



Overview

- We address **unsupervised speech enhancement (SE)**.
- A **weighted-variance variational autoencoder (VAE)** model is proposed as the speech generative model.
- A **Gamma prior** distribution is imposed on the weights, leading to a **Student's t-distribution** for time-frequency elements.
- **Efficient training and speech enhancement** algorithms are developed.
- Experimental results demonstrate the **effectiveness** and **robustness** of the proposed approach compared to the standard unweighted variance model.

Unsupervised speech enhancement



Separate the speech and noise signals *without* training on noise.

Short-time Fourier transform (STFT) domain: $\mathbf{x} = \mathbf{s} + \mathbf{b}$

- \mathbf{s} → **clean speech signal** with prior $p_\theta(\mathbf{s})$
- \mathbf{b} → **noise signal** with prior $p_\psi(\mathbf{b})$

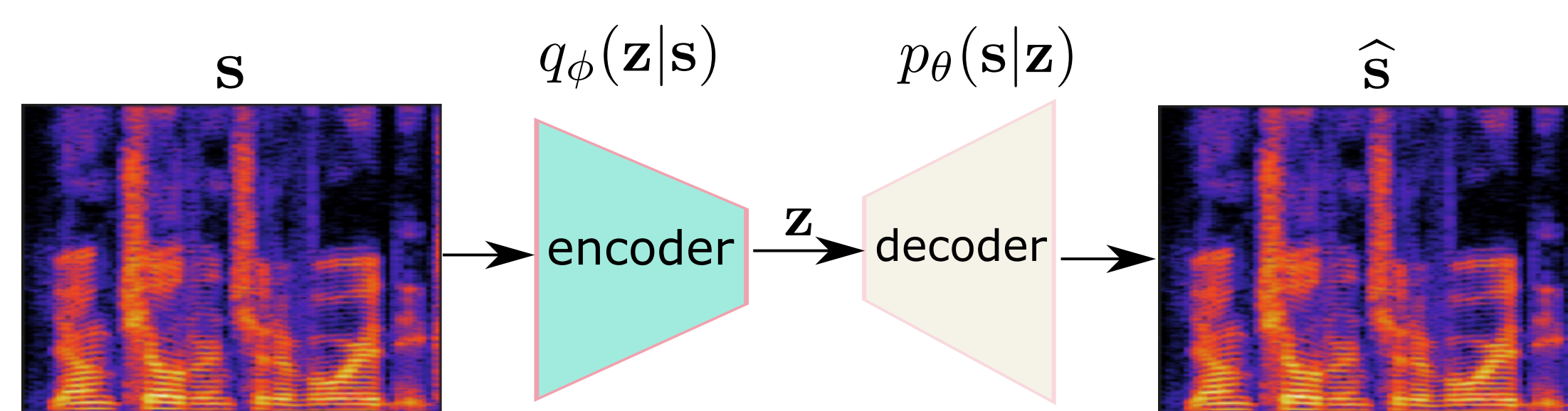
Training: Learn a parametric prior $p_\theta(\mathbf{s})$

Testing: Estimate \mathbf{s} using $p_\psi(\mathbf{s}|\mathbf{x}) \propto p_\psi(\mathbf{x}|\mathbf{s}) \times p_\theta(\mathbf{s})$

Training: learning speech prior

VAE-based speech generative model for $\mathbf{s} = \{\mathbf{s}_1, \dots, \mathbf{s}_T\}$ [1]:

$$p_\theta(\mathbf{s}) = \int p_\theta(\mathbf{s}|\mathbf{z})p(\mathbf{z})d\mathbf{z}, \quad p(\mathbf{z}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$



$$\begin{cases} p_\theta(\mathbf{s}_t|\mathbf{z}_t) = \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_\theta^2(\mathbf{z}_t))) \\ p(\mathbf{z}_t) = \mathcal{N}(\mathbf{0}, \mathbf{I}) \end{cases}$$

▷ Learn encoder-decoder parameters over *clean* speech data:

$$\mathcal{L}(\Phi; \mathbf{s}) = \mathbb{E}_{q_\psi(\mathbf{z}|\mathbf{s})} \{\log p_\theta(\mathbf{s}|\mathbf{z})\} - \mathcal{D}_{\text{KL}}(q_\psi(\mathbf{z}|\mathbf{s})||p(\mathbf{z}))$$

Testing: speech enhancement

Non-negative matrix factorization (NMF)-based noise model:

$$p_\psi(\mathbf{b}) \sim \mathcal{N}_c(\mathbf{0}, \text{diag}(\text{vec}(\mathbf{WH}))), \quad \psi = \{\mathbf{W}, \mathbf{H}\}$$

Parameter inference: Expectation-maximization (EM)

- **E-step:** compute posterior $p_\psi(\mathbf{z}|\mathbf{x})$ **Intractable** to compute! → Sample from it.
- **M-step:** update parameters:

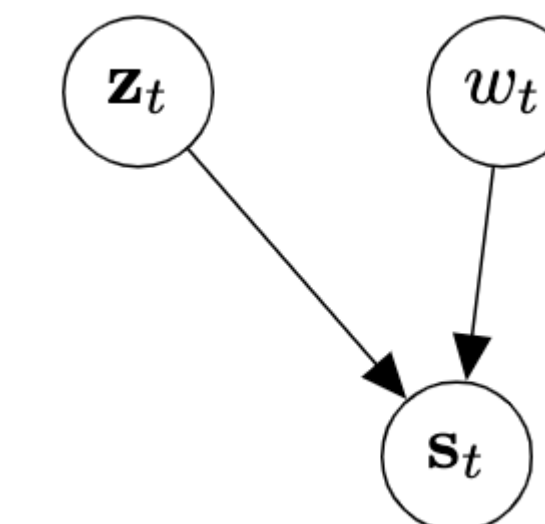
$$\max_{\psi} \mathbb{E}_{p_\psi(\mathbf{z}|\mathbf{x})} \{\log p_\psi(\mathbf{x}|\mathbf{z})\}$$

☞ *multiplicative update rules*

Proposed framework: StVAE model

- **Weighting the variance:** introduce scalar weights per variance components

$$\begin{cases} p_\theta(\mathbf{s}_t|\mathbf{z}_t, w_t) = \mathcal{N}_c(\mathbf{0}, \text{diag}(\sigma_\theta^2(\mathbf{z}_t)/w_t)) \\ p(\mathbf{z}_t) = \mathcal{N}(\mathbf{0}, \mathbf{I}) \\ p(w_t) = \mathcal{G}(w_t; \alpha, \beta) \end{cases}$$



- **Gamma prior over the weights:**

$$\mathcal{G}(w; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} w^{\alpha-1} \exp(-\beta w)$$

☞ An infinite *mixture* of Gaussian distributions: $p_\theta(\mathbf{s}_t|\mathbf{z}_t) = \int p_\theta(\mathbf{s}_t|\mathbf{z}_t, w_t)p(w_t)dw_t$

Training phase of VAE:

- **Parameters to learn:** $\tilde{\Phi} = \{\theta, \psi, \alpha, \beta\}$
- **Posterior distribution:** $p_\theta(\mathbf{z}_t, w_t|\mathbf{s}_t) = p_\theta(w_t|\mathbf{s}_t, \mathbf{z}_t) \cdot p_\theta(\mathbf{z}_t|\mathbf{s}_t)$

The first term writes $p_\theta(w_t|\mathbf{s}_t, \mathbf{z}_t) \propto p_\theta(\mathbf{s}_t|\mathbf{z}_t, w_t) \cdot p(w_t) = \mathcal{G}(\alpha'_t, \beta'_t)$, where:

$$\begin{cases} \alpha'_t = \alpha + F \\ \beta'_t = \beta + \sum_f \frac{|s_{f,t}|^2}{\sigma_{\theta,f}^2(\mathbf{z}_t)} \end{cases}$$

The intractable posterior $p_\theta(\mathbf{z}_t|\mathbf{s}_t)$ is approximated by a variational distribution: $p_\theta(\mathbf{z}_t|\mathbf{s}_t) \approx q_\psi(\mathbf{z}_t|\mathbf{s}_t)$.

$$p_\theta(\mathbf{z}, \mathbf{w}|\mathbf{s}) \approx q_\psi(\mathbf{z}, \mathbf{w}) = p_\theta(\mathbf{w}|\mathbf{s}, \mathbf{z})q_\psi(\mathbf{z}|\mathbf{s})$$

- **Training objective:**

$$\log p_\theta(\mathbf{s}) \geq \mathbb{E}_{q_\psi(\mathbf{z}, \mathbf{w})} \left\{ \log \frac{p_\theta(\mathbf{s}, \mathbf{z}, \mathbf{w})}{q_\psi(\mathbf{z}, \mathbf{w})} \right\} \triangleq \mathcal{L}(\tilde{\Phi}; \mathbf{s})$$

Speech enhancement phase:

E-step: Posterior $p_\psi(\mathbf{z}_t, w_t|\mathbf{x}_t)$ is intractable. Instead, we simply find the modes:

$$\mathbf{z}_t^*, w_t^* = \arg \max_{\mathbf{z}_t, w_t} \log p_\psi(\mathbf{z}_t, w_t|\mathbf{x}_t)$$

☞ **First-order optimization**

M-step: Update the NMF matrices using

$$\max_{\mathbf{W}, \mathbf{H}} \sum_t \log p_\psi(\mathbf{x}_t|\mathbf{z}_t^*, w_t^*)$$

☞ **Multiplicative update rules**

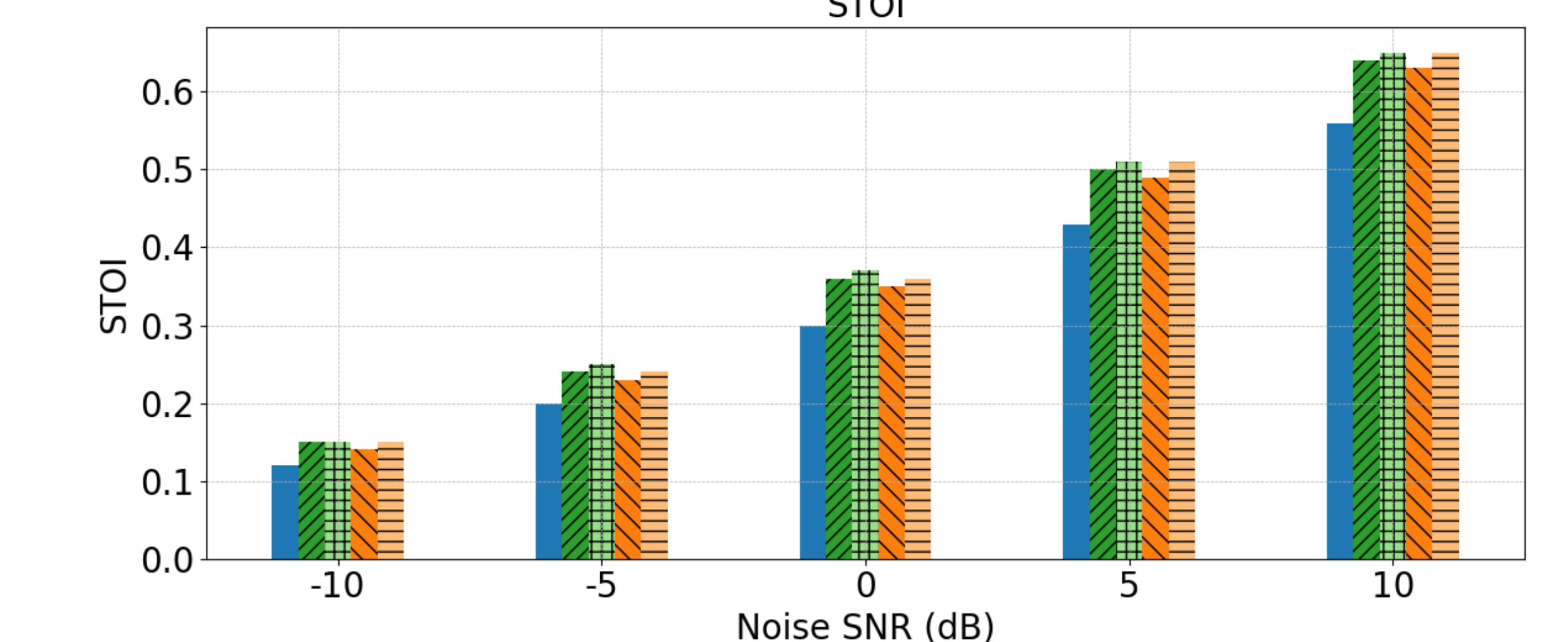
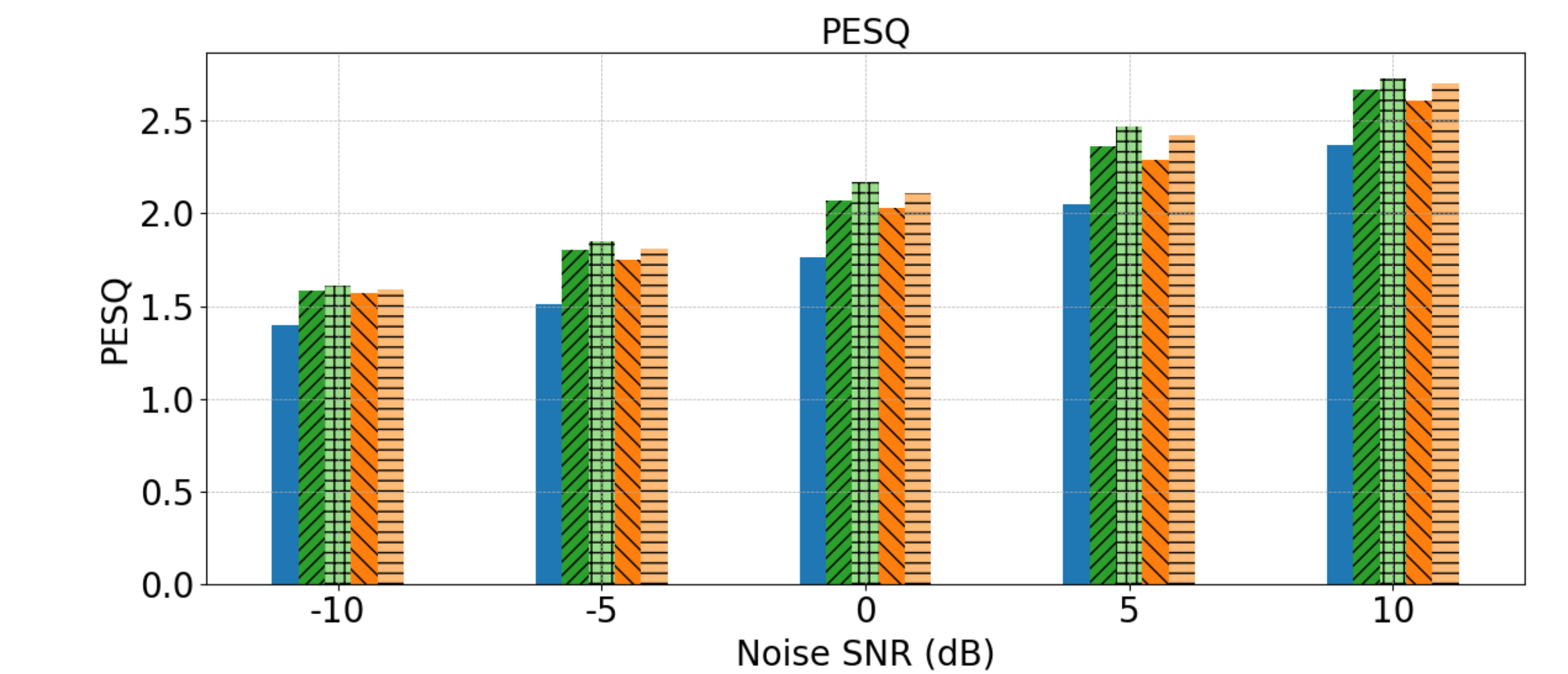
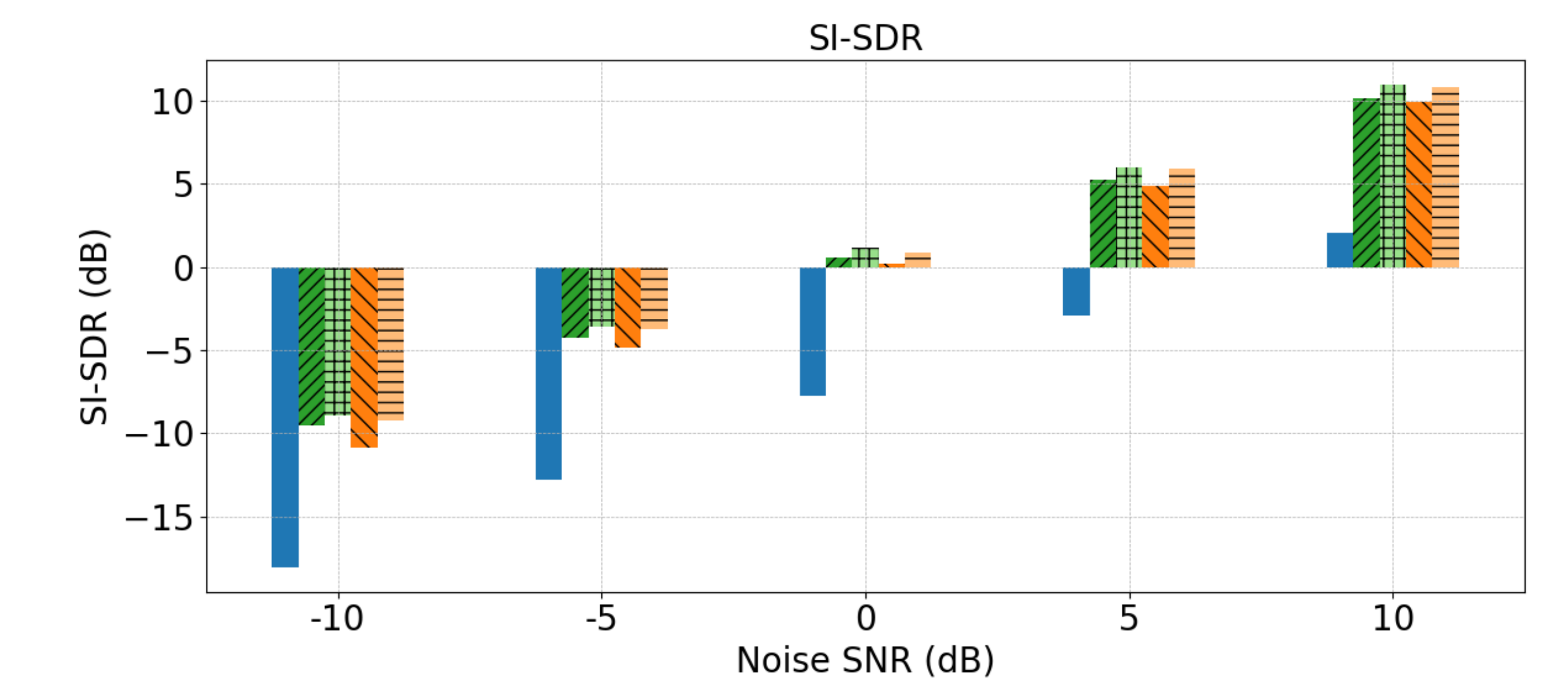
Speech estimation: Posterior mean $\hat{\mathbf{s}}_t = \mathbb{E}_{p_\psi(\mathbf{s}_t|\mathbf{x}_t)} \{\mathbf{s}_t\}, \forall t$

$$\hat{\mathbf{s}}_t = \mathbb{E}_{p_\psi(\mathbf{z}_t^*, w_t^*|\mathbf{x}_t)} \left\{ \mathbb{E}_{p_\psi(\mathbf{s}_t|\mathbf{x}_t, \mathbf{z}_t^*, w_t^*)} \{\mathbf{s}_t\} \right\} \approx \frac{(w_t^*)^{-1} \sigma_\theta^2(\mathbf{z}_t^*)}{(w_t^*)^{-1} \sigma_\theta^2(\mathbf{z}_t^*) + \mathbf{W}^* \mathbf{h}_t^*} \odot \mathbf{x}_t$$

Experiments

- **Datasets:** TCD-TIMIT (training & evaluation)
 - *Clean* setup: training on clean speech (~ 8 hrs)
 - *Outlier* setup: training on {clean, noise} data (~ 9.6 hrs)
- **Parameters:** STFT with 64 ms-long (1024 samples) sine window, 75% overlap ($F = 512$). $K = 100$ EM iterations, 10 iterations of posterior sampling. Latent dimension $L = 32$.
- **Baseline:** VAE [1] (single-layer, 128 nodes, encoder-decoder).

Legend: Input (blue), VAE (clean) (green), StVAE (clean) (dark green), VAE (outlier) (orange), StVAE (outlier) (dark orange)



▷ StVAE surpasses VAE, especially at higher SNRs, highlighting its effective weighted variance Gaussian distribution.

▷ StVAE shows robust performance on noise-contaminated data.

▷ Training StVAE on noise-contaminated data outperforms the VAE trained on clean data, demonstrating StVAE's superior robustness.

Reference

[1] S. Leglaive, et al., "A variance modeling framework based on variational autoencoders for speech enhancement," IEEE MLSP, September 2018.