

Kuan-Lin Chen Advisor: Bhaskar D. Rao

Department of Electrical and Computer Engineering, University of California, San Diego

kuc029@ucsd.edu

## 1 Complexity of deep ReLU networks

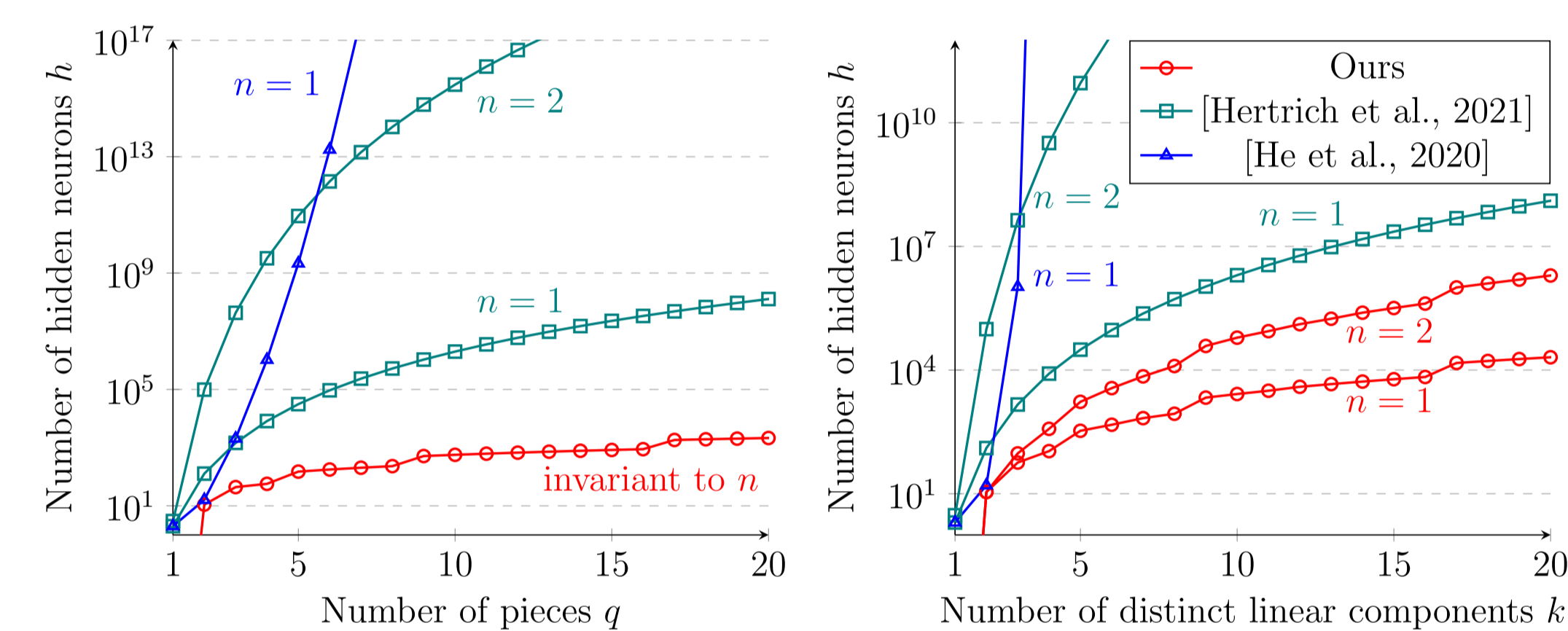
**Theorem 1** (Theorem 1 of [1]). Any continuous piecewise linear (CPWL) function  $p: \mathbb{R}^n \rightarrow \mathbb{R}$  with  $q$  pieces can be represented by a ReLU network whose number of layers  $l$ , maximum width  $w$ , and number of hidden neurons  $h$  satisfy

$$l \leq 2 \lceil \log_2 q \rceil + 1, \quad (1)$$

$$w \leq \mathbb{I}[q > 1] \left\lceil \frac{3q}{2} \right\rceil, \quad (2)$$

$$h \leq (3 \cdot 2^{\lceil \log_2 q \rceil} + 2 \lceil \log_2 q \rceil - 3)q + 3 \cdot 2^{\lceil \log_2 q \rceil} - 2 \lceil \log_2 q \rceil - 3. \quad (3)$$

Furthermore, Algorithm 1 finds such a network in  $\text{poly}(n, q, L)$  time where  $L$  is the number of bits required to represent every entry of the rational matrix  $\mathbf{A}_i$  in the polyhedron representation  $\{\mathbf{x} \in \mathbb{R}^n | \mathbf{A}_i \mathbf{x} \leq \mathbf{b}_i\}$  of the piece  $\mathcal{X}_i$  for every  $i \in [q]$ .



**Figure 1:** Any CPWL function  $\mathbb{R}^n \rightarrow \mathbb{R}$  with  $q$  pieces or  $k$  distinct linear components can be exactly represented by a ReLU network with at most  $h$  hidden neurons. The upper bounds (red) in [1] are substantially tighter than existing bounds in the literature, showing that any CPWL function can be exactly realized by a ReLU network at a much lower cost.

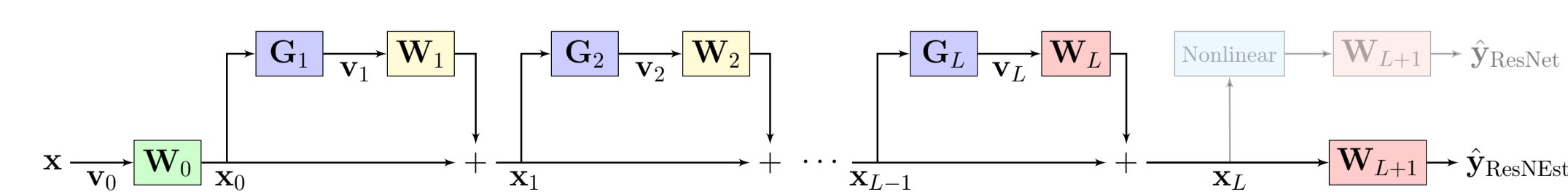
**Algorithm 1** Find a ReLU network that computes a given CPWL function

**Input:** A CPWL function  $p$  with pieces  $\{\mathcal{X}_i\}_{i \in [q]}$  of  $\mathbb{R}^n$ .

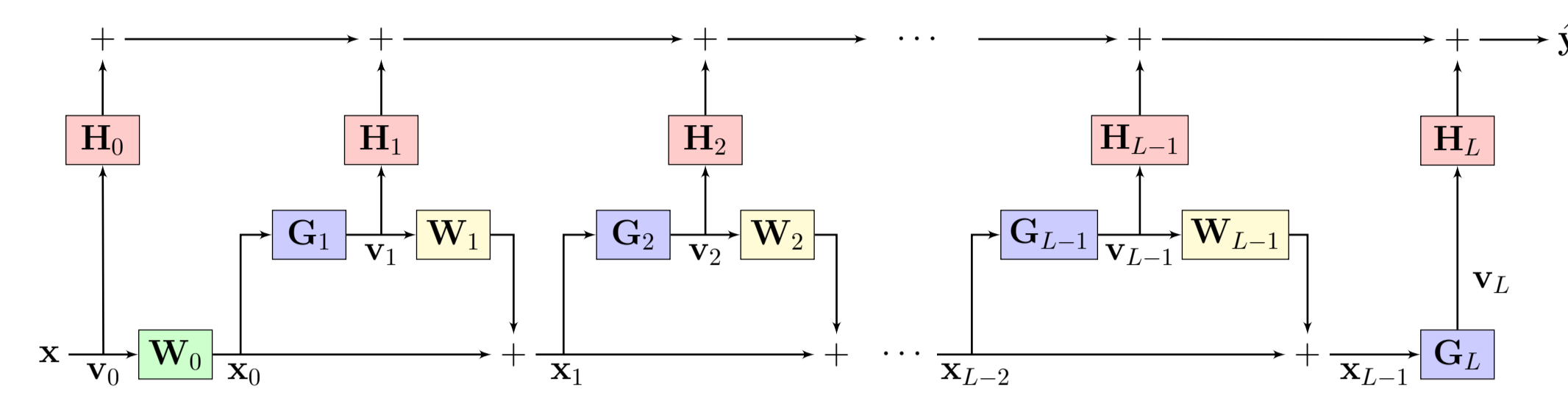
**Output:** A ReLU network  $g$  computing  $g(\mathbf{x}) = p(\mathbf{x}), \forall \mathbf{x} \in \mathbb{R}^n$ .

- 1:  $f_1, f_2, \dots, f_k \leftarrow$  Find all distinct linear components of  $p$
- 2: **for**  $i = 1, 2, \dots, q$  **do**
- 3:  $\mathcal{A}_i \leftarrow \emptyset$
- 4: **for**  $j = 1, 2, \dots, k$  **do**
- 5: **if**  $f_j(\mathbf{x}) \geq p(\mathbf{x}), \forall \mathbf{x} \in \mathcal{X}_i$  **then**
- 6:  $\mathcal{A}_i \leftarrow \mathcal{A}_i \cup \{j\}$
- 7: **end if**
- 8: **end for**
- 9:  $v_i \leftarrow$  A network representing the min-affine function of  $\{f_m\}_{m \in \mathcal{A}_i}$
- 10: **end for**
- 11:  $v \leftarrow$  Combine ReLU networks  $v_1, v_2, \dots, v_q$  in parallel
- 12:  $u \leftarrow$  A ReLU network computing the maximum of  $q$  elements
- 13:  $g \leftarrow$  A ReLU network computing the composition  $u \circ v$

## 2 Optimization of deep residual networks



**Figure 2:** A generic vector-valued ResNEst that has a chain of  $L$  residual blocks. Different from the standard ResNet architecture, our ResNEst architecture drops nonlinearities at  $\mathbf{x}_L$  so as to reveal a linear relationship between the output  $\hat{\mathbf{y}}_{\text{ResNEst}}$  and the features  $\mathbf{v}_0, \mathbf{v}_1, \dots, \mathbf{v}_L$ .



**Figure 3:** The proposed augmented ResNEst or A-ResNEst.

Because the ResNEst now reveals a linear relationship between the output and the features, we have:

$$\hat{\mathbf{y}}_{L\text{-ResNEst}}(\mathbf{x}) = \mathbf{W}_{L+1} \sum_{i=0}^L \mathbf{W}_i \mathbf{v}_i(\mathbf{x}), \quad (4)$$

$$\mathbf{v}_i(\mathbf{x}) = \mathbf{G}_i(\mathbf{x}_{i-1}; \boldsymbol{\theta}_i) = \mathbf{G}_i \left( \sum_{j=0}^{i-1} \mathbf{W}_j \mathbf{v}_j; \boldsymbol{\theta}_i \right). \quad (5)$$

We propose to utilize the basis function modeling point of view in the ResNEst and analyze the following ERM problem:

$$(\mathcal{P}_\phi) \min_{\mathbf{W}_L, \mathbf{W}_{L+1}} \mathcal{R}(\mathbf{W}_L, \mathbf{W}_{L+1}; \phi) \quad (6)$$

where  $\mathcal{R}(\mathbf{W}_L, \mathbf{W}_{L+1}; \phi) = \frac{1}{N} \sum_{n=1}^N \ell(\hat{\mathbf{y}}_{L\text{-ResNEst}}(\mathbf{x}^n), \mathbf{y}^n)$  for any fixed feature finding weights  $\phi$ . For A-ResNEst, we put

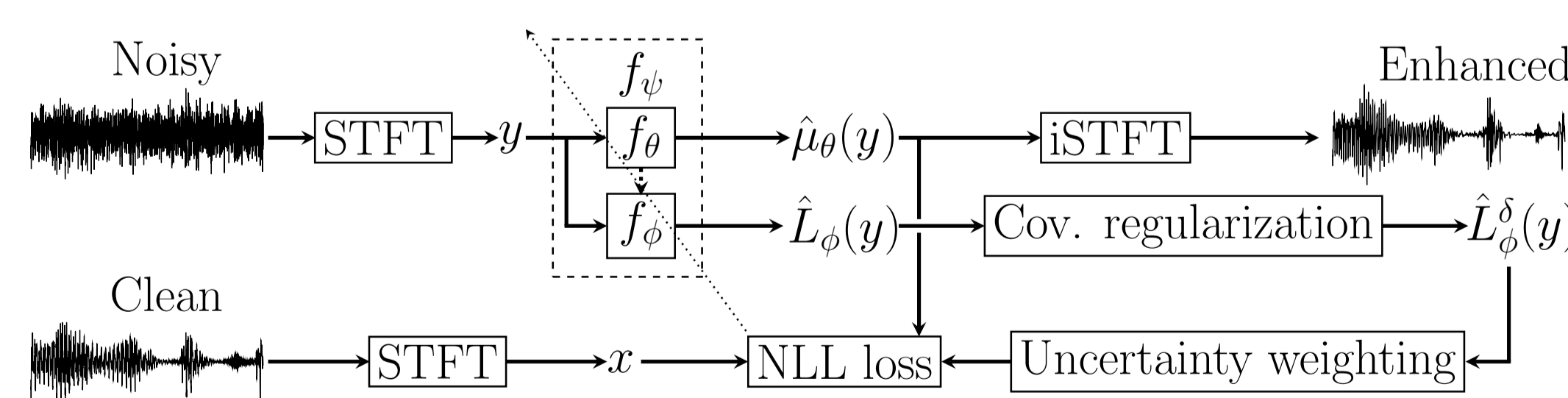
$$(\mathcal{P}_A) \min_{\mathbf{H}_0, \dots, \mathbf{H}_L} \mathcal{A}(\mathbf{H}_0, \dots, \mathbf{H}_L; \phi) \quad (7)$$

where  $\mathcal{A}(\mathbf{H}_0, \dots, \mathbf{H}_L; \phi) = \frac{1}{N} \sum_{n=1}^N \ell(\hat{\mathbf{y}}_{L\text{-A-ResNEst}}(\mathbf{x}^n), \mathbf{y}^n)$ .

- $M$  is the output dimension of  $\mathbf{W}_0$  (expansion factor).
- $N_\phi$  is the output dimension of the network.

**Theorem 2** (Theorem 1 of [2]). If the loss function  $\ell(\hat{\mathbf{y}}, \mathbf{y})$  is differentiable and convex in  $\hat{\mathbf{y}}$  for any  $\mathbf{y}$  and  $M \geq N_\phi$ , then the following two properties are true in  $(\mathcal{P}_\phi)$  under any  $\phi$  such that the linear inverse problem  $\mathbf{x}_{L-1} = \sum_{i=0}^{L-1} \mathbf{W}_i \mathbf{v}_i$  has a unique solution: (a) every critical point with full rank  $\mathbf{W}_{L+1}$  is a global minimizer; (b)  $\mathcal{R}(\mathbf{W}_L^*, \mathbf{W}_{L+1}^*; \phi) = \mathcal{A}(\mathbf{H}_0^*, \dots, \mathbf{H}_L^*; \phi)$  for every local minimizer  $(\mathbf{W}_L^*, \mathbf{W}_{L+1}^*)$  of  $(\mathcal{P}_\phi)$ .

## 3 Uncertainty in supervised speech enhancement



**Figure 4:** We augment a speech enhancement model  $f_\phi$  with a temporary sub-model  $L_\phi$  to estimate heteroscedastic uncertainty during training [3].

The problem of maximum likelihood is equivalent to minimizing the empirical risk using the multivariate Gaussian NLL loss

$$\ell_{x,y}^{\text{Full}}(\psi) = [x - \hat{\mu}_\theta(y)]^\top \hat{\Sigma}_\phi^{-1}(y) [x - \hat{\mu}_\theta(y)] + \log \det \hat{\Sigma}_\phi(y). \quad (8)$$

The number of elements in  $\hat{\Sigma}_\phi(y)$  is  $4T^2F^2$ , leading to exceedingly high training complexity. How can we reduce the complex-

ity and make the maximum likelihood tractable?

$$\ell_{x,y}^{\text{Diagonal}}(\psi) = \sum_{t,f} \sum_{k \in \{r,i\}} \left[ \frac{x_k^{t,f} - \hat{\mu}_{k;\theta}^{t,f}(y)}{\hat{\sigma}_{k;\phi}^{t,f}(y)} \right]^2 + 2 \log \hat{\sigma}_{k;\phi}^{t,f}(y). \quad (9)$$

$$\ell_{x,y}^{\text{Block}}(\psi) = \sum_{t,f} d_{\theta,x}^{t,f}(y)^\top \left[ \hat{\Sigma}_\phi^{t,f}(y) \right]^{-1} d_{\theta,x}^{t,f}(y) + \log t_{\phi}^{t,f}(y). \quad (10)$$

**Covariance regularization.** Let  $\delta > 0$  be the lower bound of the eigenvalues of the Cholesy factor of the covariance matrix.

$$\left[ \hat{L}_\phi^\delta(y) \right]_{mm} = \max \left\{ \left[ \hat{L}_\phi(y) \right]_{mm}, \delta \right\}. \quad (11)$$

**Uncertainty weighting.** Let  $\beta = 0.5$  and the loss function be a weighted average where the weight of a loss component depends on the minimum eigenvalue of the covariance matrix, i.e.,

$$\ell_{x,y}^{\beta\text{-Block}}(\psi) = \sum_{t,f} \lambda_{\min} \left[ \hat{\Sigma}_\phi^{t,f}(y) \right]^\beta z_{x,y}^{t,f}(\psi). \quad (12)$$

SNR (dB)	WB-PESQ			STOI (%)			SI-SDR (dB)			NOREQA-MOS		
	-5	0	5	-5	0	5	-5	0	5	-5	0	5
Unprocessed	1.11	1.15	1.24	69.5	77.8	85.2	-5.00	0.01	5.01	2.32	2.36	2.45
MAE	1.50	1.76	2.09	84.4	90.4	93.9	9.83	12.63	15.02	2.77	3.27	3.65
MSE	1.63	1.94	2.29	85.1	90.6	94.0	10.24	13.21	15.97	2.86	3.52	4.02
SI-SDR	1.71	2.04	2.42	86.5	91.5	94.6	<b>10.96</b>	<b>13.92</b>	<b>16.80</b>	3.05	3.65	4.20
NLL $\ell^{\text{Diagonal}}$	1.74	2.08	2.48	86.2	91.3	94.6	9.83	12.55	15.01	3.14	3.77	4.25
NLL $\ell^{\text{Block}}$	<b>1.75</b>	<b>2.10</b>	<b>2.50</b>	<b>86.7</b>	<b>91.8</b>	<b>94.9</b>	10.22	13.15	15.99	<b>3.23</b>	<b>3.89</b>	<b>4.35</b>

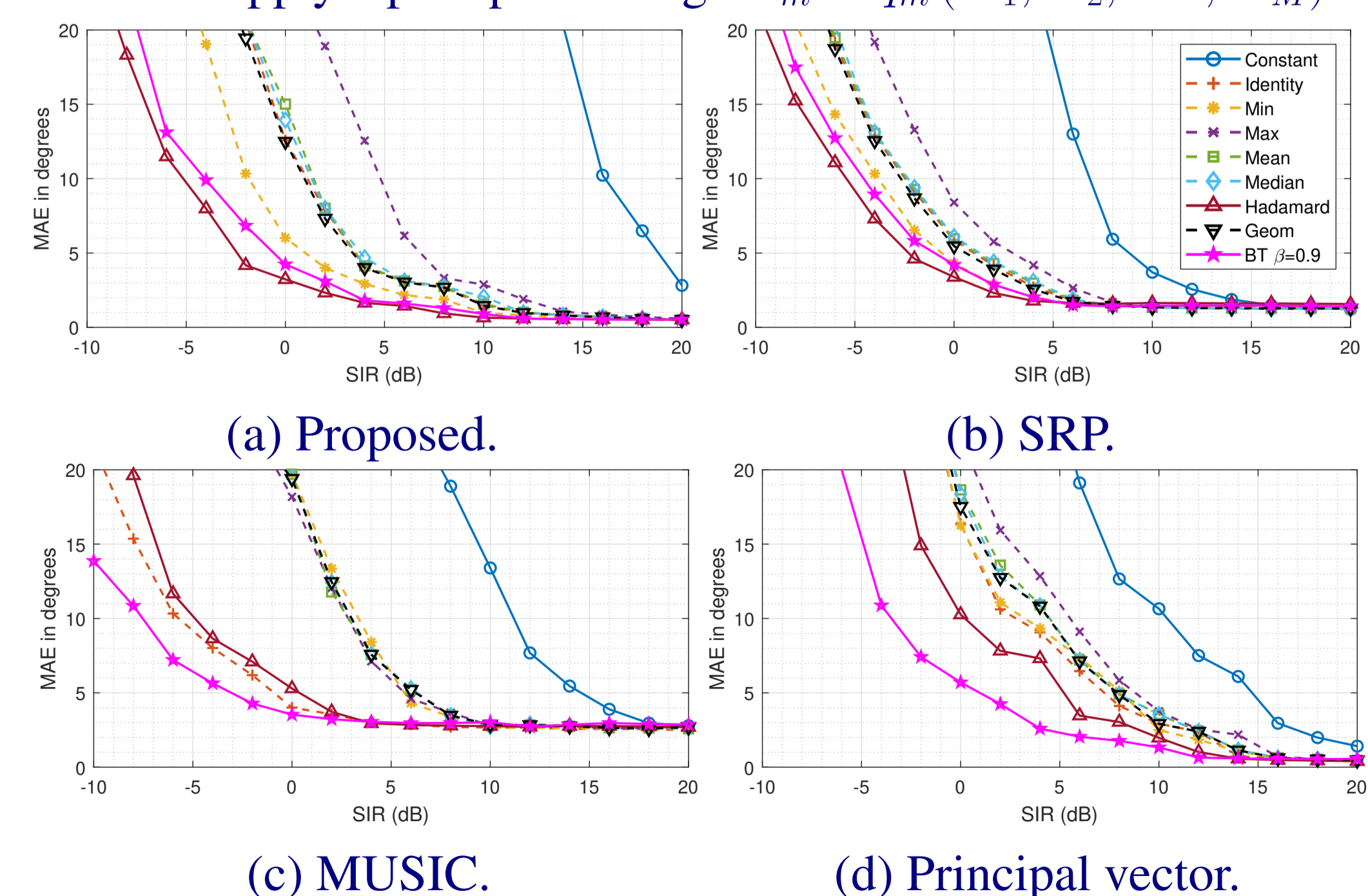
**Figure 5:** The NLL using a block diagonal covariance with suitable  $\delta$  and  $\beta$  outperforms the MAE, MSE, and SI-SDR. The DNS dataset is used. We adopt the GCRN as  $f_\phi$  for investigation.  $f_\phi$  is an additional decoder that takes the output of the in-between LSTM of the GCRN as input.

## 4 DNN based direction of arrival estimation

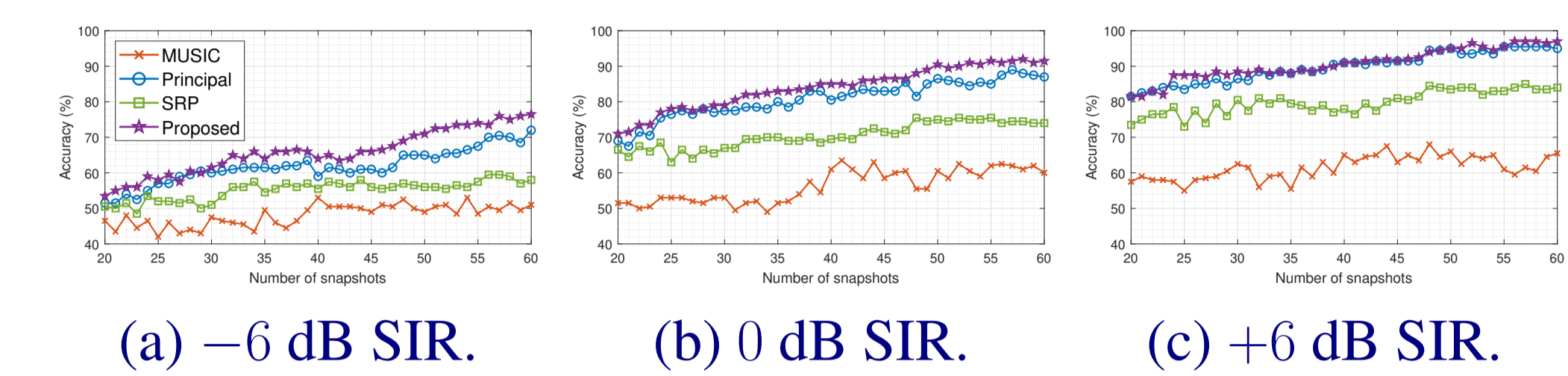
Let  $\tilde{\mathbf{y}}(t, f) = \mathbf{w}(t, f) \odot \mathbf{y}(t, f)$  be a filtered snapshot. We propose a criterion that normalizes the filtered snapshot [4].

$$\max_{\theta} \sum_f \mathbf{v}^H(\theta, f) \sum_t \frac{\tilde{\mathbf{y}}(t, f) \tilde{\mathbf{y}}^H(t, f)}{\|\mathbf{y}(t, f)\|_2^2} \mathbf{v}(\theta, f). \quad (13)$$

To find the time-frequency weights  $\mathbf{w}(t, f)$ , we first use a U-Net (0.67M params) to predict the ideal ratio mask  $\mathbf{G}_m$  on each sensor and then apply a post-processing  $\mathbf{W}_m = q_m(\mathbf{G}_1, \mathbf{G}_2, \dots, \mathbf{G}_M)$ .



**Figure 6:** MAE in degrees vs. SIR.  $\text{RT}_{60} = 0.3\text{s}$  and  $\text{SNR} = 20\text{ dB}$ .



**Figure 7:** Accuracy vs. number of snapshots.  $\text{RT}_{60} = 0.3\text{s}$  and  $\text{SNR} = 20\text{ dB}$ .

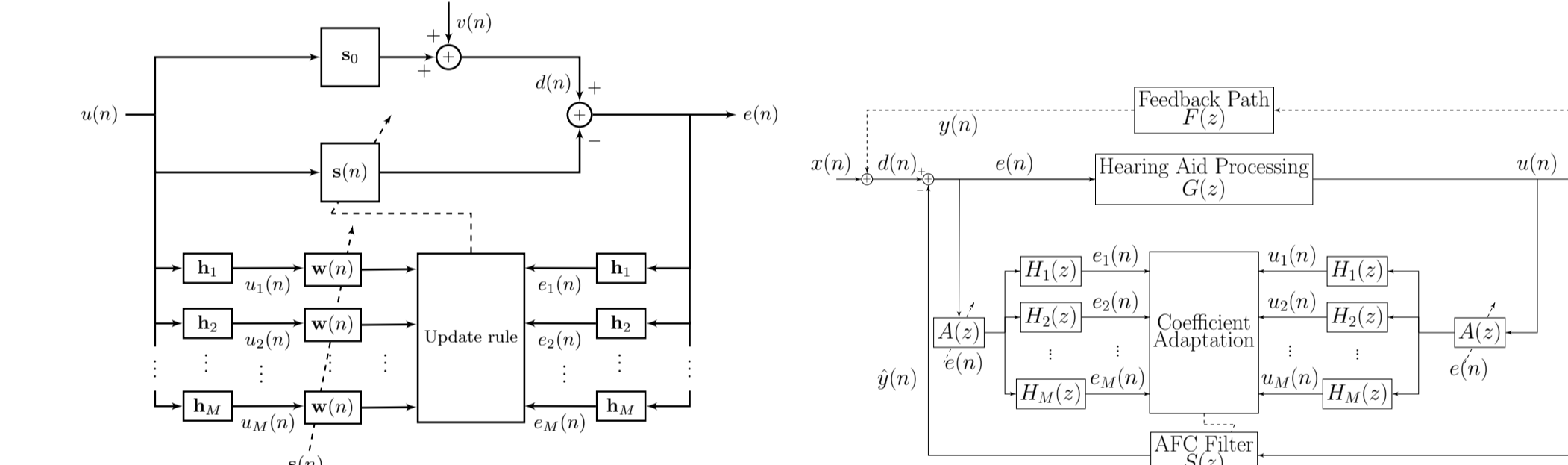
## 5 Adaptive filters and feedback cancellation

We minimize the sum of the squared error in each subband with a sparsity penalty term. We propose the following cost function:

$$J(\mathbf{s}) = \sum_{i=1}^M |e_i(n)|^2 + \tau \|\mathbf{s}\|_{\mathbf{W}^{-1}(n)}^2 \quad (14)$$

where  $e_i(n) = \mathbf{h}_i^T \mathbf{e}(n) = \mathbf{h}_i^T [\mathbf{d}(n) - \mathbf{U}^T(n)\mathbf{s}]$  is the  $i$ -th subband error and  $\mathbf{s} \in \mathbb{R}^L$  is the coefficients of the adaptive filter, leading to the generalized proportionate-type normalized subband adaptive filter (GpNSAF) [5]:  $\mathbf{s}(n+1) = \mathbf{s}(n) + \mu \mathbf{g}(n)$  where

$$\mathbf{g}(n) = \mathbf{W}(n) \mathbf{U}_b(n) \left[ \delta \mathbf{I}_M + \mathbf{U}_b^T(n) \mathbf{W}(n) \mathbf{U}_b(n) \right]^{-1} \mathbf{e}_b(n). \quad (15)$$



(a) System identification. (b) Feedback cancellation.

**Figure 8:** The GpNSAF and the feedback cancellation framework [6].

We put  $\mathbf{W}(n) = \text{diag}\{w_1(n), \dots, w_L(n)\}$  where

$$w_i(n) = (|s_i(n)| + c)^{-2-p}, \quad i = 1, 2, \dots, L, \quad (16)$$

$p \in [1.0, 2.0]$ ,  $c > 0$  for promoting different degrees of sparsity.

	$M = 1$	$M > 1, \mathbf{H} \neq \mathbf{I}$	$M > 1, \mathbf{H} = \mathbf{I}$
$p = 2$	NLMS	NSAF	APA
$2 > p > 0$	PtNLMS	PtNSAF	PtAPA

**Table 1:** Special cases of GpNSAF (or sparsity-promoting NSAF).

## References

- [1] K.-L. Chen, H. Garudadri, and B. D. Rao, "Improved bounds on neural complexity for representing piecewise linear functions," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [2] K.-L. Chen, C.-H. Lee, H. Garudadri, and B. D. Rao, "ResNEsts and DenseNEsts: Block-based DNN models with improved representation guarantees," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- [3] K.-L. Chen, D. D. E. Wong, K. Tan, B. Xu, A. Kumar, and V. K. Ithapu, "Leveraging heteroscedastic uncertainty in learning complex spectral mapping for single-channel speech enhancement," in *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023.
- [4] K.-L. Chen, C.-H. Lee, B. D. Rao, and H. Garudadri, "A DNN based normalized time-frequency weighted criterion for robust wideband DoA estimation," in *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023.
- [5] K.-L. Chen, C.-H. Lee, B. D. Rao, and H. Garudadri, "A generalized proportionate-type normalized subband adaptive filter," in *Asilomar Conference on Signals, Systems, and Computers*. IEEE, 2019.
- [6] K.-L. Chen, C.-H. Lee, B. D. Rao, and H. Garudadri, "Jointly leveraging decorrelation and sparsity for improved feedback cancellation in hearing aids," in *European Signal Processing Conference (EUSIPCO)*, 2020.