



澳門理工大學
Universidade Politécnica de Macau
Macao Polytechnic University



ICML
International Conference
On Machine Learning

2026

Spherical Procrustes Alignment for Reliable Medical Audio Diagnosis

Ying Wang¹, Guoheng Huang², Chan-Tong Lam¹, Xiaochen Yuan¹

¹*Macao Polytechnic University, xcyuan@mpu.edu.mo*

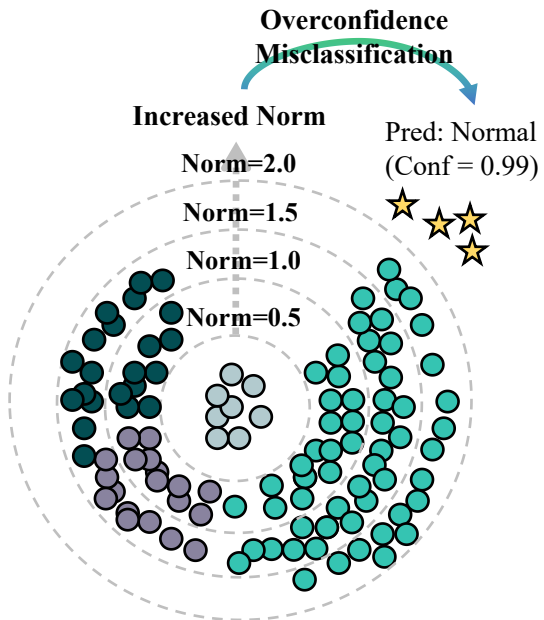
Code: github.com/wangying1586/SPA



Motivation – The Overconfidence Problem

norm-biased radial space

→ high-norm noise triggers 0.99 confidence

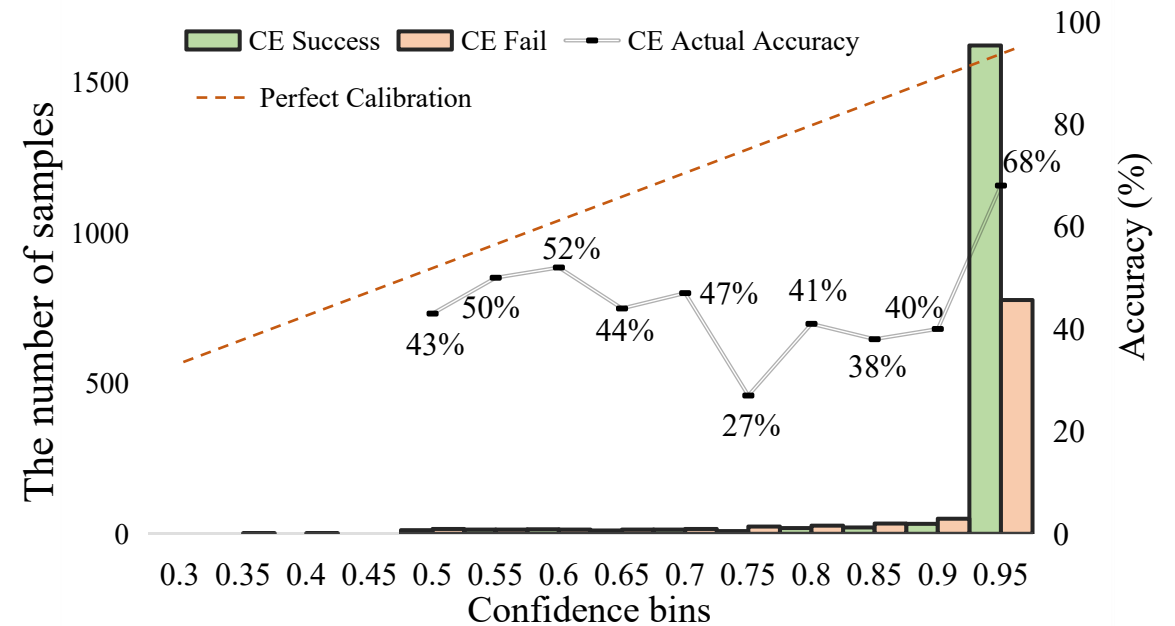


(a) CE

Reliability diagram (CE baseline)

→ far from diagonal (Perfect Calibration)

CE's true accuracy is around 60%, but most of the confidence distribution is at 0.95.



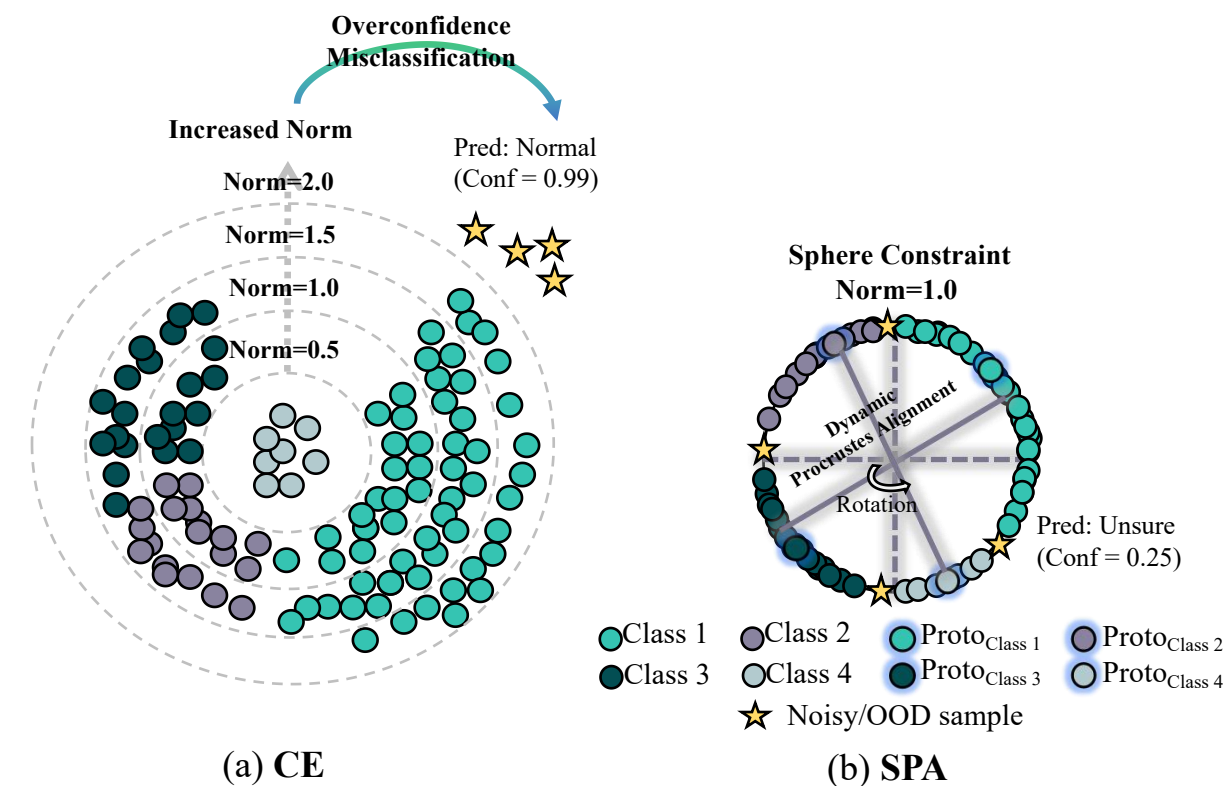
(c) CE vs. SPA

Clinical diagnosis needs both accuracy and well-calibrated confidence.

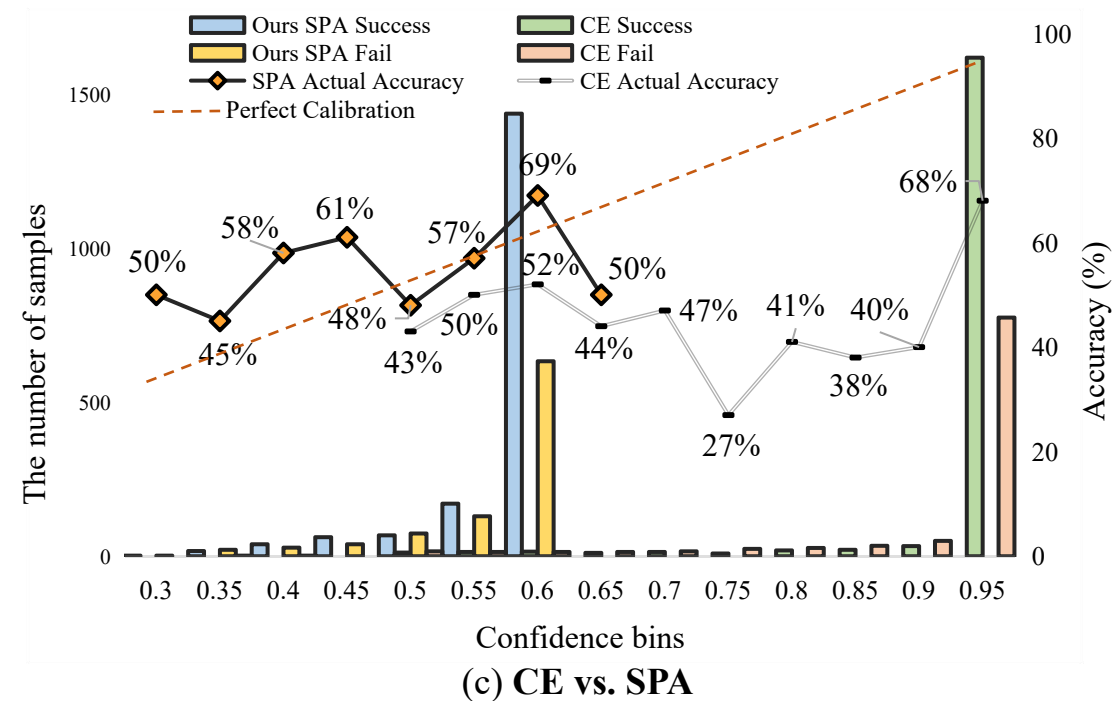
Root Cause & Our Insight

Root cause: Norm bias – logit magnitude dominates, not semantic alignment

Why existing methods fail: Loss-level or post-hoc; ignore geometric pathology



Ours SPA is closer to the diagonal (Perfect Calibration).

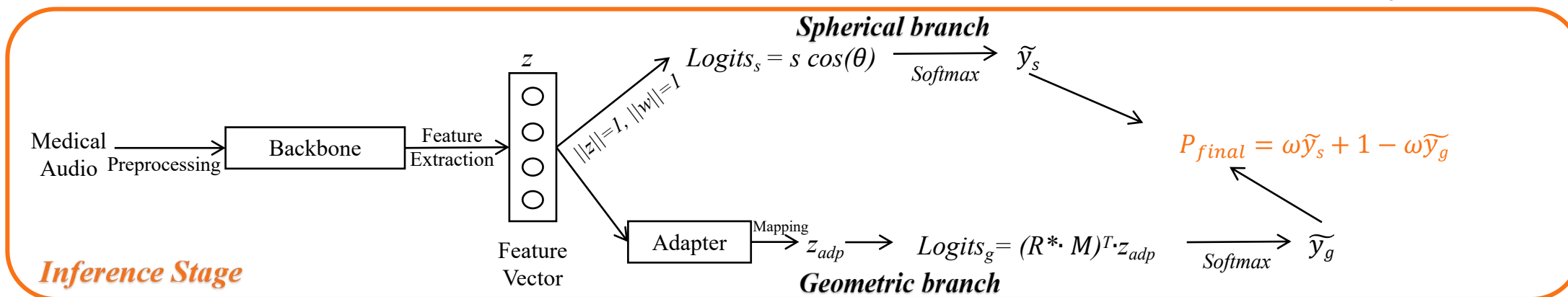
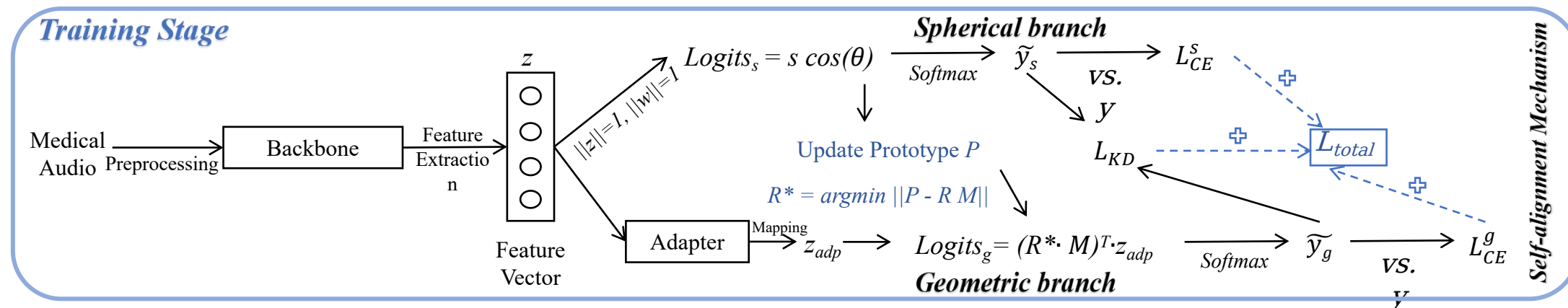


Our insight - Enforce spherical constraint plus dynamic geometric alignment with ETF



Spherical Procrustes Alignment Overview

- Training stage:**
- **Spherical branch:** unit hypersphere, eliminate norm bias
 - **Geometric branch:** fixed ETF plus dynamic Procrustes alignment
 - **Self-alignment Mechanism:** fuse via KL divergence



Inference stage: weighted fusion

Spherical Procrustes Alignment – Dual branch architecture

Spherical Branch – Decouple Norm

- Normalize features \tilde{z} and weights \tilde{W}
- Logit = $s \cdot \cos \theta \rightarrow$ confidence purely angular

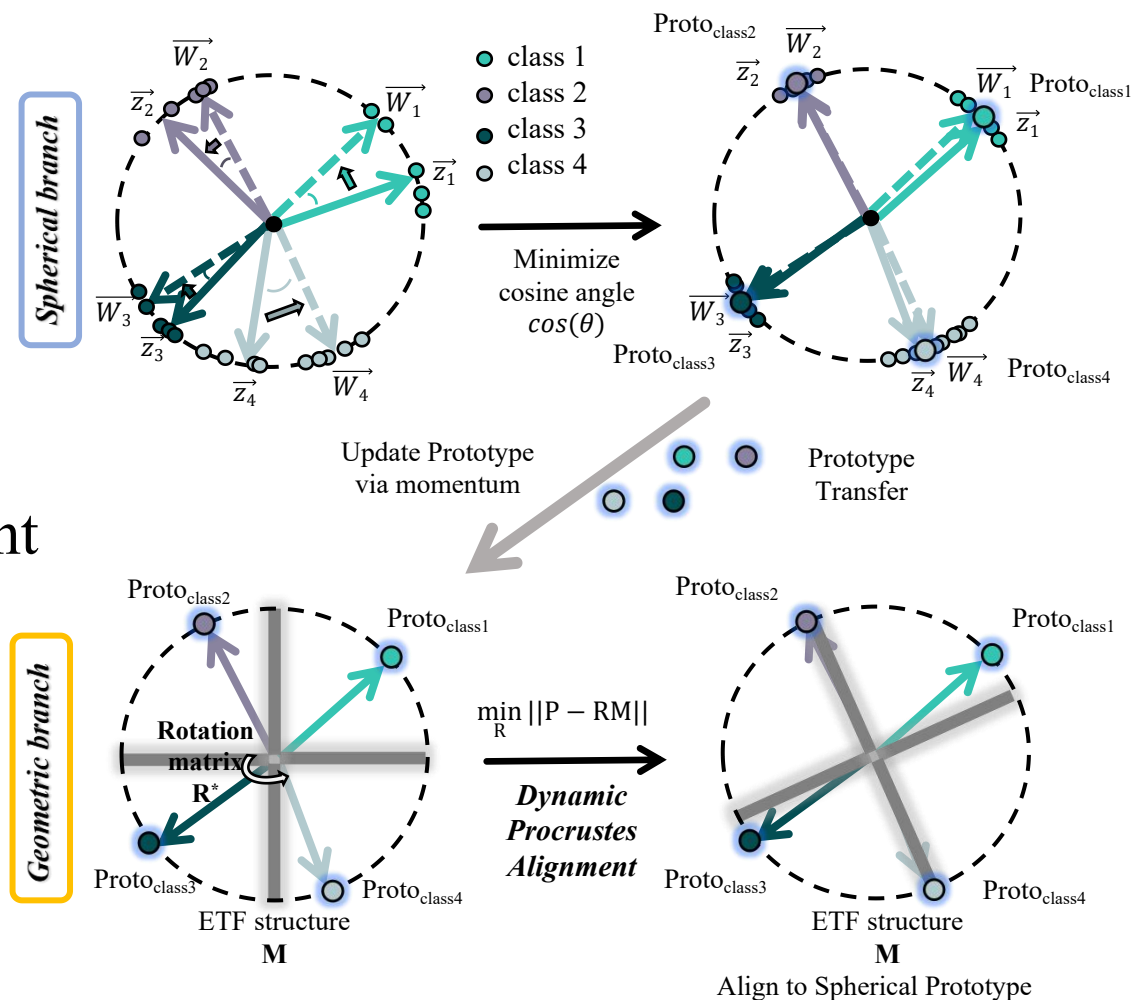
Momentum prototype update for stability

Geometric Branch – Dynamic Procrustes Alignment

- Fixed Simplex ETF M (maximally separated)
- Solve Orthogonal Procrustes:
- $\min_R \|P - RM\|_F \rightarrow R^* = UV^T$ via SVD

Avoid gradient jitter

\rightarrow stable geometry even on small/noisy data





Spherical Procrustes Alignment – Training & Inference

$$\text{Joint loss: } \mathcal{L}_{total} = \gamma(\mathcal{L}_{CE}^s + \mathcal{L}_{CE}^g) + \lambda\mathcal{L}_{KD}$$

$$\text{Inference: } P_{final} = \omega \cdot \sigma(\text{logits}_s) + (1 - \omega) \cdot \sigma(\text{logits}_g)$$

No extra inference cost



Main Results – ICBHI & CirCor & Calibration Methods

ICBHI (BEATs): AS \uparrow 65.27%, ECE \downarrow from 23.36% to **6.05%**

Backbone	Method	Classification Performance			Reliability
		Se (%) \uparrow	Sp (%) \uparrow	AS (%) \uparrow	ECE (%) \downarrow
CNN6	+ CE Loss (Kong et al., 2020b)	35.56 \pm 4.99	81.75 \pm 4.33	58.66 \pm 0.54	23.53 \pm 5.01
	+ Focal Loss (Lin et al., 2017)	35.11 \pm 5.36	79.18 \pm 5.01	57.14 \pm 0.58	17.51 \pm 2.24
	+ LDAM-DRW (Cao et al., 2019)	33.37 \pm 3.96	80.49 \pm 4.13	56.93 \pm 0.85	17.47 \pm 3.85
	+ SPA (Ours)	35.24 \pm 2.92	82.96 \pm 4.82	59.10 \pm 1.24	7.51 \pm 0.72
AST	+ CE Loss (Gong et al., 2021)	42.58 \pm 2.71	79.52 \pm 2.07	61.05 \pm 0.54	21.59 \pm 6.99
	FBS (Fraihy et al., 2025)	43.22 \pm 2.44	84.19 \pm 3.08	64.28 \pm 1.09	19.82 \pm 5.13
	FBS + SPA (Ours)	44.81 \pm 2.65	84.92 \pm 3.19	65.01 \pm 2.77	5.49 \pm 1.13
BEATs	+ CE Loss (Chen et al., 2023)	50.60 \pm 1.79	77.20 \pm 3.22	63.90 \pm 1.15	28.51 \pm 4.91
	+ SPA (Ours)	49.77 \pm 2.17	79.06 \pm 2.81	64.42 \pm 0.78	4.44 \pm 1.02
	PAFA (Jeong & Kim, 2025)	48.72 \pm 3.75	80.19 \pm 4.07	64.45 \pm 0.52	23.36 \pm 7.83
	PAFA + SPA (Ours)	49.34 \pm 4.12	81.21 \pm 4.48	65.27 \pm 0.57	6.05 \pm 3.63

CirCor (M2D): Wacc \uparrow 84.23%, ECE \downarrow to **4.63%**

Backbone	Method	Classification		Reliability
		W _{acc} (%) \uparrow	UAR (%) \uparrow	ECE (%) \downarrow
<i>Previous Studies</i>				
Wav2vec	Panah et al. (Panah et al., 2023)	80.0	70.0	-
HSM	CUED Acoustics (McDonald et al., 2022)	80.0	68.0	-
<i>Pre-trained extractors</i>				
CNN14	+ CE Loss (Kong et al., 2020a)	57.47 \pm 3.25	53.63 \pm 2.89	11.24 \pm 2.14
	+ SPA (Ours)	58.32 \pm 2.47	54.47 \pm 2.18	6.08 \pm 1.24
BYOL-A	+ CE Loss (Niizumi et al., 2021)	54.34 \pm 3.71	54.81 \pm 2.67	12.13 \pm 2.53
	+ SPA (Ours)	55.73 \pm 2.85	55.52 \pm 2.36	6.42 \pm 1.47
AST	+ CE Loss (Gong et al., 2021)	64.12 \pm 2.38	66.38 \pm 2.95	10.35 \pm 1.87
	+ SPA (Ours)	65.24 \pm 1.96	66.91 \pm 2.07	5.58 \pm 1.12
M2D	+ CE Loss (Niizumi et al., 2024)	83.51 \pm 1.92	72.14 \pm 2.14	9.17 \pm 1.65
	+ SPA (Ours)	84.23 \pm 1.53	72.96 \pm 1.78	4.63 \pm 0.89

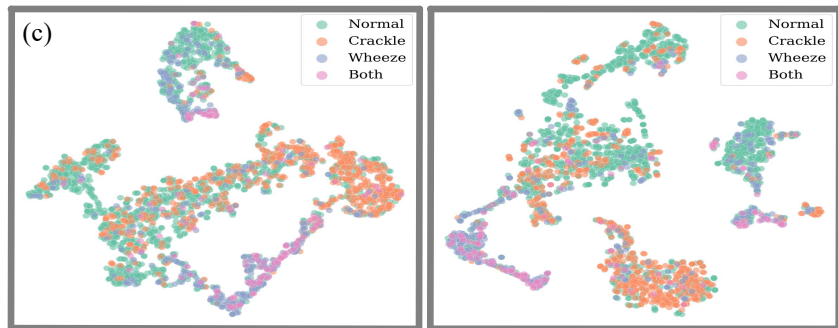
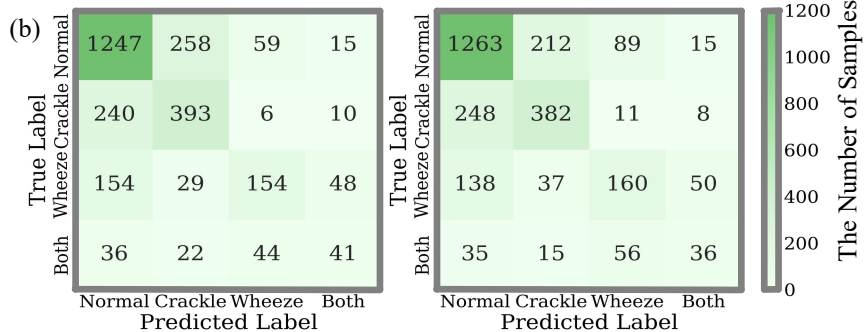
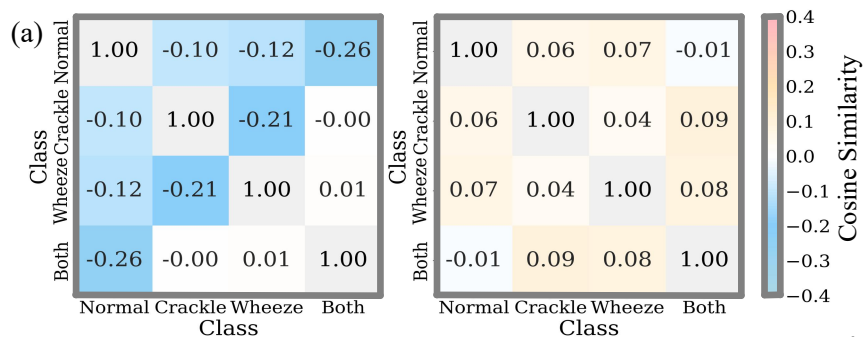
Method	AS (%) \uparrow	ECE (%) \downarrow	Inf. Time (ms)
CE Loss	63.90 \pm 1.15	28.51 \pm 4.91	4.1
Label Smoothing	62.58 \pm 0.59	33.91 \pm 4.41	4.1
Mixup	62.16 \pm 0.66	8.67 \pm 6.29	4.1
<i>Post-hoc & Sampling Methods</i>			
Temperature Scaling	63.90 \pm 1.15	21.99 \pm 6.97	4.1
MC Dropout (5 passes)	63.82 \pm 0.89	25.90 \pm 4.18	19.1
Deep Ensembles (5 models)	64.40 \pm 0.93	22.29 \pm 4.89	20.5
SPA (Ours)	64.42 \pm 0.78	4.44 \pm 1.02	4.2

Outperforms Label Smoothing, Mixup, Temperature Scaling, MC Dropout, and Deep Ensembles

without consuming extra inference

Visualization

Comparison on ICBHI



(a) Cosine similarity heatmap

SPA *enforces near-zero inter-class correlation*

(b) Confusion matrix

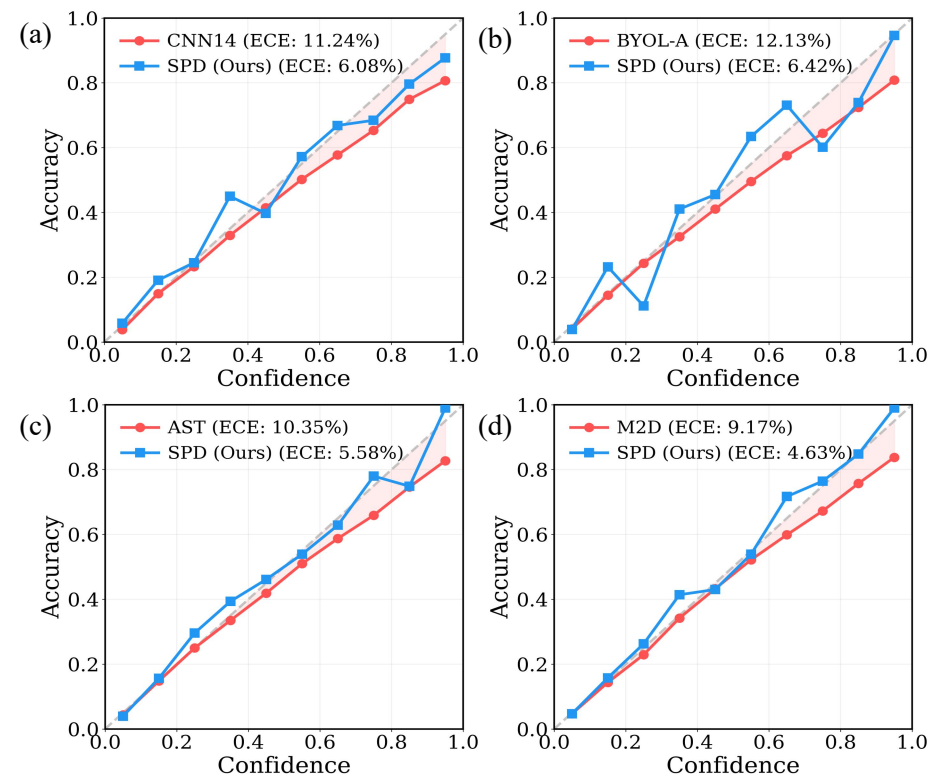
SPA *improves minority recall without hurting majority*

(c) t-SNE projection

SPA *yields compact, angularly separated clusters*

Reliability diagram on CirCor

SPA (blue) aligns near diagonal vs. baseline (red)





Conclusion

Norm bias is the geometric root of overconfidence in medical audio

SPA: spherical constraints plus dynamic Procrustes alignment

SOTA accuracy and calibration, no extra inference cost

Future Work

Robust prototype estimation, OOD detection