

# Alethia: a Foundational Encoder for Voice Deepfakes

Yi Zhu<sup>1</sup>, Brahmi Dwivedi<sup>2</sup>, Jayaram Raghuram, Surya Koppiseti

<sup>1,2</sup>Work conducted at Reality Defender

<sup>1</sup>Correspondence: yzhu00121@gmail.com



## Background

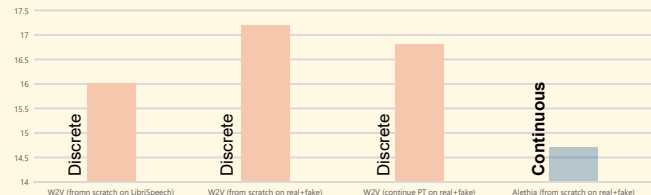
01 **Diminishing returns from post-training** on deepfake detection / localization tasks

02 Can we have a **scalable pre-training recipe** for deepfakes?

03 What are the desired representation **properties** and what **data and pretraining task** should be used?

☆ Type of tokens to predict: **Continuous vs Discrete**

Pre-training with discrete targets led to degraded performance (higher EER), we hypothesize this is due to loss of acoustic details



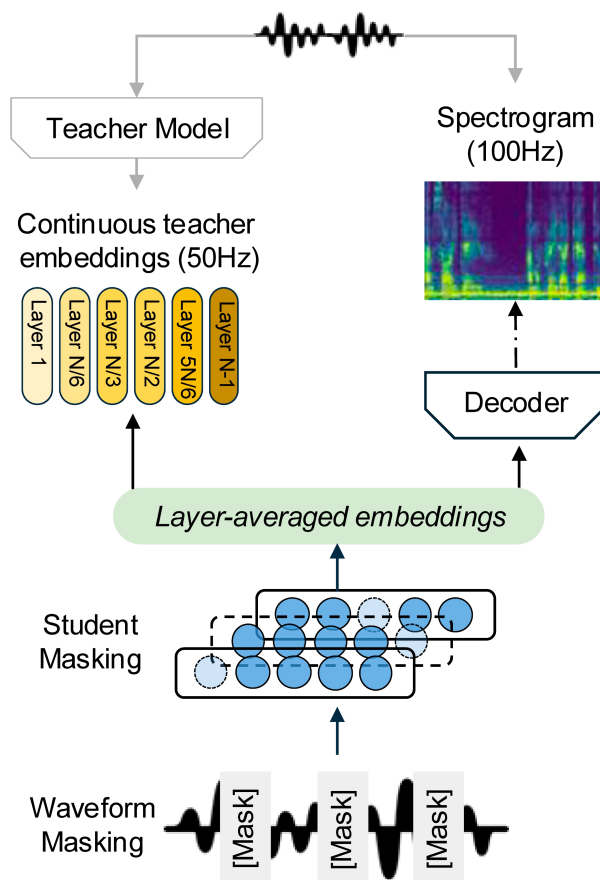
☆ A strong prior for masked latent prediction: **Bottleneck vs Layer-by-layer**

$$\min_{\theta} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\mathcal{L}_{MEP}(\{f_{\text{tea}}^m(\mathbf{x})\}_{m \in \mathcal{M}}, \Phi(f_{\text{stu}}(\tilde{\mathbf{x}}; \theta)))]$$

weighted cosine+L1 distance

M layers of teacher latent from **unmasked** speech

Φ: Layer pooling -> projection in M layers



☆ Reconstruct **Latent and Spectrogram** Simultaneously

Hypothesis: Acoustic details learned through generation are deepfake "giveaways"

Real and imaginary velocities

Bottleneck latent

$$[\hat{\mathbf{v}}_{\text{real}}, \hat{\mathbf{v}}_{\text{imag}}] = g_{\psi}(\mathbf{x}_t, t, \mathbf{z}_{\text{cond}})$$

$$\mathcal{L}_{FM} = \mathbb{E}_{t, \mathbf{x}_0, \mathbf{z}} \left[ \frac{\mathcal{L}_{\text{real}} + \mathcal{L}_{\text{imag}}}{\sigma_{\text{eps}}^2} \right]$$

⚙ Experimental Setup

Baselines: Wav2vec2.0 (0.3b, 1b), WavLM, HuBERT

Pretraining data: 30k hours of real+fake speech (19k after quality control)

Benchmark: 56 datasets, 5 tasks, 10+ languages, 100+ generative models, 10+ perturbations

