN-based Joint Acoustic Echo Cancellation and Noise Suppression. In *Proc. Interspeech*, pp. 1–5. ISCA, 2024. URL https://www.isca-archive.org/interspeech\_2024/khanagha24\_interspeech.pdf.
  - Andrey N. Kolmogorov. On the empirical determination of a distribution law. In Albert N. Shiryaev (ed.), *Selected Works of A. N. Kolmogorov, Volume II*, pp. 139–146. Springer, New York, 1986.
  - Jonathan Le Roux, Scott Wisdom, Hakan Erdogan, and John R. Hershey. SDR-half-Baked or Well Done? In *Proc. ICASSP*, pp. 626–630, 2019.
  - 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.

- Omer Moussa and Mariya Toneva. Brain-tuned speech models better reflect speech processing stages in the brain. *arXiv preprint arXiv:2506.03832*, June 2025.
- Boaz Nadler, Stéphane Lafon, Ronald R Coifman, and Ioannis G Kevrekidis. Diffusion maps, spectral clustering and reaction coordinates of dynamical systems. *Applied and Computational Harmonic Analysis*, 21(1):113–127, 2006.
  - NIST. Regularized incomplete gamma functions. http://dlmf.nist.gov/8.2, 2024.
  - Ankita Pasad, Alexei Baevski, Emmanuel Dupoux, and Karen Livescu. Comparative layer-wise analysis of self-supervised speech models. In *Proc. Interspeech*, 2022.
    - Chandan K. A. Reddy, Vishak Gopal, and Ross Cutler. DNSMOS P.835: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors. *arXiv preprint arXiv:2110.01763*, 2022.
    - Anthony W. Rix, John G. Beerends, Michael P. Hollier, and Andries P. Hekstra. Perceptual evaluation of speech quality (PESQ)—a new method for speech quality assessment of telephone networks and codecs. In *Proc. ICASSP*, volume 2, pp. 749–752, 2001.
    - Walter Rudin. *Principles of Mathematical Analysis*. International Series in Pure and Applied Mathematics. McGraw–Hill, New York, 3 edition, 1976.
    - Takaaki Saeki, Soumi Maiti, Shinnosuke Takamichi, Shinji Watanabe, and Hiroshi Saruwatari. SpeechBERTScore: Reference-aware automatic evaluation of speech generation leveraging NLP evaluation metrics. In *Proc. Interspeech*, pp. 4943–4947, 2024. doi: 10.21437/Interspeech. 2024-1508.
    - Michael Schoeffler, Sarah Bartoschek, Fabian-Robert Stoter, Marlene Roess, Susanne Westphal, Bernd Edler, and Jürgen Herre. webMUSHRA a comprehensive framework for web-based listening tests. *Journal of Open Research Software*, 6(1), 2018.
    - Manfred R. Schroeder. A new method of measuring reverberation time. *Journal of the Acoustical Society of America*, 37(2):409–412, 1965. doi: 10.1121/1.1909343.
    - Philip Sedgwick. Spearman's rank correlation coefficient. *BMJ*, 349:g7327, November 28 2014. doi: 10.1136/bmj.g7327.
    - Jiatong Shi, Hye-jin Shim, Jinchuan Tian, Siddhant Arora, Haibin Wu, Darius Petermann, Jia Qi Yip, You Zhang, Yuxun Tang, Wangyou Zhang, Dareen Safar Alharthi, Yichen Huang, Koichi Saito, Jionghao Han, Yiwen Zhao, Chris Donahue, and Shinji Watanabe. VERSA: A versatile evaluation toolkit for speech, audio, and music. *arXiv preprint arXiv:2412.17667*, 2024.
    - Nikolai V. Smirnov. Table for estimating the goodness of fit of empirical distributions. *The Annals of Mathematical Statistics*, 19(2):279–281, 1948. doi: 10.1214/aoms/1177730256.
    - Cees H. Taal, Richard C. Hendriks, Richard Heusdens, and Jesper Jensen. An algorithm for intelligibility prediction of time–frequency weighted noisy speech. *IEEE Trans. Audio, Speech and Language Processing*, 19(7):2125–2136, 2011.
    - Lukas Tamm, Sebastian Möller, and Babak Naderi. Layer-wise analysis of xls-r representations for speech quality assessment. *arXiv preprint arXiv:2308.12077*, 2023.
    - Andrei N. Tikhonov and Vasiliy Y. Arsenin. *Solutions of Ill-posed Problems*. Winston & Sons, Washington, DC, 1977.
    - Gaurav Vaidya and Alexander J E Kell. Self-supervised models of audio effectively explain human cortical responses to speech. *bioRxiv*, 2022.
  - Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2nd edition, 2024.
  - Emmanuel Vincent, R.émi Gribonval, and Cédric Févotte. Performance measurement in blind audio source separation. *IEEE Trans. Audio, Speech and Language Processing*, 14(4):1462–1469, 2006.

Emmanuel Vincent, Tuomas Virtanen, and Sharon Gannot (eds.). *Audio Source Separation and Speech Enhancement*. Wiley, 2018.

Alon Vinnikov, Amir Ivry, Aviv Hurvitz, Igor Abramovski, Sharon Koubi, Ilya Gurvich, Shai Pe'er, Xiong Xiao, Benjamin Martinez Elizalde, Naoyuki Kanda, Xiaofei Wang, Shalev Shaer, Stav Yagev, Yossi Asher, Sunit Sivasankaran, Yifan Gong, Min Tang, Huaming Wang, and Eyal Krupka. The NOTSOFAR-1 challenge: New datasets, baseline, and tasks for distant meeting transcription. In *Proc. Interspeech*, pp. 5003–5007, Kos Island, Greece, 2024. doi: 10.21437/Interspeech. 2024-1788.

Alex Wilson and Bruno Fazenda. Characterisation of distortion profiles in relation to audio quality. In *Proc. DAFx*, pp. 1–8, 2014.

Neil Zeghidour, Nicolas Usunier, Gabriel Synnaeve, and Emmanuel Dupoux. Leaf: A learnable frontend for audio classification. In *Proceedings of ICLR*, 2021.

# A PERCEPTUAL DISTORTIONS APPLIED IN THE PS AND PM CALCULATIONS

Table 3: Distortions applied to the references when calculating the PS and PM measures ( $\S4.1$ ).  $f_s$  is the sampling frequency, and  $A_{95}$  and  $A_{RMS}$  mark the 95th-percentile and RMS absolute amplitudes.

Distortion	PS	PM			
Notch Filter	Center frequencies: 500, 1000, 2000, 4000, 8000 Hz	Number of notches: $\leq 20$ Operating band: $80 \text{ Hz} - 0.45 f_s$ Notch spacing: $\geq 300 \text{ Hz}$ Bandwidth: $\pm 60 \text{ Hz}$			
Comb Filter	Delay: 2.5-15 ms Feedback gain: 0.4-0.9	<b>Delay-gain pairs:</b> (2.5 ms, 0.4), (5 ms, 0.5), (7.5 ms, 0.6), (10 ms, 0.7), (12.5 ms, 0.9)			
Tremolo	Rate: 1, 2, 4, 6 Hz Depth: 0,3-1.0	Rate: 1, 2, 4, 6 Hz Depth:			
Additive Noise	SNR: -15, -10, -5, 0, 5, 10, 15 dB Noise color: white, pink, brown	<b>SNR:</b> -15, -10, -5, 0, 5, 10, 15 dB <b>Noise color:</b> white, pink, brown			
Additive Harmonic Tone	Tone frequency: 100, 500, 1000, 4000 Hz Amplitude: 0.02-0.08 (absolute)	Tone frequency: 100, 500, 1000, 4000 Hz Amplitude: $\{0.4, 0.6, 0.8, 1\} \times A_{\rm RMS}$			
Reverberation	RT <sub>60</sub> (Schroeder, 1965): 0.3-1.1 s Early tail length: 5, 10, 15, 20 ms	Exponential tail length: 50, 100, 200, 400 ms  Decay scaling: 0.3, 0.5, 0.7, 0.9			
Noise Gate	<b>Threshold:</b> 0.005, 0.01, 0.02, 0.04 (absolute)	Threshold: $\{0.05, 0.1, 0.2, 0.4\} \times A_{95}$			
Pitch Shift	<b>Offsets:</b> -4, -2, +2, +4 semitones	<b>Offsets:</b> -4, -2, +2, +4 semitones			
Low-Pass Filter	Cutoff: 2000, 3000, 4000, 6000 Hz	Cutoff rule: spectral-energy quintiles: 50, 70, 85, 95% Rounding: nearest 100 Hz			
High-Pass Filter	<b>Cutoff:</b> 100, 300, 500, 800 Hz	Cutoff rule: spectral-energy quintiles: 5, 15, 30, 50% Rounding: nearest 100 Hz			
Echo	<b>Delay:</b> 5-20 ms <b>Gain:</b> 0.3-0.7	<b>Delay:</b> 50, 100, 150 ms <b>Gain:</b> 0.4, 0.5, 0.7			
Hard Clipping	<b>Threshold:</b> 0.3, 0.5, 0.7 (absolute)	Threshold: $\{0.3, 0.5, 0.7\} \times A_{95}$			
Vibrato	Rate: 3, 5, 7 Hz Depth: 0.001-0.003 (fractional stretch)	Rate: 3, 5, 7 Hz Depth: adaptive, clipped to 0.01-0.05			

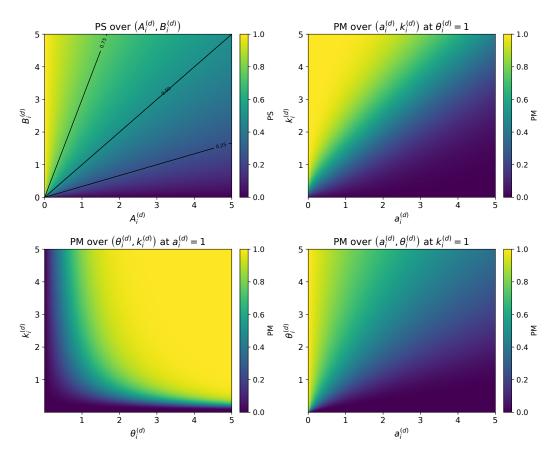


Figure 2: Functional behavior of the PS measure with 0.25, 0.5, 0.75 contour lines, and of the PM measure in three different setups of  $\theta_i^{(d)} = 1$ ,  $a_i^{(d)} = 1$ ,  $k_i^{(d)} = 1$ .

# B ADDITIONAL EXPREIMENTAL SETUP DETAILS

#### B.1 THE PS AND PM MEASURES

The functionality of the measures is demonstrated in Figure 2 and illustrates the behavior explained in  $\S 4.2$  and  $\S 4.3$ .

The empirical distributions of the frame-level values of the measures are shown in Figure 3. The PM and PS metrics exhibit contrasting distribution patterns. PM values cluster predominantly around zero with minimal density near one, while PS concentrate near one with virtually no occurrence near zero. Although frame-level human speech quality ratings are not publicly available for direct comparison, these patterns raise comparisons to how humans might perceive audio disturbances. The PM distribution aligns intuitively with human perception, as listeners typically penalize speech quality severely when disturbances occur, making ratings near the scale minimum unsurprising. However, real granular human ratings would likely show less extreme clustering around zero due to perceptual and rating scale complexities. The PS behavior presents a more complex interpretative challenge. Previous research suggests that humans perceive leakage as more quality-degrading than self-distortions, particularly in acoustic echo cancellation contexts (Khanagha et al., 2024), yet our findings here do not support this hypothesis. Whether this discrepancy stems from dataset characteristics, limitations of the PS measure itself, or the mismatch between granular PS values and aggregated human ratings remains unclear and warrants future investigation beyond the scope of this study.

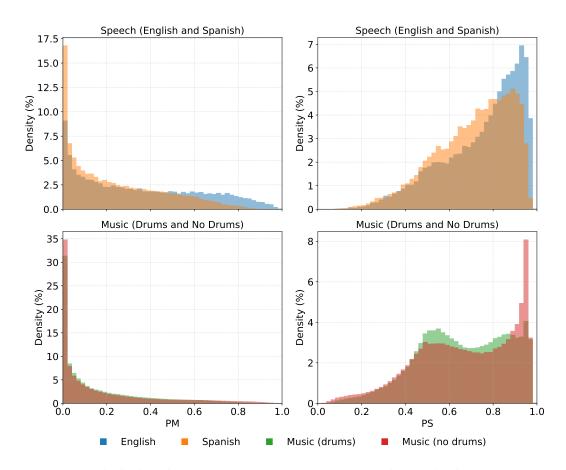


Figure 3: The distribution of PM and PS values across speech and music scenarios from the SEBASS database.

## B.2 THE SEBASS DATABASE

The SEBASS dataset suits this study for several reasons. Multilingual coverage of English and Spanish validates language-agnostic behavior, while music tests robustness to highly transient material. Large algorithmic spread creates rich output clusters that stress-test our methodology, and the dense sampling of raters allows for a more reliable estimation of the true mean-opinion score of subjective human opinion. Figure 4 shows that the speech reference signals have been recorded in a relatively clean environment with SNRs between 3.9 dB to 41.7 dB, with an average of 25 dB.

#### B.3 PRE-PROCESSING

We recognize that English and Spanish speakers rarely participate in the same conversation in real-life scenarios. To emulate realistic scenarios, we separate each 4-speaker mixture into their English and Spanish speakers, creating for each language two mixtures where the one has a pair of male speakers and the other a pair of female speakers. We acknowledge the uncertainty this step induces, as residuals of English may be present in the output signal of a Spanish speaker, and vice versa. It should be mentioned that listening tests have rendered this cross-language leakage extremely negligible. This may be since, as expected, source separation systems are able to leverage languages as a meaningful feature to recognize leakage and remove it.

Every waveform, including references, distortions, and outputs from all sources of the mixture, undergoes independent loudness normalization. We use the EBU Recommendation R-128 (European Broadcasting Union (EBU), 2011) and set the target level of each waveform to loudness units relative to full scale (LUFS) of -23. If the peak magnitude of the scaled waveform exceeds one, we attenuate it to avoid digital clipping. This step removes loudness bias, known to wrongly affect both human

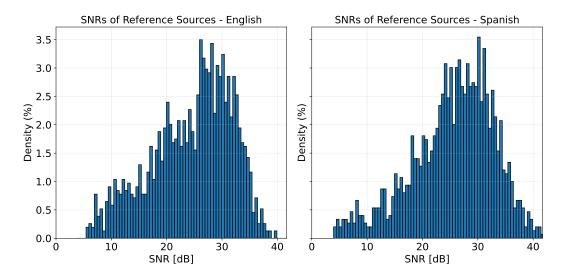


Figure 4: Frame-level SNR estimations for English and Spanish references in the SEBASS database.

and algorithmic quality judgments, while preserving inter-speaker level relations across the outputs. Since the PS and PM measures address source separation, we also filter out any frames in which there are not at least two active sources using energy-based thresholding.

When applying diffusion maps, we set  $\alpha=1$  in (3) to eliminate density-dependent bias from the embedding. This choice ensures that the PS and PM measures reflect the intrinsic geometric structure of the manifold rather than sampling density variations, which would introduce instead artificial distortions into the representation. We set t=1 in (6), to keep the diffusion operator focused on local neighborhoods and not being blurred by multi-step mixing.

#### B.4 Frame-level to Utterance-level Aggregation of the Measures

At trial l, let us denote the PM value of the i-th output of source-separation system q in time frame f as  $\mathrm{PM}_{i,f}^{q,l}$ . Let  $\mathcal{F}^l$  holds the time-frame indices with at least two active sources, and let  $\mathcal{F}_i^l \in \mathcal{F}^l$  be its subset of time-frame indices in which the i-th source is active. Then, the utterance-level PM measure after average aggregation is given by:

$$PM_{i,\text{utt}}^{q,l} = \frac{1}{|\mathcal{F}_i^l|} \sum_{f \in \mathcal{F}_i^l} PM_{i,f}^{q,l}.$$
 (45)

Although average aggregation assumes that human listeners perceive global audio quality by weighing local events equally, which is evidently not the case (Rix et al., 2001), we chose to carry it for the PM since its behavior already exhibits strong and frequent granular penalties where the score drops to around zero. Thus, it is assumed that standard human behavior that weighs negative experience heavily in the utterance-level score is implicitly carried out by the nature of the PM measure itself.

However, this is not the case for the PS measure. Here, the aggregation we applied is inspired by the window-based pooling and logistic mapping used inside PESQ (Rix et al., 2001). Again, considering only time frame indices in  $\mathcal{F}_i^l$  and dropping the rest, let us consider a window of size W frames that slides across the PS measure with a hop size of H frames. Using the p-norm, we define the following:

$$\ell_{i,m}^{q,l} = \left(\frac{1}{W} \sum_{w=1}^{W} \left| PS_{i,(m-1)H+w}^{q,l} \right|^{p} \right)^{1/p}, \tag{46}$$

where  $m \in \{1, \dots, M_i^l\}$  and  $M_i^l$  is the number of possible windows:

$$M_i^l = \max\left(1, \left|\frac{|\mathcal{F}_i^l| - W}{H}\right|\right). \tag{47}$$

 We then calculate the following root mean square expression:

$$\ell_i^{q,l} = \sqrt{\frac{1}{M_i^l} \sum_{m=1}^{M_i^l} \left(\ell_{i,m}^{q,l}\right)^2},\tag{48}$$

and eventually the aggregated PS measure is given by:

$$PS_{i,\text{utt}}^{q,l} = 0.999 + \frac{4}{1 + \exp(-1.3669 \,\ell_i^{q,l} + 3.8224)},\tag{49}$$

where the constants were chosen according to (ITU-T, 2007). Here, we penalize lower scores explicitly using the p-norm to better match human perceptual aggregation.

#### B.5 CORRELATION COEFFICIENTS BETWEEN AGGREGATED MEASURES AND MOS

At trial l, let the utterance-level MOS of the i-th output from separation system q be  $v_i^{q,l}$ . Given Q independent source separation systems such that  $q \in \{1, \ldots, Q\}$ , consider the Q-dimensional vectors:

$$\mathbf{PS}_{i,\text{utt}}^{l} = \left(\mathbf{PS}_{i,\text{utt}}^{1,l}, \dots, \mathbf{PS}_{i,\text{utt}}^{Q,l}\right)^{T},\tag{50}$$

$$\mathbf{PM}_{i,\text{utt}}^{l} = \left(\mathbf{PM}_{i,\text{utt}}^{1,l}, \dots, \mathbf{PM}_{i,\text{utt}}^{Q,l}\right)^{T},\tag{51}$$

$$\mathbf{v}_i^l = \left(v_i^{1,l}, \dots, v_i^{Q,l}\right)^T. \tag{52}$$

The PCC (Benesty et al., 2009) is measured twice, for the PS and the PM, as follows:

$$r_{i}^{\text{pcc},l}\left(\mathbf{PS}_{i,\text{utt}}^{l},\mathbf{v}_{i}^{l}\right) = \frac{\left(\overline{\mathbf{PS}}_{i,\text{utt}}^{l}\right)^{T}\overline{\mathbf{v}}_{i}^{l}}{\left\|\overline{\mathbf{PS}}_{i,\text{utt}}^{l}\right\|_{2}\left\|\overline{\mathbf{v}}_{i}^{l}\right\|_{2}},$$
(53)

$$r_{i}^{\text{pcc},l}\left(\mathbf{PM}_{i,\text{utt}}^{l},\mathbf{v}_{i}^{l}\right) = \frac{\left(\overline{\mathbf{PM}}_{i,\text{utt}}^{l}\right)^{T}\overline{\mathbf{v}}_{i}^{l}}{\left\|\overline{\mathbf{PM}}_{i,\text{utt}}^{l}\right\|_{2}\left\|\overline{\mathbf{v}}_{i}^{l}\right\|_{2}},$$
(54)

where  $\overline{\mathbf{PS}}_{i,\mathrm{utt}}^l$ ,  $\overline{\mathbf{PM}}_{i,\mathrm{utt}}^l$  and  $\overline{\mathbf{v}}_i^l$  are the centered versions of  $\mathbf{PS}_{i,\mathrm{utt}}^l$ ,  $\mathbf{PM}_{i,\mathrm{utt}}^l$  and  $\mathbf{v}_i^l$ , respectively.

Let  $\mathcal{R}: \mathbb{R}^Q \to R^Q$  be the ranking operator, which in the presence of ties assigns the average ranks. The SRCC (Sedgwick, 2014) is measured for the PS and the PM:

$$\rho_{i}^{\text{srcc},l}\left(\mathbf{PS}_{i,\text{utt}}^{l},\mathbf{v}_{i}^{l}\right) = r_{i}^{\text{pcc},l}\left(\mathcal{R}\left(\mathbf{PS}_{i,\text{utt}}^{l}\right),\mathcal{R}\left(\mathbf{v}_{i}^{l}\right)\right),\tag{55}$$

$$\rho_{i}^{\text{srcc},l}\left(\mathbf{PM}_{i,\text{utt}}^{l},\mathbf{v}_{i}^{l}\right)=r_{i}^{\text{pcc},l}\left(\mathcal{R}\left(\mathbf{PM}_{i,\text{utt}}^{l}\right),\mathcal{R}\left(\mathbf{v}_{i}^{l}\right)\right). \tag{56}$$

We report these correlation coefficients per English, Spanish, and music mixtures scenarios separately. Let us denote  $N_f^l$  the number of active sources in trial l during frame f. Then, given a scenario with  $\mathcal{L}$  independent trials such that  $l \in \{1, \dots, \mathcal{L}\}$ , we mark the maximal number of sources in trail l with  $N_{\text{max}}^l$ :

$$N_{\max}^l = \max_{f \in \mathcal{F}^l} N_f^l. \tag{57}$$

Then, for the PS and PM measures, the PCC and SRCC we report per scenario are given by:

$$PS^{pcc} = \frac{1}{\sum_{l=1}^{\mathcal{L}} N_{\text{max}}^{l}} \sum_{l=1}^{\mathcal{L}} \sum_{i=1}^{N_{\text{max}}^{l}} r_{i}^{pcc, l} \left( \mathbf{PS}_{i, \text{utt}}^{l}, \mathbf{v}_{i}^{l} \right), \tag{58}$$

$$PM^{pcc} = \frac{1}{\sum_{l=1}^{\mathcal{L}} N_{\text{max}}^{l}} \sum_{l=1}^{\mathcal{L}} \sum_{i=1}^{N_{\text{max}}^{l}} r_{i}^{\text{pcc},l} \left( \mathbf{PM}_{i,\text{utt}}^{l}, \mathbf{v}_{i}^{l} \right), \tag{59}$$

$$PS^{\text{srcc}} = \frac{1}{\sum_{l=1}^{\mathcal{L}} N_{\text{max}}^{l}} \sum_{l=1}^{\mathcal{L}} \sum_{i=1}^{N_{\text{max}}^{l}} \rho_{i}^{\text{srcc},l} \left( \mathbf{PS}_{i,\text{utt}}^{l}, \mathbf{v}_{i}^{l} \right), \tag{60}$$

$$PM^{\text{srcc}} = \frac{1}{\sum_{l=1}^{\mathcal{L}} N_{\text{max}}^{l}} \sum_{l=1}^{\mathcal{L}} \sum_{i=1}^{N_{\text{max}}^{l}} \rho_{i}^{\text{srcc},l} \left( \mathbf{PM}_{i,\text{utt}}^{l}, \mathbf{v}_{i}^{l} \right).$$
 (61)

# C ADDITIONAL EXPERIMENTAL RESULTS

Table 4: Self-supervised architectures, their pre-trained checkpoints, scenarios, and number of transformer layers.

Architecture	Checkpoint	Scenario	Transformer Layers
WavLM Large	microsoft/wavlm-large	English	24
WavLM Base	microsoft/wavlm-base	English	12
wav2vec 2.0 Large	facebook/wav2vec2-large-lv60	English	24
wav2vec 2.0 Base	facebook/wav2vec2-base	English	12
HuBERT Large	facebook/hubert-large-ll60k	English	24
HuBERT Base	facebook/hubert-base-ls960	English	12
wav2vec 2.0 Large	facebook/wav2vec2-large-xlsr-53	Spanish	24
MERT	m-a-p/MERT-v1-95M	Music	12

We begin by analyzing how performance depends on the choice of the pre-trained self-supervised model, the purpose of which is encoding waveforms into perceptual representations before they are fed into the diffusion maps. Table 4 lists the models we examine in this study. We consider six different models for English mixtures, based on the wav2vec 2.0 (Baevski et al., 2020), WavLM (Chen et al., 2022), and HuBERT (Hsu et al., 2021) backbones, with Figure 5 demonstrating their layer-wise performance. When using "Large" versions of the models, for both PCC and SRCC values, earlier layers frequently produce representations that allow superior results that gradually decline toward deeper layers, showing approximately 10% average absolute degradation between extremes. Existing layer-wise analysis already reported that acoustic and phonetic content is richly represented in intermediate layers, while deeper layers shift toward semantic abstraction (Pasad et al., 2022; Vaidya & Kell, 2022). Additional work confirms that distortion sensitivity peaks in the lower or middle layers and diminishes in deeper ones (Tamm et al., 2023; Hung et al., 2022), and that pretrained models tend to lose low-level signal fidelity in their deepest layers (Moussa & Toneva, 2025). A notable data point appears in the final layers of wav2vec 2.0 with a sharp drop in performance, especially for PS. This is likely due to its contrastive learning pretraining objective, which drives later layers to specialize in predicting quantized latent codes rather than preserving acoustic detail. For the "Base" versions of the models, we observe a somewhat different behavior. At low and middle layers, their performance is often quite competitive with the "Large" variants, and in several cases the former even outperforms the latter in deeper layers. However, for WavLM, the gap widens toward the final layers, with the "Large" version consistently outperforming. Interestingly, wav2vec 2.0 Base does not exhibit the sharp degradation observed in its counterpart and instead its deeper layers remain stable and even show improvements for PS, suggesting that the absence of over-specialization to quantized prediction in the "Base" model preserves sensitivity to perceptual distortions.

Table 5 narrows these models down to their top performing layer, chosen by the max-min criteria of the PCC and SRCC values, across all layers. A no-encoding option is also reported, where

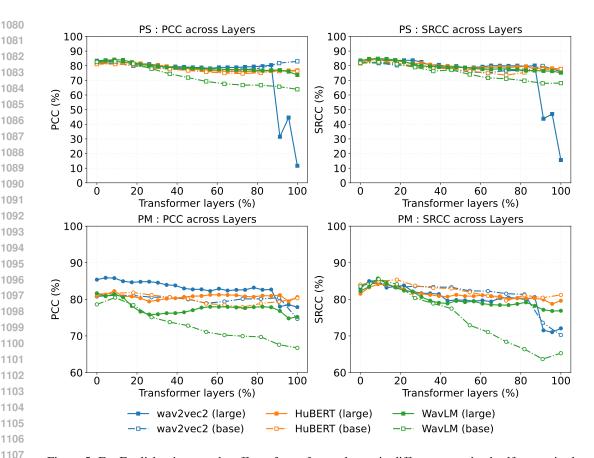


Figure 5: For English mixtures, the effect of transformer layers in different pretrained self-supervised models on the PCC and SRCC values for the PS and PM measures.

waveforms are skip-connected directly into the diffusion maps, which under-performs compared to encoded modes and emphasizes the effectiveness of the waveform encoding in the proposed pipeline. These results reaffirm that shallow layers achieve optimal performance. Although point-by-point comparisons show that "Base" models perform comparably to, or occasionally slightly exceed "Large" models, applying the max-min criteria across the models in the table reveals that 'Large" models are preferable when jointly optimizing for PS and PM. For once, wav2vec 2.0 Large achieves for the PM a PCC and SRCC differences from its "Base" counterpart of absolute 6% and 2%, respectively, even when the PS case shows a negligible gap. Among the "Large" model variants, wav2vec 2.0 Large with transformer layer 2 emerges as the ideal configuration and we carry it forward as a case study we investigate. It should be noted that among "Large" models, the PS very slightly changes with roughly 1% and 0.5% gaps between extremes for the PCC and SRCC, respectively, while the PM gaps are more meaningful. This suggests that the choice of model may mainly affect the PM scores.

Next, we analyze all scenarios with wav2vec 2.0 "Large" encoders for speech and the MERT encoder for music representations, and analyze the effect of their transformer layer on performance. The results are shown in Figure 6. English and Spanish mixtures, both evaluated with wav2vec 2.0 backbones, show broadly similar trends across layers, with Spanish exhibiting a sharper decline in deeper layers. This can be explained by the XLSR pretraining data being relatively scarcer in Spanish than in English, leading later layers to emphasize cross-lingual abstractions over fine acoustic detail (Conneau et al., 2020). Music mixtures with drums show the lowest performance among scenarios, which we attribute to the dominance of strong percussive transients. Self-supervised models have demonstrated less stability in these highly non-stationary regions, reducing the ability of PS and PM to capture perceptual degradations (Zeghidour et al., 2021). In contrast, music mixtures without drums demonstrate consistently high performance, in most layers even surpassing speech mixtures. This likely stems from the MERT backbone being particularly suited in capturing harmonic and timbral structure, allowing the measures to remain faithful to perceptual cues such as instrument

Measure	Representation	Transformer Layer	SRCC	PCC
PS	wav2vec 2.0 (Large)	2	84.12%	83.74%
PS	wav2vec 2.0 (Base)	2	84.25%	83.23%
PM	wav2vec 2.0 (Large)	2	84.69%	86.36%
PM	wav2vec 2.0 (Base)	2	82.79%	80.07%
PS	WavLM (Large)	3	84.80%	84.16%
PS	WavLM (Base)	2	84.84%	84.19%
PM	WavLM (Large)	3	<b>85.71</b> %	81.44%
PM	WavLM (Base)	2	82.82%	77.51%
PS	HuBERT (Large)	3	84.48%	83.09%
PS	HuBERT (Base)	2	84.83%	82.73%
PM	HuBERT (Large)	3	84.12%	82.24%
PM	HuBERT (Base)	2	81.37%	79.47%
PS	Waveform (raw)	-	73.42%	71.04%
PM	Waveform (raw)	-	69.30%	66.62%

Table 5: For English mixtures, comparing PCC and SRCC values between best-layer performance of "Large" and "Base" models. A raw waveform option, i.e. no encoding, is also reported. The highest SRCC and PCC are in bold per PS and PM.

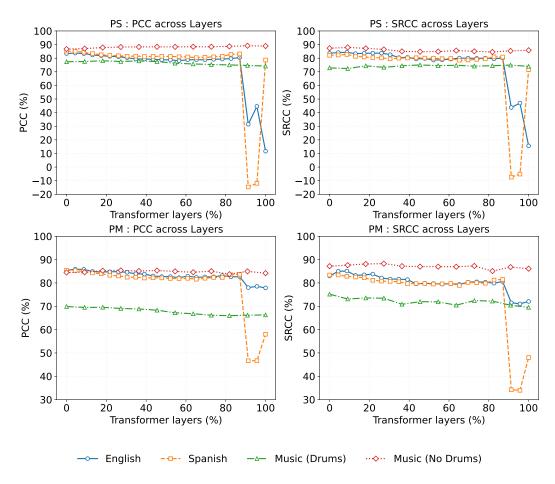


Figure 6: For all scenarios, the effect of transformer layers in their respective pretrained self-supervised architectures on the PCC and SRCC values for the PS and PM measures.

texture and vocal clarity (Li et al., 2023). Interestingly, whether drums are present or not, MERT-based performance demonstrates a steady behavior across all layers, suggesting the MERT representations are not vulnerable to degradation across processing stages. The max-min criteria across all layers,

per scenario, shows that the ideal layers for English, Spanish, drums, and no-drums music mixtures are layers 2, 2, 1, and 3, respectively.

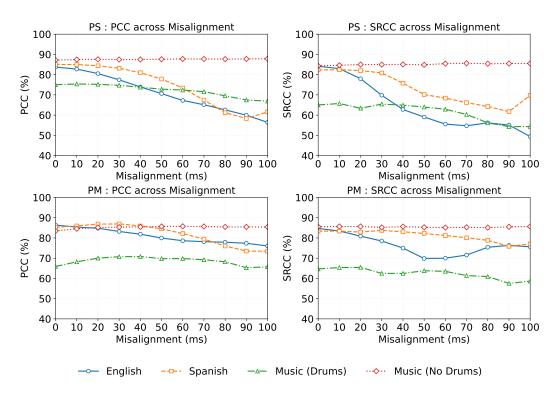


Figure 7: The effect of temporal misalignment between references and outputs of the separation system on the PS and PM measures.

We employed these layers to construct Table 1 in the main text and now we extend the discussion on it. The advantage of PESQ can be attributed to its long-standing perceptual model, which explicitly encodes aspects of loudness perception, asymmetry, and time-alignment penalties, features that directly penalize separation artifacts. In Spanish mixtures, the PS and PM are most performant in terms of PCC, but fall behind PESO and SDR-based metrics in SRCC. One possible explanation is that the syllable-timed rhythm and steady vowels of Spanish make fidelity-driven metrics such as SI-SDR, CI-SDR, and SDR more predictive of listener rankings, as these metrics emphasize reconstruction accuracy at the waveform level. For music mixtures, PS and PM achieve the strongest overall correlations across both drums and no-drums conditions for both PCC and SRCC. Even though SpeechBERTscore has shown impressive results and is also based on a self-supervised backbone, it is mostly not competitive with our measures, and notably even performing worse than our raw waveform version at times, which projects on the importance of the diffusion maps in the pipeline. We emphasize that unlike English, we only inspected one backbone model for Spanish or music mixtures. In addition, the aggregation strategies we applied were not data-driven but a heuristic and reasoning-based choice. Consequently, while the proposed measures already demonstrate strong alignment with human perception, these low-hanging fruits may potentially boost performance. Quite surprisingly, the first group of STOI, PESQ, and SDR-based measures is consistently preferable to the second group consisting of DNSMOS, speechBERTscore, UTMOS, and others, which rarely achieve more than 70% in performance. One crucial conclusion this table suggests is that measures originally developed for a certain audio application, should not be zero-shot adapted into other applications, and in that case into source separation evaluation. Otherwise, values that drift from human opinion may be reported, which may spiral the development of audio technologies instead of accelerating it.

An additional stress test for our measures concerns their robustness to temporal misalignment between the input and output streams of the separator, a phenomenon commonly introduced by modern communication systems or, e.g., when dealing with references obtained from different, per-speaker microphones, such as in meeting datasets (Carletta et al., 2005; Vinnikov et al., 2024).

 Figure 7 illustrates the effect of artificial delays ranging from 0 ms to 100 ms across English, Spanish, and music scenarios. While performance gradually degrades for speech scenarios as misalignment grows, as expected, a 20 ms delay or less still preserves coefficients higher than 80%. Surpassing this threshold, however, often causes a pronounced drop that underscores this weakness in our measures, since human ratings are insensitive to these short latencies. Music mixtures exhibit a different pattern, as performance remains largely stable across delays, with the presence of drums introducing more variability than its counterpart.

Table 6 provides complementary information to the NMI results in Figure 1 but listing the frame counts used for the PS and PM measures for every examined threshold. Even the lowest threshold of 0.1, had a minimum of 481 for calculations, rendering its results statistically reliable. An interesting observation in the music scenario, without drums, is that it exhibits significantly more time frames in which there are at least two active sources, compared to all other scenarios.

Threshold	English		Spa	Spanish   Music		(Drums)	Music (No Drums)	
	PS≤th	PM≤th	PS≤th	PM≤th	PS≤th	PM≤th	PS≤th	PM≤th
0.1	583	7426	622	10721	481	13987	622	22859
0.2	1546	12086	1591	15756	1536	16126	1832	28828
0.3	3191	16350	3350	19714	3327	17846	4226	33231
0.4	5753	20118	6091	23054	5861	19492	8821	36953
0.5	9115	23697	9725	25964	8927	21033	15748	40426
0.6	13364	27232	13904	28627	12037	22522	24958	43769
0.7	18465	30758	18592	30885	15373	23864	35434	47186
0.8	24477	34076	23871	32703	19498	25168	46078	50682
0.9	31572	36507	29589	33902	25748	26614	57795	54939
1.0	37888	37888	34496	34496	34688	34688	66528	66528

Table 6: Frame counts used for NMI computation at each threshold, denoted 'th' in the table. Columns show counts of frames per scenario, split by PS and PM subsets.

In the next phase, we investigate the deterministic error radius and probabilistic CIs derived for the PS and PM measures in Appendix E. Figure 8 shows histograms of the frame-level error distributions for speech and music mixtures. As expected, the radius caused by the spectral truncation in the diffusion maps process is typically an order of magnitude smaller than the 95% probabilistic width, which is originated from finite-sample clusters on the manifold. The error radius is also concentrated mostly near zero, which further confirms its negligibility. CIs typically span 10-50% of the dynamic range of the measures at the frame level, but surprisingly in the PM, between 10-15% of frame-level instances have probabilistic tails that approach zero across scenarios. The immediate contribution of these results are by making development of source separation systems more reliable and informed at the frame-level.

For illustration, Figure 9 shows reference and output spectrograms from an English mixture, timealigned to corresponding PM and PS values over a 10-second utterance using the "Large" models, with layers specified in Table 5. While a single example cannot be over-interpreted, the latest observation about the PS gaps across layers is visually supported here, with very similar behavior of all models. The PM shows highly correlated behavior, but with noticeable different values by wav2vec 2.0, which exhibits the highest PM value for PCC. An interesting visual example is shown at approximately the 9 seconds mark, when both speakers exhibit visible self-distortion artifacts accompanied by sharp drops in their PM measures. Listening tests confirmed that leakage is indeed more present in "Speaker 2" than in "Speaker 1", as supported by the PS plot.

Finally, to give the reader an intuitive grasp of how the two error terms evolve in a time-aligned manner with the PS and PM measures, Figure 10 illustrates a representative example.

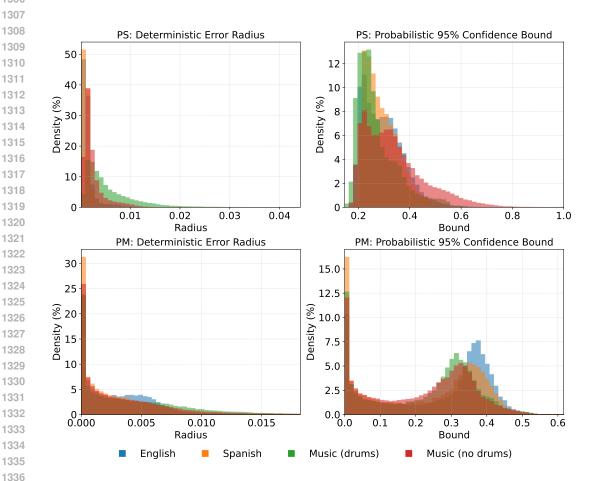


Figure 8: An histogram view of the frame-level deterministic error radius and the 95% probabilistic tail in the PS and PM measures across scenarios.

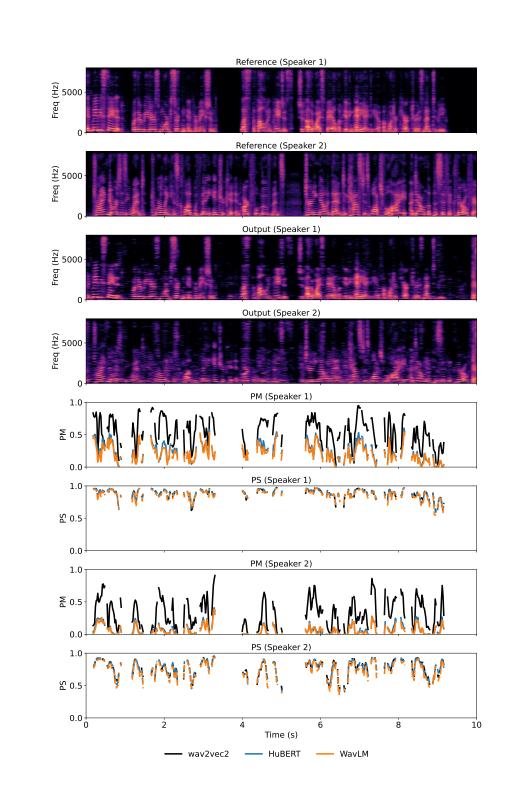


Figure 9: For an English mixture with two speakers, a spectral view of the system signals and aligned with a time-series view of the PS and PM measures of each speaker across different self-supervised architectures. Blank time intervals remain whenever speech does not overlap.

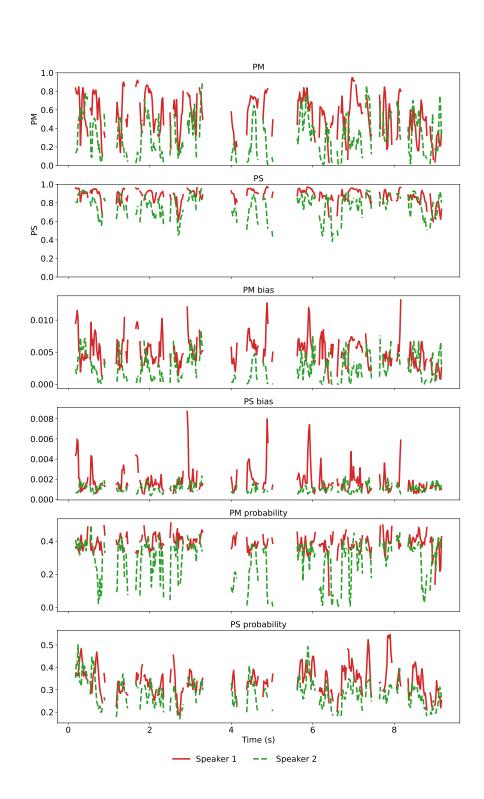


Figure 10: Time-aligned view of the PM and PS measures and their deterministic error radius and probabilistic tail with 95%, of two English speakers. Time indices where speech does not overlap remain blank.

# D EXPECTATION AND PROBABILISTIC CONFIDENCE BOUND OF THE TRUNCATION ERROR

Truncating the spectrum to d dimensions breaks the equality in Equation (9), and leads to a truncation error. Here, we derive the expectation and probabilistic tail bound for this truncation error. Assume a point  $\mathbf{x}_i \in \mathcal{X}$  is drawn from the stationary distribution  $\pi$  7 of the diffusion process, where  $i \in \{1, \dots, N\}$ . This assumption is supported by (Hein et al., 2005, Lem. 1), and by showing empirically on 5,000 graphs that the corresponding eigenvector matches the theoretical stationary distribution up to statistical fluctuations. Given that the N-1-dimensional embedding of  $\mathbf{x}_i$  is truncated to dimension d, then the truncation error is expressed as (10):

$$E(\mathbf{x}_i) = \left(\sum_{\ell=d+1}^{N-1} \lambda_{\ell}^{2t} \mathbf{u}_{\ell}^2(i)\right)^{1/2}.$$
 (62)

We define the squared truncation error and analyze it:

$$T(\mathbf{x}_i) = E^2(\mathbf{x}_i) = \sum_{\ell=d+1}^{N-1} \lambda_{\ell}^{2t} \mathbf{u}_{\ell}^2(i).$$
 (63)

Since the eigenvectors  $\{\mathbf{u}_{\ell}\}_{\ell=0}^{N-1}$  are orthonormal under  $\pi$ , then:

$$\mathbb{E}_{\boldsymbol{\pi}}\left[\mathbf{u}_{\ell}^{2}(i)\right] = \sum_{i=1}^{N} \boldsymbol{\pi}_{i} \mathbf{u}_{\ell}^{2}(i) = 1, \tag{64}$$

from which we derive the expectation of  $T(\mathbf{x}_i)$  under  $\pi$ :

$$\mathbb{E}_{\pi} [T(\mathbf{x}_i)] = \mathbb{E}_{\pi} \left( \sum_{\ell=d+1}^{N-1} \lambda_{\ell}^{2t} u_{\ell}^{2}(i) \right) = \sum_{\ell=d+1}^{N-1} \lambda_{\ell}^{2t}.$$
 (65)

Thus, the expectation of the truncation error is given directly by:

$$\mathbb{E}_{\boldsymbol{\pi}}\left[E(\mathbf{x}_i)\right] = \left(\sum_{\ell=d+1}^{N-1} \lambda_{\ell}^{2t}\right)^{1/2}.$$
(66)

This term decays monotonically as d grows and is typically lower than  $10^{-3}$ . To obtain a non-asymptotic and high-probability confidence bound on the truncation error, we derive  $\forall \ell, i$  (64):

$$\left|\mathbf{u}_{\ell}(i)\right| \le \pi_{\min}^{-1/2}, \quad \pi_{\min} = \min_{i \in \{1, \dots, N\}} \boldsymbol{\pi}.$$
 (67)

Any bounded variable is sub-Gaussian, and its  $\psi_2$ -norm is at most the bound divided by  $\sqrt{\ln 2}$  (Vershynin, 2024, Example 2.6.5):

$$\|\mathbf{u}_{\ell}(i)\|_{\psi_{2},\pi} \le \frac{\pi_{\min}^{-1/2}}{\sqrt{\ln 2}} := K.$$
 (68)

Let m = N - 1 - d, so we define  $\mathbf{z}_i \in \mathbb{R}^m$  as:

$$\mathbf{z}_{i} = \left(\mathbf{u}_{d+1}(i), \dots, \mathbf{u}_{N-1}(i)\right)^{T},\tag{69}$$

and the diagonal matrix of weights  $\mathbf{D} \in \mathbb{R}^{m \times m}$  as:

$$\mathbf{D} = \operatorname{diag}(\lambda_{d+1}^t, \dots, \lambda_{N-1}^t). \tag{70}$$

Then  $E(\mathbf{x}_i)$  and  $T(\mathbf{x}_i)$  can be rewritten as:

$$T(\mathbf{x}_i) = \left\| \mathbf{D} \mathbf{z}_i \right\|_2^2,\tag{71}$$

$$E(\mathbf{x}_i) = \left\| \mathbf{D} \mathbf{z}_i \right\|_2. \tag{72}$$

For  $\ell > 0$ ,  $\mathbf{u}_{\ell}(i)$  is zero-mean under  $\pi$ . Consequently, the vector  $\mathbf{z}_{i}$  is zero-mean and by definition satisfies  $\|\mathbf{z}_{i}\|_{\psi_{2},\pi} \leq K\sqrt{m}$ . We also notice that multiplication by a fixed matrix scales the sub-Gaussian norm linearly, and since  $\mathbf{D}$  is symmetric and positive:

$$\|\mathbf{D}\mathbf{z}(i)\|_{\psi_2,\pi} \le K\sqrt{m}\|\mathbf{D}\|_2 = K\sqrt{m}\max_{\ell>d}\lambda_{\ell}^t = K\sqrt{m}\lambda_{d+1}^t.$$
 (73)

According to (Vershynin, 2024, Prop. 6.2.1), for an m-dimensional, zero-mean and sub-Gaussian vector  $\mathbf{Y}$  with  $\|\mathbf{Y}\|_{\psi_2,\pi} \leq \kappa$ , it holds:

$$\mathbb{P}_{\pi}\{\|\mathbf{Y}\|_{2} \ge C\kappa(\sqrt{m} + t)\} \le e^{-t^{2}},\tag{74}$$

where  $t \ge 0$  and C > 0 is a constant. Setting  $\mathbf{Y} = \mathbf{D}\mathbf{z}_i$  and  $\kappa = K\sqrt{m}\lambda_{d+1}^t$  gives:

$$\mathbb{P}_{\pi}\left\{T(\mathbf{x}_i) > C^2 \lambda_{d+1}^{2t} K^2 m \left(\sqrt{m} + t\right)^2\right\} \le e^{-t^2}.$$
 (75)

Let  $\delta \in (0,1)$  and set  $t = \sqrt{\ln(1/\delta)}$ . We can rewrite (75) as:

$$\mathbb{P}_{\pi} \left\{ T(\mathbf{x}_i) \le C^2 \lambda_{d+1}^{2t} K^2 m \left( \sqrt{m} + \sqrt{\ln \frac{1}{\delta}} \right)^2 \right\} \ge 1 - \delta. \tag{76}$$

Thus, the desired confidence bound on the truncation error is:

$$\mathbb{P}_{\pi} \left\{ E(\mathbf{x}_i) \le C \lambda_{d+1}^t K \left( m + \sqrt{m \ln \frac{1}{\delta}} \right) \right\} \ge 1 - \delta. \tag{77}$$

The choice of d dimensions affects both m that shrinks linearly with d and  $\lambda_{d+1}^t$  that falls monotonically with d. K is affected by the minimal stationary probability  $\pi_{\min}$ , so if the graph contains rare points then  $\pi_{\min}$  may be tiny, while a well-balanced graph derives  $K \sim \sqrt{N}$  and tightens the bound.

# E DETERMINISTIC ERROR RADIUS AND PROBABILISTIC TAIL BOUND OF THE MEASURES

We derive a deterministic error radius and a high-probability confidence bound on the frame-level PS and PM measures by considering: (i) spectral truncation error due to retaining d diffusion coordinates, which is separately developed in Appendix D; (ii) finite-sample uncertainty in estimating the cluster centroid and covariance. We then combine these via union bounds. In this section, we consider a fixed trail l, separation system q, and time frame f.

#### E.1 THE PS MEASURE

Considering source indices  $i,j \in \{1,\ldots,N_f\}$  (§2), we begin by analyzing the effect of the truncation error, assuming access to cluster statistics. The difference between the embedding of  $\hat{\mathbf{x}}_i$  and the centroid of cluster j can be expressed in the truncated subspace  $\mathbb{R}^d$  and in its complement subspace  $\mathbb{R}^m$ , respectively denoted  $\mathbf{\Delta}_{i,j}^{(d)}$  and  $\mathbf{\Delta}_{i,j}^{(m)}$ . Using (10), (15):

$$\boldsymbol{\Delta}_{i,j}^{(d)} = \boldsymbol{\Psi}_t^{(d)}(\hat{\mathbf{x}}_i) - \boldsymbol{\mu}_j^{(d)} \in \mathbb{R}^d, \tag{78}$$

$$\Delta_{i,j}^{(m)} = \Psi_t^{(m)}(\hat{\mathbf{x}}_i) - \mu_j^{(m)} \in \mathbb{R}^m,$$
 (79)

where m = N - d - 1. For completion, for every  $\mathbf{x} \in \mathcal{X}$  and its global index  $k \in \{1, \dots, N\}$ :

$$\boldsymbol{\mu}_{j}^{(m)} = \frac{1}{\left|\mathcal{C}_{j}^{(m)}\right|} \sum_{\boldsymbol{\psi} \in \mathcal{C}_{j}^{(m)}} \boldsymbol{\psi} \tag{80}$$

$$C_j^{(m)} = \left\{ \mathbf{\Psi}_t^{(m)}(\mathbf{x}_j), \mathbf{\Psi}_t^{(m)}(\mathbf{x}_{j,p}) \mid p = 1, \dots, N_p \right\}$$
(81)

$$\Psi_t^{(m)}(\mathbf{x}) = (\lambda_{d+1}^t \mathbf{u}_{d+1}(k), \dots, \lambda_{N-1}^t \mathbf{u}_{N-1}(k)).$$
(82)

In the full, N-1-dimensional space, the cluster  $C_j$  is given by:

$$C_j = \{ \mathbf{\Psi}_t(\mathbf{x}_j), \mathbf{\Psi}_t(\mathbf{x}_{j,p}) \mid p = 1, \dots, N_p \},$$
(83)

with mean  $\mu \in \mathbb{R}^{N-1}$ , difference  $\Delta_{i,j} \in \mathbb{R}^{N-1}$  and covariance  $\Sigma_j \in \mathbb{R}^{(N-1)\times(N-1)}$  that hold:

$$\boldsymbol{\mu}_{j} = \begin{bmatrix} \boldsymbol{\mu}_{j}^{(d)} \\ \boldsymbol{\mu}_{j}^{(m)} \end{bmatrix}, \quad \boldsymbol{\Delta}_{i,j} = \begin{bmatrix} \boldsymbol{\Delta}_{i,j}^{(d)} \\ \boldsymbol{\Delta}_{i,j}^{(m)} \end{bmatrix}, \quad \boldsymbol{\Sigma}_{j} = \begin{bmatrix} \boldsymbol{\Sigma}_{j}^{(d)} & \boldsymbol{C}_{j} \\ \boldsymbol{C}_{j}^{T} & \boldsymbol{\Sigma}_{j}^{(m)} \end{bmatrix}, \tag{84}$$

where  $\Sigma_i^{(m)} \in \mathbb{R}^{m \times m}$  and  $C_j \in \mathbb{R}^{d \times m}$  are:

$$\Sigma_{j}^{(m)} = \frac{1}{\left|\mathcal{C}_{j}^{(m)}\right| - 1} \sum_{\boldsymbol{\psi} \in \mathcal{C}_{j}^{(m)}} \left(\boldsymbol{\psi} - \boldsymbol{\mu}_{j}^{(m)}\right) \left(\boldsymbol{\psi} - \boldsymbol{\mu}_{j}^{(m)}\right)^{T}, \tag{85}$$

$$C_{j} = \frac{1}{\left|\mathcal{C}_{j}^{(m)}\right| - 1} \sum_{p=0}^{N_{p}} \left(\boldsymbol{\Psi}_{t}^{(d)}(\mathbf{x}_{j,p}) - \boldsymbol{\mu}_{j}^{(d)}\right) \left(\boldsymbol{\Psi}_{t}^{(m)}(\mathbf{x}_{j,p}) - \boldsymbol{\mu}_{j}^{(m)}\right)^{T}.$$
 (86)

According to 16, the squared Mahalanobis distance from  $\Psi_t(\hat{\mathbf{x}}_i)$  to  $C_j$  is:

$$d_M^2\left(\mathbf{\Psi}_t(\hat{\mathbf{x}}_i); \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j\right) = \boldsymbol{\Delta}_{i,j}^T \left(\boldsymbol{\Sigma}_j + \epsilon I^{(N-1)}\right)^{-1} \boldsymbol{\Delta}_{i,j}, \tag{87}$$

where inversion is empirically obtained by taking  $\epsilon = 10^{-6}$  with  $I^{(N-1)}$  being the N-1-dimensional identity matrix. To evaluate the truncation effect, we perform blockwise inversion on (87) via the Schur complement (Horn & Johnson, 2013):

$$d_M^2\left(\boldsymbol{\Psi}_t(\hat{\mathbf{x}}_i); \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j\right) = \left(\boldsymbol{\Delta}_{i,j}^{(d)}\right)^T \left(\boldsymbol{\Sigma}_j^{(d)} + \epsilon I^{(d)}\right)^{-1} \boldsymbol{\Delta}_{i,j}^{(d)} + \boldsymbol{r}_{i,j}^T \boldsymbol{S}_j^{-1} \boldsymbol{r}_{i,j}, \tag{88}$$

where  $r_{i,j} \in \mathbb{R}^m$  and the Schur complement  $S_j \in \mathbb{R}^{m \times m}$  hold:

$$r_{i,j} = \Delta_{i,j}^{(m)} - C_j^T \left( \Sigma_j^{(d)} + \epsilon I^{(d)} \right)^{-1} \Delta_{i,j}^{(d)},$$
 (89)

$$S_j = \Sigma_j^{(m)} - C_j^T \left( \Sigma_j^{(d)} + \epsilon I^{(d)} \right)^{-1} C_j.$$

$$(90)$$

We now utilize the inequality:

$$\forall a, b \ge 0: \quad \left| \sqrt{a+b} - \sqrt{a} \right| \le \sqrt{b}, \tag{91}$$

obtained by the mean-value theorem for  $f(\cdot) = \sqrt{\cdot}$  (Rudin, 1976, Ch. 5). Let us set:

$$a = \left(\boldsymbol{\Delta}_{i,j}^{(d)}\right)^T \left(\boldsymbol{\Sigma}_j^{(d)} + \epsilon I^{(d)}\right)^{-1} \boldsymbol{\Delta}_{i,j}^{(d)}, \tag{92}$$

$$b = \boldsymbol{r}_{i,j}^T \boldsymbol{S}_j^{-1} \boldsymbol{r}_{i,j}, \tag{93}$$

to obtain:

$$|\delta_{i,j}| = \left| d_M\left(\mathbf{\Psi}_t(\hat{\mathbf{x}}_i); \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j\right) - d_M\left(\mathbf{\Psi}_t^{(d)}(\hat{\mathbf{x}}_i); \boldsymbol{\mu}_j^{(d)}, \boldsymbol{\Sigma}_j^{(d)}\right) \right| \le \sqrt{\boldsymbol{r}_{i,j}^T \boldsymbol{S}_j^{-1} \boldsymbol{r}_{i,j}}.$$
 (94)

Namely,  $|\delta_{i,j}|$  is the truncation error of this Mahalanobis distance. From (17), it holds that:

$$\delta_{i,i} = A_i - A_i^{(d)}, \quad \delta_{i,j^*} = B_i - B_i^{(d)},$$
(95)

where  $j^*$  is defined in (17) and:

$$A_i = d_M \left( \mathbf{\Psi}_t(\hat{\mathbf{x}}_i); \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i \right), \tag{96}$$

$$B_i = d_M \left( \mathbf{\Psi}_t(\hat{\mathbf{x}}_i); \boldsymbol{\mu}_{i^*}, \boldsymbol{\Sigma}_{i^*} \right). \tag{97}$$

Consider the N-1-dimensional representation of  $PS_i^{(d)}$ , i.e.  $PS_i$  (18):

$$PS_i = 1 - \frac{A_i}{A_i + B_i}, (98)$$

which is smooth and differentiable in  $A_i$  and  $B_i$ , since by definition  $A_i + B_i > 0$ . We assume and have empirically validated that truncation introduces a small relative change, i.e.:

$$|\delta_{i,i}| \ll A_i^{(d)} + B_i^{(d)},$$
 (99)

$$|\delta_{i,j^*}| \ll A_i^{(d)} + B_i^{(d)},$$
 (100)

making the first-order Taylor expansion of  $\mathrm{PS}_i$  around  $\left(A_i^{(d)}, B_i^{(d)}\right)$  valid. We can therefore write:

$$PS_{i} - PS_{i}^{(d)} \simeq \frac{\partial PS_{i}}{\partial A_{i}} \left( A_{i}^{(d)}, B_{i}^{(d)} \right) \delta_{i,i} + \frac{\partial PS_{i}}{\partial B_{i}} \left( A_{i}^{(d)}, B_{i}^{(d)} \right) \delta_{i,j^{*}}$$

$$= -\frac{B_{i}^{(d)}}{\left( A_{i}^{(d)} + B_{i}^{(d)} \right)^{2}} \delta_{i,i} + \frac{A_{i}^{(d)}}{\left( A_{i}^{(d)} + B_{i}^{(d)} \right)^{2}} \delta_{i,j^{*}},$$
(101)

where the quadratic remainder in the expansion is empirically one order smaller than the first-order term and can be safely dropped. Applying the triangle inequality and (94) yields the deterministic error radius in the PS measure:

$$\left| \text{PS}_{i} - \text{PS}_{i}^{(d)} \right| \leq \frac{B_{i}^{(d)} |\delta_{i,i}| + A_{i}^{(d)} |\delta_{i,j^{*}}|}{\left(A_{i}^{(d)} + B_{i}^{(d)}\right)^{2}} = \frac{B_{i}^{(d)} \sqrt{\boldsymbol{r}_{i,i}^{T} \boldsymbol{S}_{i}^{-1} \boldsymbol{r}_{i,i}} + A_{i}^{(d)} \sqrt{\boldsymbol{r}_{i,j^{*}}^{T} \boldsymbol{S}_{j^{*}}^{-1} \boldsymbol{r}_{i,j^{*}}}}{\left(A_{i}^{(d)} + B_{i}^{(d)}\right)^{2}}. \quad (102)$$

We now quantify the uncertainty in  $\widehat{\mathrm{PS}}_i^{(d)}$  due to finite-sample cluster statistics. Empirically, we observe that cluster coordinates exhibit weak dependence between one another and derive from (Bartlett, 1946) the following cut-off rule for the effective sample size of cluster  $\mathcal{C}_i^{(d)}$ :

$$n_{j,\text{eff}} = \frac{n_j}{1 + 2\sum_{\ell=1}^{L_j} \hat{\rho}_{j,\ell}}, \quad n_j = \left| \mathcal{C}_j^{(d)} \right|$$
 (103)

where  $\hat{\rho}_{j,\ell}$  is the empirical average Pearson auto-correlation of coordinates at lag  $\ell$ , and:

$$L_j = \arg\min_{\ell} \left\{ |\hat{\rho}_{j,\ell}| < \frac{z_{0.975}}{\sqrt{n_j - \ell}} \right\}.$$
 (104)

Empirical evidence across 5,000 graphs suggest that on average  $\sum_{\ell} \hat{\rho}_{j,\ell} \simeq 0.2$ , and so we set  $n_{j,\text{eff}} = 0.7 \, n_j$  for all clusters.

To bound the deviation between the estimated and true cluster mean and covariance, we employ the vector and matrix Bernstein (Vershynin, 2024, Props. 2.8.1, 4.7.1) and the dependent Hanson-Wright inequalities (Adamczak, 2015, Thm. 2.5). For every  $\delta_{j,\mu}^{PS}$ ,  $\delta_{j,\Sigma}^{PS} \in (0,1/2)$ , with respective least probabilities  $1 - \delta_{j,\mu}^{PS}$  and  $1 - \delta_{j,\Sigma}^{PS}$ :

$$\left| \boldsymbol{\mu}_{j} - \hat{\boldsymbol{\mu}}_{j} \right| \leq \sqrt{\frac{2\lambda_{\max}\left(\widehat{\boldsymbol{\Sigma}}_{j}^{(d)}\right)\ln\left(2/\delta_{j,\boldsymbol{\mu}}^{\mathrm{PS}}\right)}{n_{j,\mathrm{eff}}}} := \Delta_{j,\boldsymbol{\mu}}$$
 (105)

$$\left\| \mathbf{\Sigma}_{j}^{(d)} - \widehat{\mathbf{\Sigma}}_{j}^{(d)} \right\|_{2} \le C \lambda_{\max} \left( \widehat{\mathbf{\Sigma}}_{j}^{(d)} \right) \left( \frac{r_{j}}{n_{j,\text{eff}}} + \frac{r_{j} + \ln\left(2/\delta_{j,\mathbf{\Sigma}}^{\text{PS}}\right)}{n_{j,\text{eff}}} \right) := \Delta_{j,\mathbf{\Sigma}}, \tag{106}$$

with an absolute constant C > 0 and the ratio:

$$r_{j} = \frac{\operatorname{tr}\left(\widehat{\Sigma}_{j}^{(d)}\right)}{\lambda_{\max}\left(\widehat{\Sigma}_{j}^{(d)}\right)}.$$
(107)

Let us integrate the definitions of  $\widehat{A}_i^{(d)}$  and  $\widehat{B}_i^{(d)}$  (17) with (105)-(106). Then, with probability of at least  $1 - \delta_{i, \mu}^{\rm PS} - \delta_{i, \Sigma}^{\rm PS}$ ,  $\widehat{A}_i^{(d)}$  and  $\widehat{B}_i^{(d)}$  deviate from their true versions by  $\varepsilon^{\rm PS}\left(\widehat{A}_i^{(d)}\right)$  and  $\varepsilon^{\rm PS}\left(\widehat{B}_i^{(d)}\right)$ ,

bounded by:

$$\varepsilon^{\text{PS}}\left(\widehat{A}_{i}^{(d)}\right) \leq 2\sqrt{\widehat{A}_{i}^{(d)}}\Delta_{i,\mu}\sqrt{\frac{\lambda_{\max}\left(\widehat{\Sigma}_{i}^{(d)}\right)}{\widetilde{\lambda}_{\min}\left(\widehat{\Sigma}_{i}^{(d)}\right)}} + \widehat{A}_{i}^{(d)}\frac{\Delta_{i,\Sigma}}{\lambda_{\max}\left(\widehat{\Sigma}_{i}^{(d)}\right)},\tag{108}$$

$$\varepsilon^{\text{PS}}\left(\widehat{B}_{i}^{(d)}\right) \leq 2\sqrt{\widehat{B}_{i}^{(d)}} \Delta_{j^{*}, \boldsymbol{\mu}} \sqrt{\frac{\lambda_{\max}\left(\widehat{\boldsymbol{\Sigma}}_{j^{*}}^{(d)}\right)}{\widetilde{\lambda}_{\min}\left(\widehat{\boldsymbol{\Sigma}}_{j^{*}}^{(d)}\right)} + \widehat{B}_{i}^{(d)} \frac{\Delta_{j^{*}, \boldsymbol{\Sigma}}}{\lambda_{\max}\left(\widehat{\boldsymbol{\Sigma}}_{j^{*}}^{(d)}\right)}}.$$
(109)

We avoid extremely loose bounds by replacing tiny, rarely observable eigenvalues, by a robust floor eigenvalue. Given a matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$ , we define  $\tilde{\lambda}_{\min}(\mathbf{A})$  as (Horn & Johnson, 2013, Thm. 4.3.1):

$$\tilde{\lambda}_{\min}(\mathbf{A}) = \lambda_{\min}\left(\mathbf{A} + \epsilon_r \lambda_{\max}(\mathbf{A}) I^{(d)}\right), \tag{110}$$

where  $\epsilon_r = 0.05$  is typically taken and  $I^{(d)}$  is the identity matrix. Ultimately, let us define the Euclidean Lipschitz constant  $L_{i,\text{lip}}$  as:

$$L_{i}^{\text{PS}} = \sqrt{\left(\frac{\partial \text{PS}_{i}}{\partial A_{i}} \left(A_{i}^{(d)}, B_{i}^{(d)}\right)\right)^{2} + \left(\frac{\partial \text{PS}_{i}}{\partial B_{i}} \left(A_{i}^{(d)}, B_{i}^{(d)}\right)\right)^{2}} = \frac{\sqrt{\left(\widehat{A}_{i}^{(d)}\right)^{2} + \left(\widehat{B}_{i}^{(d)}\right)^{2}}}{\left(\widehat{A}_{i}^{(d)} + \widehat{B}_{i}^{(d)}\right)^{2}}, \quad (111)$$

which enables us to bound the finite-sample deviation of  $\widehat{\mathrm{PS}}_i^{(d)}$  with:

$$\left| \widehat{\mathrm{PS}}_{i}^{(d)} - \mathrm{PS}_{i}^{(d)} \right| \le L_{i}^{\mathrm{PS}} \sqrt{\varepsilon^{\mathrm{PS}} \left( \widehat{A}_{i}^{(d)} \right) + \varepsilon^{\mathrm{PS}} \left( \widehat{B}_{i}^{(d)} \right)}. \tag{112}$$

Finally, we employ the triangle inequality on both error sources (102) and (112) and obtain for  $\delta_i^{\text{PS}} = \delta_{i,\mu}^{\text{PS}} + \delta_{i,\Sigma}^{\text{PS}}$ , with  $\delta_i^{\text{PS}} \in (0,1)$ :

$$\mathbb{P}_{\pi} \left\{ \left| \widehat{\mathsf{PS}}_{i}^{(d)} - \mathsf{PS}_{i} \right| \le (113)$$

$$\frac{\hat{B}_{i}^{(d)}\sqrt{\boldsymbol{r}_{i,i}^{T}\boldsymbol{S}_{i}^{-1}\boldsymbol{r}_{i,i}}+\hat{A}_{i}^{(d)}\sqrt{\boldsymbol{r}_{i,j^{*}}^{T}\boldsymbol{S}_{j^{*}}^{-1}\boldsymbol{r}_{i,j^{*}}}}{\left(\hat{A}_{i}^{(d)}+\hat{B}_{i}^{(d)}\right)^{2}}+L_{i}^{\mathrm{PS}}\sqrt{\varepsilon^{\mathrm{PS}}\left(\widehat{A}_{i}^{(d)}\right)+\varepsilon^{\mathrm{PS}}\left(\widehat{B}_{i}^{(d)}\right)}\right\}\geq1-\delta_{i}^{\mathrm{PS}}.$$

We now analyze the obtained expression separately for the deterministic and probabilistic terms. In the former term, the two square-root terms are energies that leak into the truncated complement after regressing out the retained d diffusion coordinates. Intuitively,  $\Sigma_j^{(d)}$  encodes the local anisotropy of cluster j in the kept coordinates,  $C_j$  represents coupling of residual energy in the truncated block, and  $\Sigma_j^{(m)}$  is the spread that remains in the truncated block. Therefore, larger  $S_j$  down-weights complement deviations, reducing the bias, while  $r_{i,j}$  and  $S_j$  co-vary through the cross-covariance  $C_j$ . Namely, increasing  $C_j$  shrinks both  $r_{i,j}$  and  $S_j$ , while decreasing  $C_j$  does the opposite. The practical rule is to prevent tiny  $\lambda_{\min}(S_j)$ , e.g., by promoting such directions into the kept set via a local choice of d, and shape distortions so the complement is predictable from the kept coordinates, keeping  $r_{i,j}$  small.

For the probabilistic part, its width reflects uncertainty in the empirical centroid and covariance of the attributed and nearest foreign clusters, where  $\lambda_{\max}(\widehat{\Sigma}_j^{(d)})$  and  $n_{j,\text{eff}}$  determine the width primarily. Practically, correlated distortions shrink  $n_{j,\text{eff}}$  and widen the bound, and a smaller spectral flatness ratio  $r_j$  yields tighter matrix concentration. As expected, the probabilistic piece dominates the deterministic as shown in Figure 8, which underlines the importance of cluster construction and dependence control.

# E.2 THE PM MEASURE

As in the PS case, we start with the truncation error and assume access to cluster statistics. Let us reconsider (78)-(90), but with two adjustments. First, the cluster coordinates are centered around the cluster reference embedding and not the cluster mean. Given m = N - d - 1, we define:

$$\boldsymbol{\Delta}_{i,p}^{(d)} = \boldsymbol{\Psi}_t^{(d)}(\mathbf{x}_{i,p}) - \boldsymbol{\Psi}_t^{(d)}(\mathbf{x}_i) \in \mathbb{R}^d, \tag{114}$$

$$\Delta_{i,p}^{(m)} = \Psi_t^{(m)}(\mathbf{x}_{i,p}) - \Psi_t^{(m)}(\mathbf{x}_i) \in \mathbb{R}^m.$$
(115)

Second, the cluster is now absent the reference embedding. Namely, the full N-1-dimensional cluster is (83):

$$\tilde{\mathcal{C}}_i = \mathcal{C}_i \setminus \Psi_t(\mathbf{x}_i). \tag{116}$$

The cluster  $\tilde{\mathcal{C}}_i$  has difference  $\Delta_{i,p} \in \mathbb{R}^{N-1}$  for every  $p \in \{1,\dots,N_p\}$  and covariance  $\tilde{\Sigma}_i \in \mathbb{R}^{(N-1)\times(N-1)}$  that hold (19):

$$\boldsymbol{\Delta}_{i,p} = \begin{bmatrix} \boldsymbol{\Delta}_{i,p}^{(d)} \\ \boldsymbol{\Delta}_{i,p}^{(m)} \end{bmatrix}, \quad \tilde{\boldsymbol{\Sigma}}_{i} = \begin{bmatrix} \tilde{\boldsymbol{\Sigma}}_{i}^{(d)} & \tilde{\boldsymbol{C}}_{i} \\ \tilde{\boldsymbol{C}}_{i}^{T} & \tilde{\boldsymbol{\Sigma}}_{i}^{(m)} \end{bmatrix}, \tag{117}$$

with  $ilde{m{\Sigma}}_i^{(m)} \in \mathbb{R}^{m imes m}$  and  $ilde{m{C}}_i \in \mathbb{R}^{d imes m}$  being:

$$\tilde{\Sigma}_{i}^{(m)} = \frac{1}{\left|\tilde{\mathcal{C}}_{i}^{(m)}\right| - 1} \sum_{\boldsymbol{\psi} \in \tilde{\mathcal{C}}_{i}^{(m)}} \left(\boldsymbol{\psi} - \Psi_{t}^{(m)}(\mathbf{x}_{i})\right) \left(\boldsymbol{\psi} - \Psi_{t}^{(m)}(\mathbf{x}_{i})\right)^{T}, \tag{118}$$

$$\tilde{C}_{i} = \frac{1}{\left|\tilde{C}_{i}^{(m)}\right| - 1} \sum_{p=1}^{N_{p}} \left(\Psi_{t}^{(d)}(\mathbf{x}_{i,p}) - \Psi_{t}^{(d)}(\mathbf{x}_{i})\right) \left(\Psi_{t}^{(m)}(\mathbf{x}_{i,p}) - \Psi_{t}^{(m)}(\mathbf{x}_{i})\right)^{T}.$$
 (119)

In N-1 dimensions, the squared Mahalanobis distance from  $\Psi_t(\mathbf{x}_{i,p})$  to  $\tilde{\mathcal{C}}_i$  is given by (16):

$$d_M^2\left(\mathbf{\Psi}_t(\mathbf{x}_{i,p}); \mathbf{\Psi}_t(\mathbf{x}_i), \tilde{\mathbf{\Sigma}}_i\right) = \mathbf{\Delta}_{i,p}^T \left(\tilde{\mathbf{\Sigma}}_i + \epsilon I^{(N-1)}\right)^{-1} \mathbf{\Delta}_{i,p}, \tag{120}$$

where as in (87), inversion has been empirically obtained with  $\epsilon = 10^{-6}$  and the N-1-dimensional identity matrix  $I^{(N-1)}$ . We again turn to the Schur complement (Horn & Johnson, 2013) and decompose (120):

$$d_M^2\left(\boldsymbol{\Psi}_t(\mathbf{x}_{i,p});\boldsymbol{\Psi}_t(\mathbf{x}_i),\tilde{\boldsymbol{\Sigma}}_i\right) = \left(\boldsymbol{\Delta}_{i,p}^{(d)}\right)^T \left(\tilde{\boldsymbol{\Sigma}}_i^{(d)} + \epsilon I^{(d)}\right)^{-1} \boldsymbol{\Delta}_{i,p}^{(d)} + \boldsymbol{r}_{i,p}^T \boldsymbol{S}_i^{-1} \boldsymbol{r}_{i,p}, \tag{121}$$

with  $r_{i,p} \in \mathbb{R}^m$  and the Schur complement  $S_i \in \mathbb{R}^{m \times m}$  being:

$$\boldsymbol{r}_{i,p} = \boldsymbol{\Delta}_{i,p}^{(m)} - \tilde{\boldsymbol{C}}_i^T \left( \tilde{\boldsymbol{\Sigma}}_i^{(d)} + \epsilon \boldsymbol{I}^{(d)} \right)^{-1} \boldsymbol{\Delta}_{i,p}^{(d)}, \tag{122}$$

$$S_i = \tilde{\Sigma}_i^{(m)} - \tilde{C}_i^T \left( \tilde{\Sigma}_i^{(d)} + \epsilon I^{(d)} \right)^{-1} \tilde{C}_i.$$
 (123)

Let us define the set of squared Mahalanobis distances of cluster  $\tilde{C}_i$  in dimension N-1 as:

$$\mathcal{G}_i = \left\{ d_M^2 \left( \mathbf{\Psi}_t(\mathbf{x}_{i,p}); \mathbf{\Psi}_t(\mathbf{x}_i), \widetilde{\mathbf{\Sigma}}_i \right) \mid p = 1, \dots, N_p \right\},$$
 (124)

in accordance to the truncated version of  $\mathcal{G}_i^{(d)}$  in (20). By employing (121), for every  $p \in \{1, \dots, N_p\}$ , we can bound the truncation error of the squared Mahalanobis distance as follows:

$$d_M^2\left(\boldsymbol{\Psi}_t(\mathbf{x}_{i,p}); \boldsymbol{\Psi}_t(\mathbf{x}_i), \tilde{\boldsymbol{\Sigma}}_i\right) - d_M^2\left(\boldsymbol{\Psi}_t^{(d)}(\mathbf{x}_{i,p}); \boldsymbol{\Psi}_t^{(d)}(\mathbf{x}_i), \tilde{\boldsymbol{\Sigma}}_i^{(d)}\right) = \mathbf{r}_{i,p}^T \mathbf{S}_i^{-1} \mathbf{r}_{i,p} := \delta_{\mathcal{G}_{i,p}},$$
(125)

and the difference between the mean of the elements in  $\mathcal{G}_i$  and  $\mathcal{G}_i^{(d)}$  can be expressed as:

$$\mu_{\mathcal{G}_i} - \mu_{\mathcal{G}_i^{(d)}} = \frac{1}{|\mathcal{G}_i|} \sum_{g \in \mathcal{G}_i} g - \frac{1}{|\mathcal{G}_i^{(d)}|} \sum_{g \in \mathcal{G}_i^{(d)}} g = \frac{1}{N_p} \sum_{p=1}^{N_p} \mathbf{r}_{i,p}^T \mathbf{S}_i^{-1} \mathbf{r}_{i,p} =$$
(126)

$$\frac{1}{N_p} \sum_{p=1}^{N_p} \delta_{\mathcal{G}_i,p} := \delta_{\mathcal{G}_i,\mu}.$$

Similarly, we can express the deviation of the variance:

$$\sigma_{\mathcal{G}_{i}}^{2} - \sigma_{\mathcal{G}_{i}^{(d)}}^{2} = \frac{1}{|\mathcal{G}_{i}| - 1} \sum_{g \in \mathcal{G}_{i}} (g - \mu_{\mathcal{G}_{i}})^{2} - \frac{1}{|\mathcal{G}_{i}^{(d)}| - 1} \sum_{g \in \mathcal{G}_{i}^{(d)}} \left(g - \mu_{\mathcal{G}_{i}^{(d)}}\right)^{2}, \tag{127}$$

and with (126) and the Cauchy-Schwartz inequality, we can obtain:

$$\left| \sigma_{\mathcal{G}_{i}}^{2} - \sigma_{\mathcal{G}_{i}^{(d)}}^{2} \right| \leq \frac{N_{p}}{N_{p} - 1} \left( 2\delta_{\mathcal{G}_{i}, p}^{\max} \left( \sigma_{\mathcal{G}_{i}} + \sigma_{\mathcal{G}_{i}^{(d)}} \right) + \left( \delta_{\mathcal{G}_{i}, p}^{\max} \right)^{2} \right), \tag{128}$$

where  $\delta_{\mathcal{G}_i,p}^{\max} = \max_p \delta_{\mathcal{G}_i,p}$ . The Gamma-matching parameters in the truncated and full dimensions are (22):

$$k_i^{(d)} = \frac{\mu_{\mathcal{G}_i^{(d)}}^2}{\sigma_{\mathcal{G}_i^{(d)}}^2}, \quad k_i = \frac{\mu_{\mathcal{G}_i}^2}{\sigma_{\mathcal{G}_i}^2},$$
 (129)

$$\theta_i^{(d)} = \frac{\sigma_{\mathcal{G}_i^{(d)}}^2}{\mu_{\mathcal{G}^{(d)}}}, \quad \theta_i = \frac{\sigma_{\mathcal{G}_i}^2}{\mu_{\mathcal{G}_i}}, \tag{130}$$

and their deviations can be bounded by considering (126), (128):

$$\left| k_i - k_i^{(d)} \right| \le C_1 \delta_{\mathcal{G}_i, p}^{\max} \frac{N_p}{N_p - 1} \frac{\mu_{\mathcal{G}_i} + \mu_{\mathcal{G}_i^{(d)}}}{\sigma_{\mathcal{G}_i^{(d)}}^2} := \delta_{\mathcal{G}_i, k},$$
 (131)

$$\left|\theta_{i} - \theta_{i}^{(d)}\right| \leq C_{2} \delta_{\mathcal{G}_{i}, p}^{\max} \frac{N_{p}}{N_{p} - 1} \frac{\sigma_{\mathcal{G}_{i}}^{2} + \sigma_{\mathcal{G}_{i}^{(d)}}^{2}}{\mu_{\mathcal{G}_{i}^{(d)}}^{2}} := \delta_{\mathcal{G}_{i}, \theta}, \tag{132}$$

with universal constants  $C_1, C_2 > 0$ . Let the squared Mahalanobis distance from the output embedding to the cluster be:

$$d_M^2\left(\mathbf{\Psi}_t^{(d)}(\hat{\mathbf{x}}_i); \mathbf{\Psi}_t^{(d)}(\mathbf{x}_i), \tilde{\mathbf{\Sigma}}_i^{(d)}\right) := a_i, \tag{133}$$

and employing (125) for the output embedding yields:

$$d_M^2\left(\boldsymbol{\Psi}_t(\hat{\mathbf{x}}_i);\boldsymbol{\Psi}_t(\mathbf{x}_i),\tilde{\boldsymbol{\Sigma}}_i\right) - d_M^2\left(\boldsymbol{\Psi}_t^{(d)}(\hat{\mathbf{x}}_i);\boldsymbol{\Psi}_t^{(d)}(\mathbf{x}_i),\tilde{\boldsymbol{\Sigma}}_i^{(d)}\right) = \boldsymbol{r}_{i,a}^T\boldsymbol{S}_i^{-1}\boldsymbol{r}_{i,a} := \delta_{\mathcal{G}_i,a}. \quad (134)$$

As in (23), the PM definition in dimension N-1 can be expressed using the regularized upper incomplete gamma function  $Q(k,x) = \Gamma(k,x)/\Gamma(k)$ :

$$PM_i = Q\left(k_i, \frac{a_i}{\theta_i}\right). \tag{135}$$

Consider the truncation-induced ellipsoid:

$$\mathcal{B}_{i} = \left\{ \left( k_{i}^{\prime}, \theta_{i}^{\prime}, a_{i}^{\prime} \right) : \left| k_{i}^{\prime} - k_{i}^{(d)} \right| \le \delta_{\mathcal{G}_{i}, k}, \left| \theta_{i}^{\prime} - \theta_{i}^{(d)} \right| \le \delta_{\mathcal{G}_{i}, \theta}, \left| a_{i}^{\prime} - a_{i}^{(d)} \right| \le \delta_{\mathcal{G}_{i}, a} \right\}, \tag{136}$$

For  $F(k, \theta, a) = Q(k, a/\theta)$ , the gradient with the partial derivatives with respect to  $k, \theta$  and a is:

$$\nabla F(k, \theta, a) = \begin{pmatrix} \frac{1}{\Gamma(k)} \int_{x}^{\infty} t^{k-1} e^{-t} \ln t \, dt - \psi(k) \, Q(k, x) \\ \frac{a}{\theta^{2}} \frac{x^{k-1} e^{-x}}{\Gamma(k)} \\ -\frac{1}{\theta} \frac{x^{k-1} e^{-x}}{\Gamma(k)} \end{pmatrix}, \tag{137}$$

where  $x = a/\theta$  and  $\psi(\cdot)$  is the digamma function. Since  $\nabla F(k, \theta, a)$  is continuous and bounded on the compact set  $\mathcal{B}_i$ , we set:

$$L_{\mathcal{B}_i} = \sup_{(k,\theta,a)\in\mathcal{B}_i} \|\nabla F(k,\theta,a)\|_2 < \infty.$$
 (138)

To yield the bound on the PM measure due to truncation, we notice that both  $(k_i, \theta_i, a_i)$  and  $(k_i^{(d)}, \theta_i^{(d)}, a_i^{(d)})$  lie in  $\mathcal{B}_i$ , and apply the multivariate mean-value theorem to yield the following:

$$\left| \operatorname{PM}_{i} - \operatorname{PM}_{i}^{(d)} \right| \leq L_{\mathcal{B}_{i}} \left( \delta_{\mathcal{G}_{i},k}^{2} + \delta_{\mathcal{G}_{i},\theta}^{2} + \delta_{\mathcal{G}_{i},a}^{2} \right)^{1/2}.$$
(139)

However, this bound can be tightened. We notice that Q(k,x) (23) is monotonically increasing in k and decreasing in k, for k, k = 0 (137). We assume that on  $\mathcal{B}_i$ , and for all k = k = k does not change signs, or otherwise we fallback to (139). Consequently, the maximal change of k = k is attained at one of its eight corners, and (139) can be tightened to this PM error radius:

$$\left| \operatorname{PM}_{i} - \operatorname{PM}_{i}^{(d)} \right| \leq \max_{(k_{c}, \theta_{c}, a_{c}) \in \partial \mathcal{B}_{i}} \left| Q\left(k_{c}, a_{c} / \theta_{c}\right) - Q\left(\hat{k}_{i}^{(d)}, \hat{a}_{i}^{(d)} / \hat{\theta}_{i}^{(d)}\right) \right|. \tag{140}$$

As in the PS case, we now analyze how the finite number of coordinates in a cluster leads to uncertainty in the PM evaluation. Let  $R_i$  be the maximal squared Mahalanobis distance in  $\mathcal{G}_i$ , namely:

$$R_i = \max_{g \in \mathcal{G}_i} g. \tag{141}$$

Again, similarly to the PS case, we utilize the vector and matrix Bernstein (Vershynin, 2024, Props. 2.8.1, 4.7.1) and the dependent Hanson-Wright inequalities (Adamczak, 2015, Thm. 2.5). Let us consider the confidence parameters  $\delta^{\text{PM}}_{i,\mu}$ ,  $\delta^{\text{PM}}_{i,\sigma}$ ,  $\delta^{\text{PM}}_{i,\sigma}$ ,  $\delta^{\text{PM}}_{i,a}$   $\in$  (0,1/3), so with respective least probabilities of  $1-\delta^{\text{PM}}_{i,\mu}$ ,  $1-\delta^{\text{PM}}_{i,\sigma}$ ,  $1-\delta^{\text{PM}}_{i,\sigma}$ :

$$\left| \mu_{\mathcal{G}_{i}^{(d)}} - \hat{\mu}_{\mathcal{G}_{i}^{(d)}} \right| \leq \sqrt{\frac{2\hat{\sigma}_{\mathcal{G}_{i}^{(d)}}^{2} \ln\left(2/\delta_{i,\mu}^{\text{PM}}\right)}{N_{p}}} + \frac{3R_{i} \ln\left(2/\delta_{i,\mu}^{\text{PM}}\right)}{N_{p}} := \Delta_{i,\mu}, \tag{142}$$

$$\left| \sigma_{\mathcal{G}_i^{(d)}} - \hat{\sigma}_{\mathcal{G}_i^{(d)}} \right| \le \sqrt{\frac{2R_i^2 \ln\left(2/\delta_{i,\sigma}^{\text{PM}}\right)}{N_p}} + \frac{3R_i^2 \ln\left(2/\delta_{i,\sigma}^{\text{PM}}\right)}{N_p} := \Delta_{i,\sigma},\tag{143}$$

$$\left| a_i - \hat{a}_i \right| \le R_i \sqrt{\frac{\ln\left(2/\delta_{i,a}^{\text{PM}}\right)}{N_p}} := \Delta_{i,a}. \tag{144}$$

Recalling the definition of  $k_i^{(d)}$ ,  $\theta_i^{(d)}$  from (129), (130), since by design  $\hat{\mu}_{\mathcal{G}_i^{(d)}}$ ,  $\hat{\sigma}_{\mathcal{G}_i^{(d)}} > 0$ , and since we empirically validate that  $\Delta_{i,\mu} \ll \hat{\mu}_{\mathcal{G}_i^{(d)}}$ ,  $\Delta_{i,\sigma} \ll \hat{\sigma}_{\mathcal{G}_i^{(d)}}$ , we can apply the first-order Taylor expansions to  $k_i^{(d)}$ ,  $\theta_i^{(d)}$  around  $\hat{k}_i^{(d)}$ ,  $\hat{\theta}_i^{(d)}$ , respectively. Apply the triangle inequality to it gives:

$$\left| k_i^{(d)} - \hat{k}_i^{(d)} \right| \le \left| \frac{\partial k_i^{(d)}}{\partial \mu} \right| \Delta_{i,\mu} + \left| \frac{\partial k_i^{(d)}}{\partial \sigma} \right| \Delta_{i,\sigma} = \left| \frac{2\hat{\mu}_{\mathcal{G}_i^{(d)}}}{\hat{\sigma}_{\mathcal{G}_i^{(d)}}^2} \right| \Delta_{i,\mu} + \left| \frac{-2\hat{\mu}_{\mathcal{G}_i^{(d)}}^2}{\hat{\sigma}_{\mathcal{G}_i^{(d)}}^3} \right| \Delta_{i,\sigma} := \Delta_{i,k}, \quad (145)$$

$$\left| \theta_i^{(d)} - \hat{\theta}_i^{(d)} \right| \le \left| \frac{\partial \theta_i^{(d)}}{\partial \mu} \right| \Delta_{i,\mu} + \left| \frac{\partial \theta_i^{(d)}}{\partial \sigma} \right| \Delta_{i,\sigma} = \left| \frac{-\hat{\sigma}_{\mathcal{G}_i^{(d)}}^2}{\hat{\mu}_{\mathcal{G}_i^{(d)}}^2} \right| \Delta_{i,\mu} + \left| \frac{2\hat{\sigma}_{\mathcal{G}_i^{(d)}}}{\hat{\mu}_{\mathcal{G}_i^{(d)}}^2} \right| \Delta_{i,\sigma} := \Delta_{i,\theta}. \tag{146}$$

Empirically, rarely  $\Delta_{i,k}$ ,  $\Delta_{i,\theta}$  or  $\Delta_{i,a}$  become extremely loose. To avoid this behavior, we practically regularize the box by setting:

$$\Delta_{i,k} \to \min\left(\Delta_{i,k}, 0.5k_i^{(d)}\right),$$
 (147)

$$\Delta_{i,\theta} \to \min\left(\Delta_{i,\theta}, 0.5\theta_i^{(d)}\right),$$
 (148)

$$\Delta_{i,a} \to \min\left(\Delta_{i,a}, 0.5a_i^{(d)}\right).$$
 (149)

Let us consider the local box of values:

$$\mathcal{B}_{i}^{\text{loc}} = \left\{ \left( k_{i}', \theta_{i}', a_{i}' \right) : k_{i}' \in \left\lceil k_{i}^{(d)} \pm \Delta_{i,k} \right\rceil, \theta_{i}' \in \left\lceil \theta_{i}^{(d)} \pm \Delta_{i,\theta} \right\rceil, a_{i}' \in \left\lceil a_{i}^{(d)} \pm \Delta_{i,a} \right\rceil \right\}. \tag{150}$$

As discussed earlier, Q(k,x) is monotonically increasing in k and decreasing in x, for k, x > 0 (137). Consequently, the maximal change of Q(k,x) inside  $\mathcal{B}_i^{\mathrm{loc}}$  is attained at one of its eight corners. Thus, the finite-sample error of the PM measure in dimension d is bounded by:

$$\left| PM_i^{(d)} - \widehat{PM}_i^{(d)} \right| \le \max_{(k_c, \theta_c, a_c) \in \partial \mathcal{B}_i^{loc}} \left| Q\left(k_c, a_c/\theta_c\right) - Q\left(\hat{k}_i^{(d)}, \hat{a}_i^{(d)}/\hat{\theta}_i^{(d)}\right) \right|. \tag{151}$$

Ultimately, we combine the deterministic error radius with the probabilistic width. Let  $\delta_i^{\text{PM}} = \delta_{i,\mu}^{\text{PM}} + \delta_{i,\sigma}^{\text{PM}} + \delta_{i,a}^{\text{PM}}$ , which yields for  $\delta_i^{\text{PM}} \in (0,1)$ :

$$\mathbb{P}_{\pi} \left\{ \left| \widehat{PM}_{i}^{(d)} - PM_{i} \right| \leq \frac{1}{(k_{c}, \theta_{c}, a_{c}) \in \partial \mathcal{B}_{i}} \left| Q(k_{c}, a_{c}/\theta_{c}) - Q(\hat{k}_{i}^{(d)}, \hat{a}_{i}^{(d)}/\hat{\theta}_{i}^{(d)}) \right| + \frac{1}{(k_{c}, \theta_{c}, a_{c}) \in \partial \mathcal{B}_{i}^{\text{loc}}} \left| Q(k_{c}, a_{c}/\theta_{c}) - Q(\hat{k}_{i}^{(d)}, \hat{a}_{i}^{(d)}/\hat{\theta}_{i}^{(d)}) \right| \right\} \geq 1 - \delta_{i}^{\text{PM}}.$$
(152)

In the deterministic term, large cross-block coupling  $\widetilde{C}_i$  or residual spread  $\widetilde{\Sigma}_i^{(m)}$  again directly inflate the error radius via the Schur complement.

In the probabilistic part, the local box  $\mathcal{B}_i^{\mathrm{loc}}$  aggregates two finite-sample pieces. The first is the uncertainty of the moment, with  $\Delta_{i,\mu}$  and  $\Delta_{i,\sigma}$  scale as  $1/N_p$  but are amplified by the maximal radius of Mahalanobis within the cluster  $R_i$ . Heavy outliers increase  $R_i$  and widen both bounds. The second is the uncertainty of the distance of the output, contributed by  $\Delta_{i,a}$  which is also proportional to  $R_i$  but scales by  $1/\sqrt{N_p}$ . Again, this emphasizes the importance of the design of distortions.

# F ERROR RADIUS AND PROBABILISTIC CONFIDENCE BOUND OF THE PCC AND SRCC

In this Appendix, we propagate the frame-level error radius and probabilistic widths developed in Appendix E.1 and E.2 to the reported PCC and SRCC values.

We start by fixing a trial l, a source separation system q, and a time frame f. Let the indices of the active sources in frame f be  $\mathcal{S}_f^l$  and consider a source  $i \in \mathcal{S}_f^l$ . The observation of measure  $\mathcal{P} \in \{\mathrm{PS},\mathrm{PM}\}$ , denoted  $\widehat{v}_{i,f}^{q,l,\mathcal{P}}$ , can be decomposed as:

$$\widehat{v}_{i,f}^{q,l,\mathcal{P}} = v_{i,f}^{q,l,\mathcal{P}} + \widetilde{\beta}_{i,f}^{q,l,\mathcal{P}} + \zeta_{i,f}^{q,l,\mathcal{P}}, \tag{153}$$

where:

$$\widetilde{\beta}_{i,f}^{q,l,\mathcal{P}} = \beta_{i,f}^{q,l,\mathcal{P}} + \mu_{i,f}^{q,l,\mathcal{P}},\tag{154}$$

and  $\beta_{i,f}^{q,l,\mathcal{P}}$  is an unknown deterministic bias with a provided radius  $b_{i,f}^{q,l,\mathcal{P}}$ , such that:

$$\left|\beta_{i,f}^{q,l,\mathcal{P}}\right| \le b_{i,f}^{q,l,\mathcal{P}},\tag{155}$$

with  $b_{i,f}^{q,l,\mathcal{P}}$  given by either (102) or (140). Regarding the probabilistic side, we define:

$$\zeta_{i,f}^{q,l,\mathcal{P}} = \varepsilon_{i,f}^{q,l,\mathcal{P}} - \mu_{i,f}^{q,l,\mathcal{P}},\tag{156}$$

where:

$$\varepsilon_{i,f}^{q,l,\mathcal{P}} = \widehat{v}_{i,f}^{q,l,\mathcal{P}} - v_{i,f}^{q,l,\mathcal{P}},\tag{157}$$

and  $\mathbb{E}_{\pi}\left(\zeta_{i,f}^{q,l,\mathcal{P}}\right)=0$ . Thus, the two-sided probabilistic half-width  $p_{i,f}^{q,l,\mathcal{P}}\geq 0$  can be interpreted as:

$$\mathbb{P}_{\pi}\left(\left|\varepsilon_{i,f}^{q,l,\mathcal{P}} - \mu_{i,f}^{q,l,\mathcal{P}}\right| \le p_{i,f}^{q,l,\mathcal{P}}\right) \ge 1 - \delta^{\mathcal{P}} := c^{\mathcal{P}},\tag{158}$$

with  $\delta^{\mathcal{P}}$  and the probabilistic bounds defined in (112) and (151). We abbreviate  $c^{\mathcal{P}}$  as c from now on. Consider  $z_{c^*}$  the normal quantile at level  $c^* = (1+c)/2$ , so we calibrate the half-widths scale to be:

$$\sigma_{i,f}^{q,l,\mathcal{P}} = \frac{c}{z_{c^*}},\tag{159}$$

with tails still reported back as half-widths at the original confidence c.

We now propagate these errors from frame to utterance level, based on the aggregations we introduced in (45) and (49). On average, experiments showed that frames more than g=4 apart are effectively independent both for speech and music mixtures. Given the set  $\mathcal{F}^l$  of time frames with two or more active sources, the standard Bartlett block-decimation (Bartlett, 1946) yields the conservative inflation:

$$\operatorname{std}\left(\frac{1}{\mathcal{F}^{l}}\sum_{f=1}^{\mathcal{F}^{l}}\zeta_{i,f}^{q,l,\mathcal{P}}\right) \leq \frac{\sqrt{g+1}}{\sqrt{\mathcal{F}^{l}}}\left(\frac{1}{\mathcal{F}^{l}}\sum_{f=1}^{\mathcal{F}^{l}}\left(\sigma_{i,f}^{q,l,\mathcal{P}}\right)^{2}\right)^{1/2}.$$
(160)

Let the radius error and the p-level probabilistic half-width obtained at the utterance-level using average pooling equal, respectively:

$$b_{i,\text{average}}^{q,l,\mathcal{P}} = \frac{1}{\mathcal{F}^l} \sum_{f=1}^{\mathcal{F}^l} b_{i,f}^{q,l,\mathcal{P}},\tag{161}$$

$$h_{i,\text{average}}^{q,l,\mathcal{P}} = z_{c^*} \frac{\sqrt{g+1}}{\sqrt{\mathcal{F}^l}} \left( \frac{1}{\mathcal{F}^l} \sum_{f=1}^{\mathcal{F}^l} \left( \sigma_{i,f}^{q,l,\mathcal{P}} \right)^2 \right)^{1/2}. \tag{162}$$

For the PESQ-like aggregation, let us denote its aggregation function from (49) as:

$$s(u) = 0.999 + 4(1 + \exp(-1.3669 u + 3.8224))^{-1}.$$
 (163)

Let W be the window and H the hop of frame used for aggregation, then  $M^l$  is the maximal number of possible windows. By norm submultiplicativity and the mean-value theorem (Horn & Johnson, 2013, Sec. 5.6):

$$b_{i,\text{pesq}}^{q,l,\mathcal{P}} = \frac{C_{\text{OL}}}{\sqrt{M^l}} \left( \frac{1}{\mathcal{F}^l} \sum_{f=1}^{\mathcal{F}^l} \left( b_f^{q,l,\mathcal{P}} \right)^2 \right)^{1/2} \frac{\partial s}{\partial u}, \tag{164}$$

$$h_{i,\text{pesq}}^{q,l,\mathcal{P}} = z_{c^*} \frac{C_{\text{OL}}}{\sqrt{M^l}} \left( \frac{1}{\mathcal{F}^l} \sum_{f=1}^{\mathcal{F}^l} \left( \sigma_f^{q,l,\mathcal{P}} \right)^2 \right)^{1/2} \frac{\partial s}{\partial u}, \tag{165}$$

where  $C_{\text{OL}} = \lceil W/H \rceil$  and by construction  $\partial s/\partial u \leq 1.3669$  when evaluated at point u.

To translate utterance-level errors to source-based PCC and SRCC values, let the integration of utterance-level MOS ratings from every system  $q \in \{1, ..., Q\}$  be:

$$\mathbf{v}_{i,\text{MOS}}^l = \left(v_{i,\text{MOS}}^{1,l}, \dots, v_{i,\text{MOS}}^{Q,l}\right),\tag{166}$$

and similarly, denoting  $\hat{v}_{i,\mathcal{A}}^{q,l,\mathcal{P}}$  as the estimated aggregated measure across systems, where  $\mathcal{A}$  is either average or PESQ-like aggregation (§B.4), then its integration is given by:

$$\hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}} = \left(\hat{v}_{i,\mathcal{A}}^{1,l,\mathcal{P}}, \dots, \hat{v}_{i,\mathcal{A}}^{Q,l,\mathcal{P}}\right). \tag{167}$$

For every vector  $\mathbf{v}$ , we denote its centered version by  $\tilde{\mathbf{v}}$ . Let us denote the PCC value between an observation vector  $\mathbf{v}$  and a MOS vector  $\mathbf{m}$  as  $r^{\text{PCC}}(\mathbf{v}, \mathbf{m})$ , according to (53) and (54). Its gradient with respect to  $\mathbf{v}$  at point  $\hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}$  is (Benesty et al., 2009):

$$\frac{\partial r^{\text{PCC}}}{\partial \mathbf{v}} \bigg|_{\mathbf{v} = \hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}} = \frac{\mathbf{v}_{i,\text{MOS}}^{l}}{\left\|\hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}\right\|_{2} \left\|\mathbf{v}_{i,\text{MOS}}^{l}\right\|_{2}} - \frac{r^{\text{PCC}}\left(\hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}, \mathbf{v}_{i,\text{MOS}}^{l}\right)}{\left\|\hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}\right\|_{2}^{2}} \hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}.$$
(168)

Consider  $\mathbf{b}_{i}^{l,\mathcal{P}}$  the utterance-level bias radii from (161) or (164) across all systems:

$$\mathbf{b}_{i,\mathcal{A}}^{l,\mathcal{P}} = \left(b_{i,\mathcal{A}}^{1,l,\mathcal{P}}, \dots, b_{i,\mathcal{A}}^{Q,l,\mathcal{P}}\right). \tag{169}$$

Then, the induced PCC bias can be bounded by:

$$b_{i,\mathcal{A}}^{l,\mathcal{P},PCC} \leq \left\| \frac{\partial r^{PCC}}{\partial \hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}} \right\|_{2} \left\| \tilde{\mathbf{b}}_{i,\mathcal{A}}^{l,\mathcal{P}} \right\|_{2}. \tag{170}$$

For the probabilistic half-width, we model independent Gaussian jitters across systems with scales fixed by the utterance half-widths. Consider the Q-dimensional Gaussian vector:

$$\eta \sim \mathcal{N}\left(\mathbf{0}, \operatorname{diag}\left(\left(\frac{h_{i,\mathcal{A}}^{1,l,\mathcal{P}}}{z_{c^*}}\right)^2, \dots, \left(\frac{h_{i,\mathcal{A}}^{Q,l,\mathcal{P}}}{z_{c^*}}\right)^2\right)\right),$$
(171)

with  $\mathbf{0} \in \mathbb{R}^Q$ . Using the delta method, first-order error propagation gives:

$$h_{i,\mathcal{A}}^{l,\mathcal{P},PCC} = z_{c^*} \sqrt{\left(\frac{\partial r^{PCC}}{\partial \hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}}\right)^T \operatorname{diag}\left(\left(\frac{h_{i,\mathcal{A}}^{1,l,\mathcal{P}}}{z_{c^*}}\right)^2, \dots, \left(\frac{h_{i,\mathcal{A}}^{Q,l,\mathcal{P}}}{z_{c^*}}\right)^2\right) \left(\frac{\partial r^{PCC}}{\partial \hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}}\right)}.$$
 (172)

Turning to the SRCC, let  $\rho^{\rm SRCC}(\cdot,\cdot)$  denote Spearman's rank correlation between two vectors (Kendall & Gibbons, 1990), as defined in (55) and (56). Because ranks are piecewise-constant, a safe deterministic error radius is obtained by checking the two extreme bias orientations:

$$b_{i,\mathcal{A}}^{l,\mathcal{P},SRCC} = \max\left(\left|\rho^{SRCC}\left(\hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}} + \mathbf{b}_{i,\mathcal{A}}^{l,\mathcal{P}}, \mathbf{v}_{i,MOS}^{l}\right) - \rho^{SRCC}\left(\hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}, \mathbf{v}_{i,MOS}^{l}\right)\right|,$$

$$\left|\rho^{SRCC}\left(\hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}} - \mathbf{b}_{i,\mathcal{A}}^{l,\mathcal{P}}, \mathbf{v}_{i,MOS}^{l}\right) - \rho^{SRCC}\left(\hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}, \mathbf{v}_{i,MOS}^{l}\right)\right|\right).$$

$$(173)$$

For the probabilistic half-width we jitter  $\hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}$  with the same independent Gaussian model in (171) and report the empirical  $c^*$  quantile from Monte Carlo of the following:

$$h_{i,\mathcal{A}}^{l,\mathcal{P},SRCC} = \text{Quantile}_{c^*} \left( \left| \rho^{SRCC} \left( \hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}} + \boldsymbol{\eta}, \mathbf{v}_{i,MOS}^l \right) - \rho^{SRCC} \left( \hat{\mathbf{v}}_{i,\mathcal{A}}^{l,\mathcal{P}}, \mathbf{v}_{i,MOS}^l \right) \right| \right), \tag{174}$$

where we used  $10^4$  draws for estimation, in the spirit of quantile bootstrap (Efron & Tibshirani, 1994).

Lastly, we consider the error propagation across all trials and their sources in a given scenario, e.g., English mixtures. Let  $\mathcal L$  denote the number of trials in a scenario, and for each trial  $l \in \{1,\dots,\mathcal L\}$ , assume the number of total speakers in the trial is  $N_{\max}^l$  (57). The values we report average across all  $\mathcal L$  trials and  $N_{\max}^l$  speakers, following (58)-(61).

The deterministic error radius of the PCC and SRCC per scenario are respectively given by:

$$b^{\text{PCC}} = \frac{1}{\sum_{l=1}^{\mathcal{L}} N_{\text{max}}^{l}} \sum_{l=1}^{\mathcal{L}} \sum_{i=1}^{N_{\text{max}}^{l}} b_{i,\mathcal{A}}^{l,\mathcal{P},\text{PCC}},$$
(175)

$$b^{\text{SRCC}} = \frac{1}{\sum_{l=1}^{\mathcal{L}} N_{\text{max}}^{l}} \sum_{l=1}^{\mathcal{L}} \sum_{i=1}^{N_{\text{max}}^{l}} b_{i,\mathcal{A}}^{l,\mathcal{P},\text{SRCC}}.$$
 (176)

To yield the probabilistic term, we assume that within any fixed trial l, the pairwise correlation between the source jitters has been empirically estimated and is denoted by  $\rho_l$ , while jitters from different trials are independent. This assumption holds by the construction of our trials in every scenario. Consequently, the c-level probabilistic half-width on the scenario mean equals:

$$h^{\text{PCC}} = \tag{177}$$

$$z_{c^*} \sqrt{\frac{1}{\left(\sum_{l=1}^{\mathcal{L}} N_{\text{max}}^l\right)^2} \sum_{l=1}^{\mathcal{L}} \left(\sum_{i=1}^{N_{\text{max}}^l} \left(\frac{h_{i,\mathcal{A}}^{l,\mathcal{P},\text{PCC}}}{z_{c^*}}\right)^2 + 2\rho_l \sum_{\substack{i,j=1\\i < j}}^{N_{\text{max}}^l} \left(\frac{h_{i,\mathcal{A}}^{l,\mathcal{P},\text{PCC}}}{z_{c^*}}\right) \left(\frac{h_{j,\mathcal{A}}^{l,\mathcal{P},\text{PCC}}}{z_{c^*}}\right)\right)}.$$

$$h^{\text{SRCC}} = \tag{178}$$

$$z_{c^*} \sqrt{\frac{1}{\left(\sum_{l=1}^{\mathcal{L}} N_{\max}^l\right)^2} \sum_{l=1}^{\mathcal{L}} \left(\sum_{i=1}^{N_{\max}^l} \left(\frac{h_{i,\mathcal{A}}^{l,\mathcal{P},SRCC}}{z_{c^*}}\right)^2 + 2\rho_l \sum_{\substack{i,j=1\\i < j}}^{N_{\max}^l} \left(\frac{h_{i,\mathcal{A}}^{l,\mathcal{P},SRCC}}{z_{c^*}}\right) \left(\frac{h_{j,\mathcal{A}}^{l,\mathcal{P},SRCC}}{z_{c^*}}\right)\right)}.$$

Ultimately, for each scenario and each measure  $\mathcal{P}$  that uses aggregation technique  $\mathcal{A}$ , we report the deterministic envelope and probabilistic half-width  $b^{\text{PCC}}$  and  $h^{\text{PCC}}$  for PCC values and  $b^{\text{SRCC}}$  and  $h^{\text{SRCC}}$  for SRCC values.

## G FURTHER DISCUSSIONS

#### G.1 LIMITATIONS

Our validation depends exclusively on the SEBASS database, the only public corpus that provides human ratings for source separation systems, which limits the diversity in acoustic and linguistic traits that multiple dataset usually carry together. Moreover, the listening tests in SEBASS ask the human raters a generic quality question, rather than questions that isolate leakage versus self-distortion. This design choice may attenuate the ground-truth sensitivity to the specific error modes that PS and PM are intended to disentangle, and can introduce a systematic bias that even multi-rater averaging cannot fully cancel. Another noticeable limitation of this research concerns the aggregation techniques we employ to convert frame-level to utterance-level scores. Since neither granular human ratings exist nor is there any documented data-driven mapping from granular to global human ratings, we limit the capability of the PS and PM measures by merely approximating aggregation functions.

On a single NVIDIA A6000 GPU paired with 32 CPU cores with 64 GB of memory, our implementation achieves a real-time factor of 1.2, e.g., when analyzing a 25 ms frame in 30 ms on average. While this enables offline evaluation and hyper-parameter sweeps, it falls short of strict real-time monitoring and may limit large-scale neural-architecture searches and limit using the PS and PM measures inside loss function during training sessions. Profiling reveals that the dominant bottlenecks are diffusion-map eigendecompositions and repeated Mahalanobis distance computations with per-frame covariance estimation for all distortions points in every cluster. We plan to introduce more efficient implementations as we maintain our code repository.

We also point out that in music mixtures, 0.5% of frames exhibit for the PM measure an error radius that exceeds 1, rendering these observations irrelevant. These cases should be ignored completely, and future work that focuses on the separation of music sources will further investigate this phenomenon.

#### G.2 Positioning Our Work as a Catalyst

The absence of large, diverse datasets annotated with fine-grained human scores remains a critical gap in source separation research. We argue that introducing perceptually grounded measures is precisely what enables this gap to be closed. By releasing PS and PM as open-source tools, we provide the community with a foundation on which richer benchmark datasets can be built, rather than waiting for such datasets to exist before new measures are introduced. Their availability can catalyze the creation of corpora that include human annotations at both frame-level and utterance-level resolutions. Such resources would support systematic, fine-grained comparisons between objective measures and human perception, stimulate the development of new evaluation metrics and systems, and allow researchers to study the relationship between granular and global ratings, an aspect currently reduced to heuristic aggregation. In this way, PS and PM act as a gateway toward more rigorous and perceptually aligned evaluation standards in source separation.

## H LLM USAGE

We used a large language model (LLM) as a general-purpose assistant in three ways:

- 1. Language polishing to improve clarity. Every word was read and proofed by the authors.
- 2. Exploration of literature. All cited literature was validated by the authors.
- Coding assistance. All code was reviewed, rewritten as needed, and tested by the authors before use.

We did not delegate authorship decisions or scientific claims to the LLM. We manually verified all content, checked citations, and validated all results. No confidential or identity-revealing information was provided to the LLM, and use complied with dataset licenses and the ICLR Code of Ethics.