# DEEP LEARNING
# FOR MONAURAL SPEECH SEPARATION

Po-Sen Huang, Minje Kim,

Mark Hasegawa-Johnson, Paris Smaragdis

ILLINOIS

# Motivation

- Source separation is important for several real-world applications
  - Monaural speech separation is more difficult
- Previous approaches – linear models
  - Non-negative matrix factorization (NMF), probabilistic latent semantic indexing (PLSI)
  - Similar to a one layer linear network with non-negative weights and coefficients
- Representation
  - NMF based models – spectral representation
  - Deep learning models – learning optimal representation
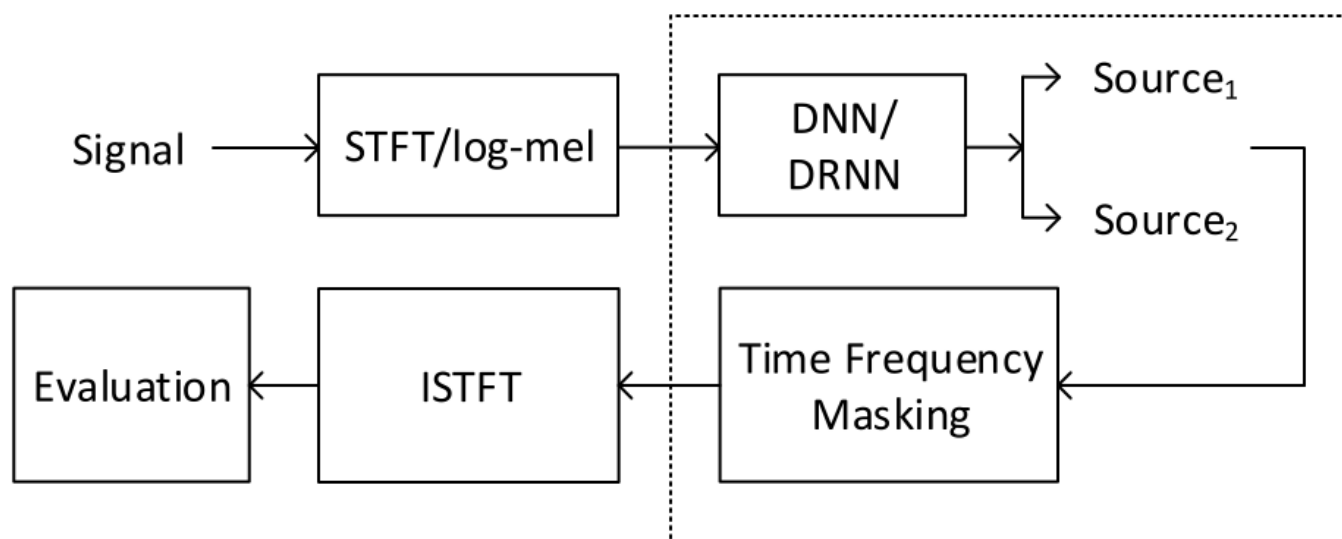- ➢ Explore deep learning models

# Overview

- Previous Work

- Proposed Framework

- Proposed methods

  - Model architecture

  - Joint time frequency masking

  - Discriminative training

- Experiment

- Demo

- Conclusion

# Previous Work – Deep Learning

- Two-stage framework for predicting ideal binary mask [Narayanan and Wang, 2013]
  - First stage: Train one DNN per output dimension
  - Second stage: train another one layer perceptron or SVM for refinement
  - ➢ Impractical for high dimensional output
- Robust ASR [Mass et al. 2012]
  - Given noisy speech, train DRNN to predict clean speech
  - ➢ Suboptimal in the source separation scenario to model only one source

# Proposed Framework

- Given spectral or log-mel features, use DNN or DRNN to predict spectral targets (multiple sources)
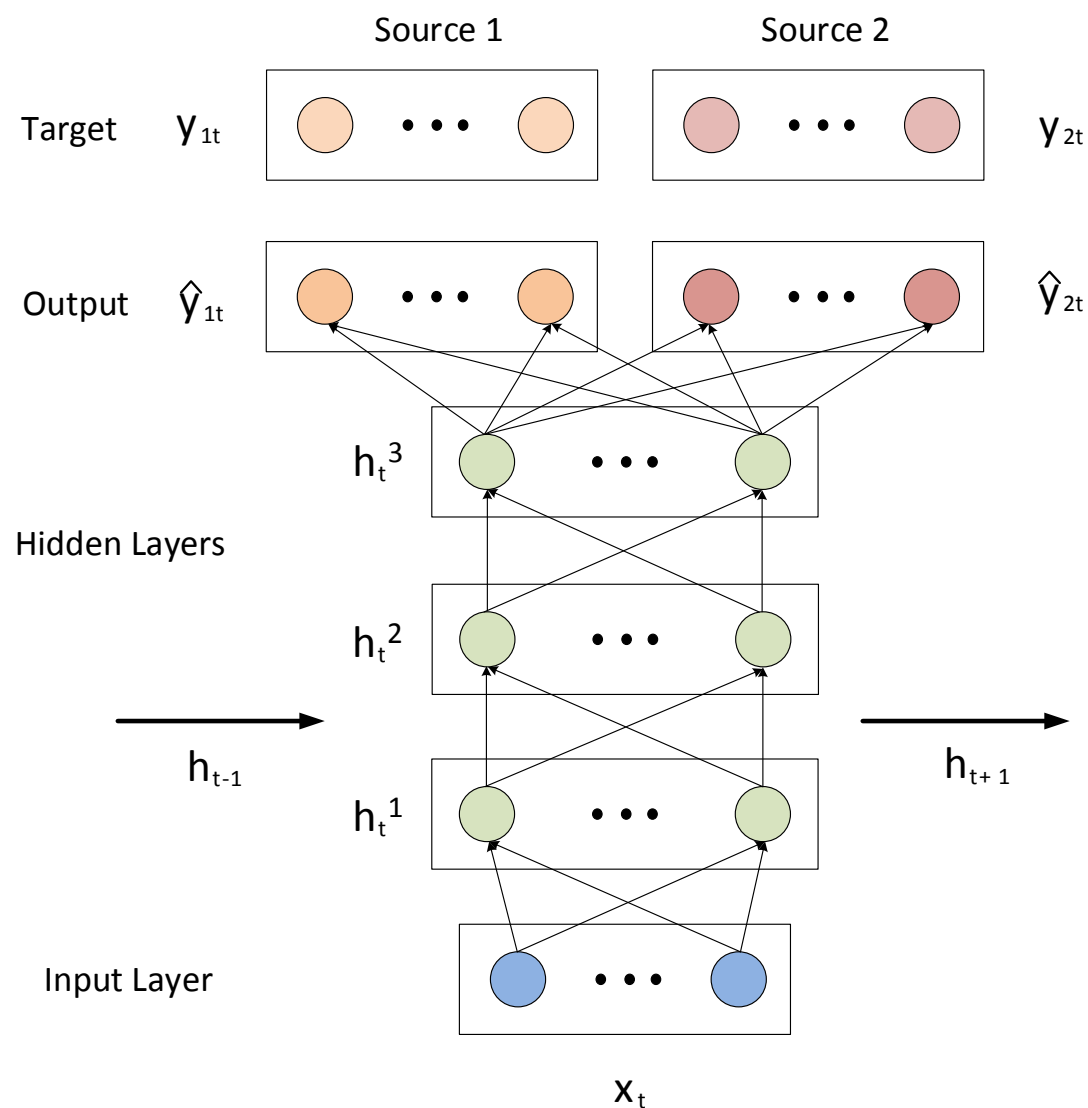- Apply time-frequency masking
- ISTFT

# Proposed Methods

- Model architecture
- Joint time frequency masking
- Discriminative training

# Model Architecture

- Jointly model two sources as targets
- Enable to use different features – log-mel, spectra
- Explore DNN, DRNN
- ➤ Time frequency masking

Source 1    Source 2

Target    $y_{1t}$    $y_{2t}$

Output    $\hat{y}_{1t}$    $\hat{y}_{2t}$

Hidden Layers

$h_t^3$

$h_t^2$

$h_{t-1}$    $h_{t+1}$

$h_t^1$

Input Layer

$x_t$

# Time Frequency Masking

- Masking – enforce the constraints that sum of the predictions equals to the original mixture

- Binary mask

$$\mathbf{M_b}(f) = \begin{cases} 1 & |\hat{\mathbf{y}}_{\mathbf{1}_t}(f)| > |\hat{\mathbf{y}}_{\mathbf{2}_t}(f)| \\ 0 & \text{otherwise} \end{cases}$$

- Soft mask

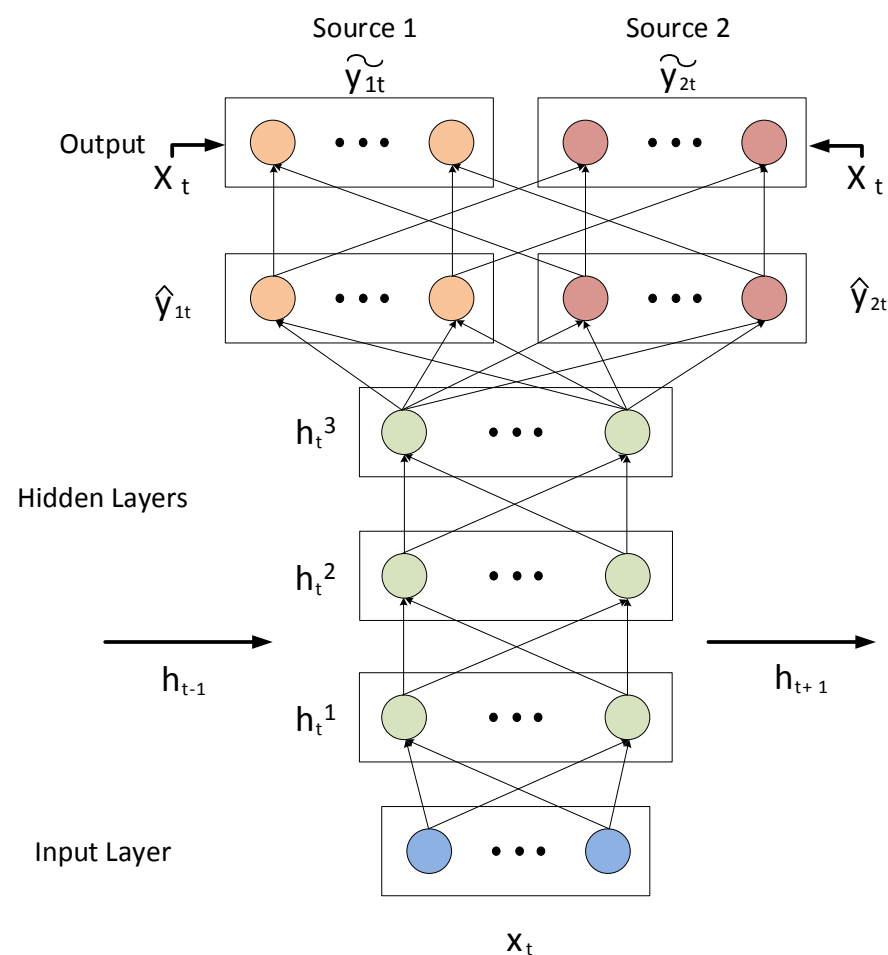$$\mathbf{M_s}(f) = \frac{|\hat{\mathbf{y}}_{\mathbf{1}_t}(f)|}{|\hat{\mathbf{y}}_{\mathbf{1}_t}(f)| + |\hat{\mathbf{y}}_{\mathbf{2}_t}(f)|}$$

- Apply the mask to predicted results

$$\hat{\mathbf{s}}_{\mathbf{1}_t}(f) = \mathbf{M}(f)\mathbf{X}_t(f)$$
$$\hat{\mathbf{s}}_{\mathbf{2}_t}(f) = (1 - \mathbf{M}(f))\,\mathbf{X}_t(f)$$
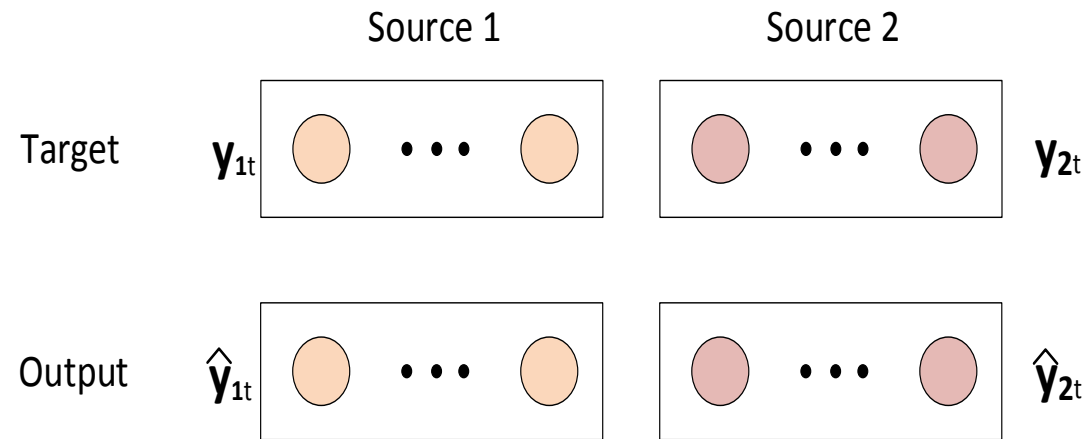
# Joint Time Frequency Masking

- Viewed the masking operation as a deterministic layer
- Train the model with masking jointly

$$\tilde{\mathbf{y}}_{1_t} = \frac{\left|\hat{\mathbf{y}}_{1_t}\right|}{\left|\hat{\mathbf{y}}_{1_t}\right| + \left|\hat{\mathbf{y}}_{2_t}\right|} \odot \mathbf{X}_t$$

$$\tilde{\mathbf{y}}_{2_t} = \frac{\left|\hat{\mathbf{y}}_{2_t}\right|}{\left|\hat{\mathbf{y}}_{1_t}\right| + \left|\hat{\mathbf{y}}_{2_t}\right|} \odot \mathbf{X}_t$$

# Discriminative training



Source 1     Source 2

Target $\mathbf{y_{1t}}$   $\mathbf{y_{2t}}$

Output $\hat{\mathbf{y}}_{1t}$   $\hat{\mathbf{y}}_{2t}$

- Minimize squared error

$$||\hat{\mathbf{y}}_{\mathbf{1}_t} - \mathbf{y}_{\mathbf{1}_t}||_2^2 + ||\hat{\mathbf{y}}_{\mathbf{2}_t} - \mathbf{y}_{\mathbf{2}_t}||_2^2$$

- Enforce source to interference ratio

$$||\hat{\mathbf{y}}_{\mathbf{1}_t} - \mathbf{y}_{\mathbf{1}_t}||_2^2 - \gamma ||\hat{\mathbf{y}}_{\mathbf{1}_t} - \mathbf{y}_{\mathbf{2}_t}||_2^2 + ||\hat{\mathbf{y}}_{\mathbf{2}_t} - \mathbf{y}_{\mathbf{2}_t}||_2^2 - \gamma ||\hat{\mathbf{y}}_{\mathbf{2}_t} - \mathbf{y}_{\mathbf{1}_t}||_2^2$$

ILLINOIS

# Experimental Setting

- TIMIT dataset
  - Mixed the speech from a male and a female speaker at 0 dB
- Circular shift to increase the variety of training samples
- Deep learning models
  - RELU
  - L-BFGS for optimization
  - Use 2 hidden layers with 150 hidden units
- BSS EVAL metric (SDR, SIR, SAR)
  - SIR - Suppression of interference
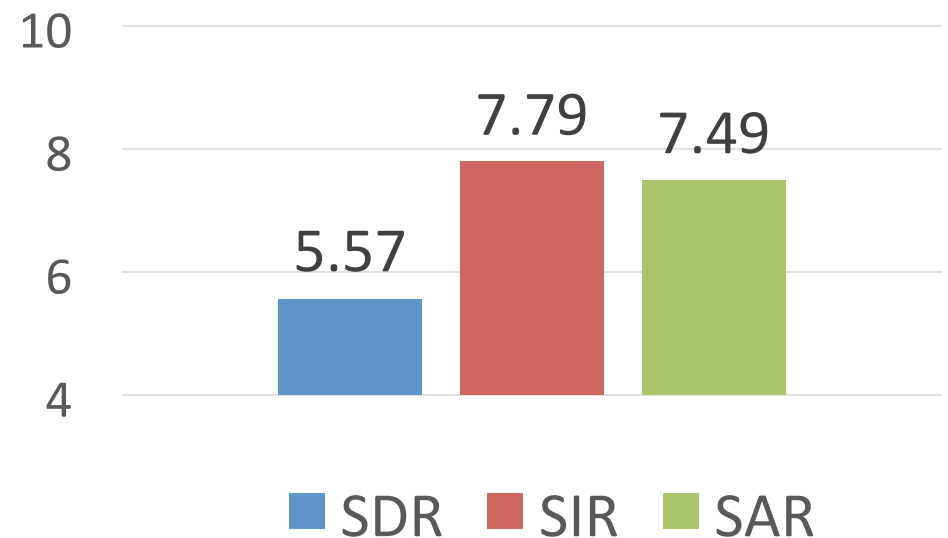  - SAR - Artifacts introduced by the separation process
  - SDR - Overall performance

ILLINOIS

# Baseline - NMF

- 512-point STFT
- Generalized KL-divergence metric

NMF with Binary Masking
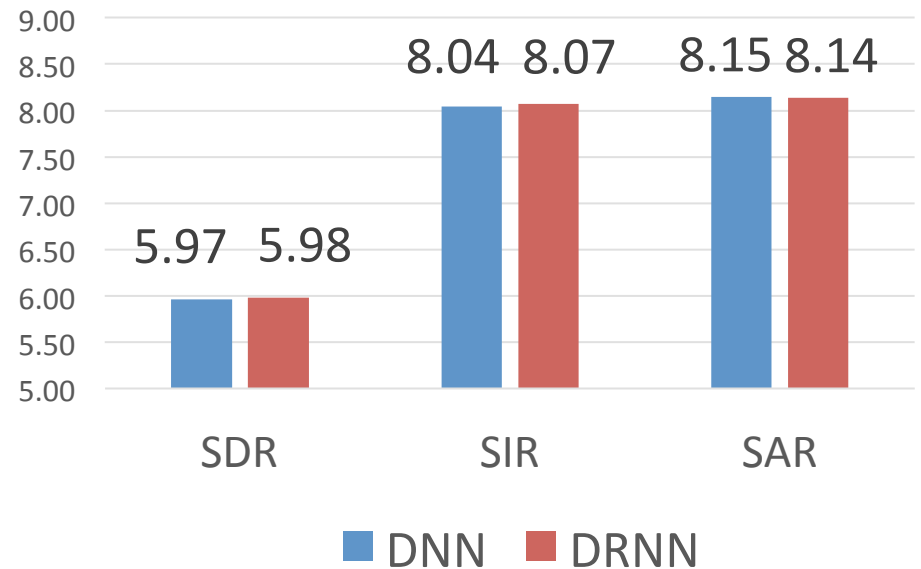


NMF with Soft Masking

# Experimental Results – DNN vs. DRNN

- Spectral features
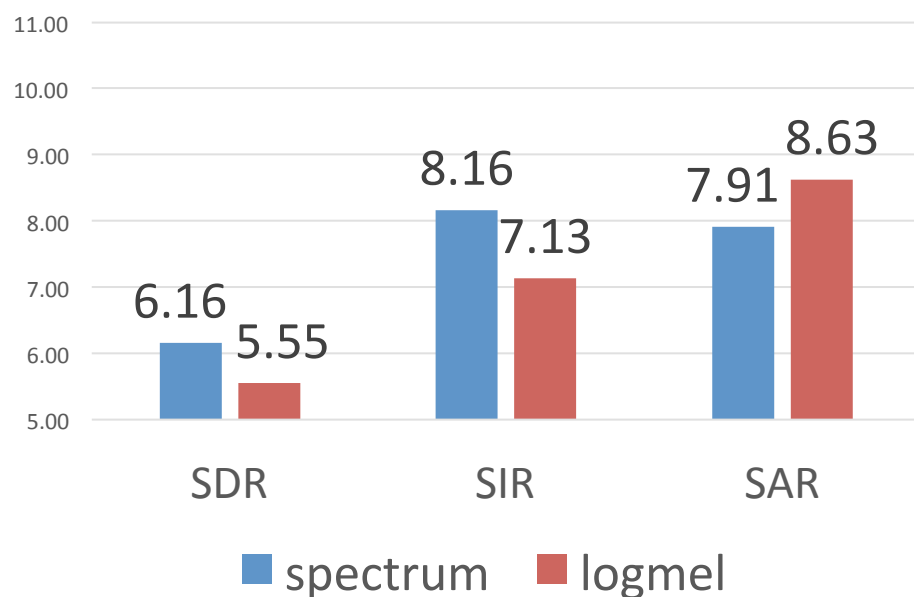- No significant difference



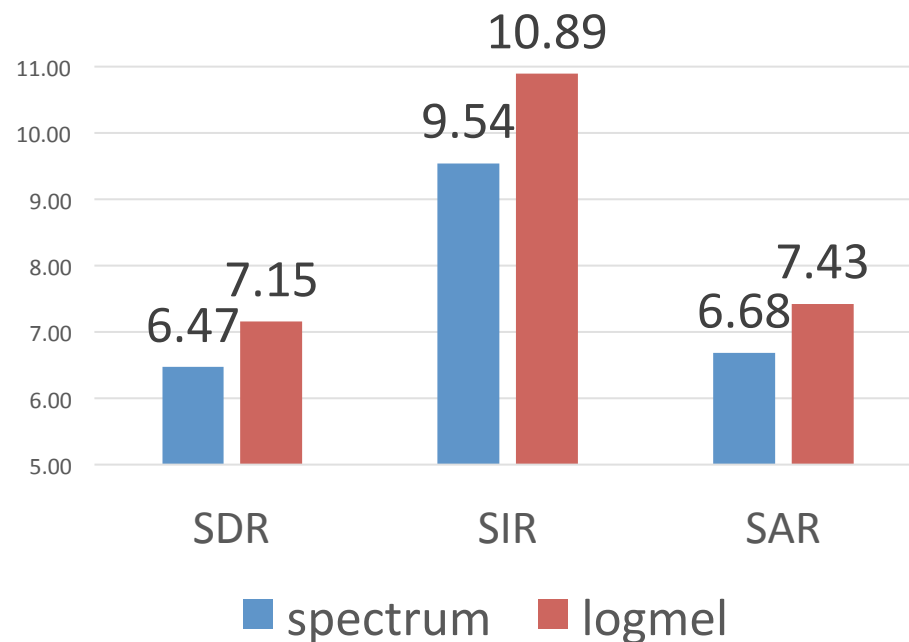win=1, soft mask

win=3, soft mask

# Experimental Results – Features, Joint Mask Training

- Log-mel features perform better with joint mask training
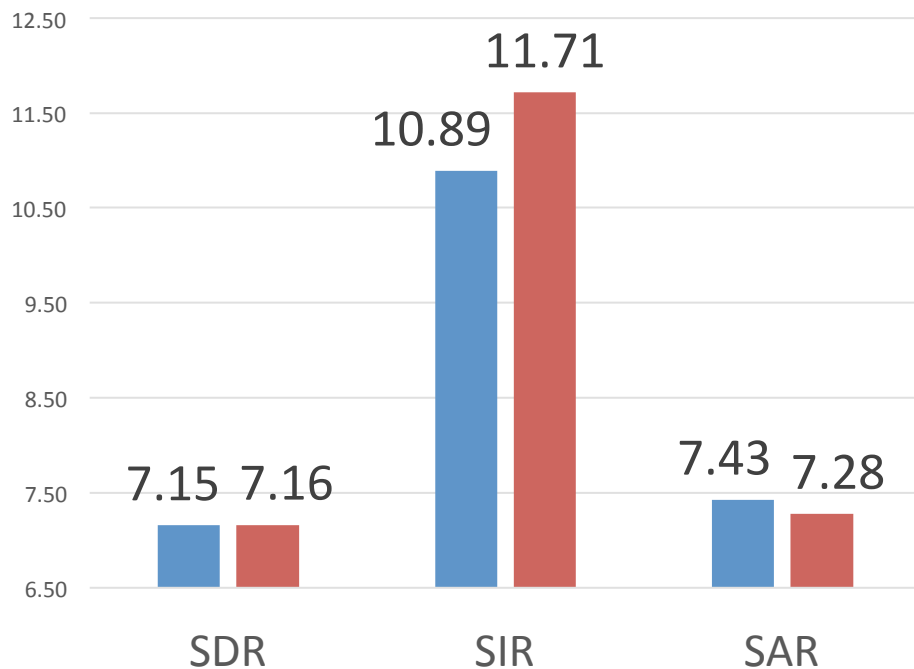


win=1, without joint mask training

win=1, with joint mask training

ILLINOIS

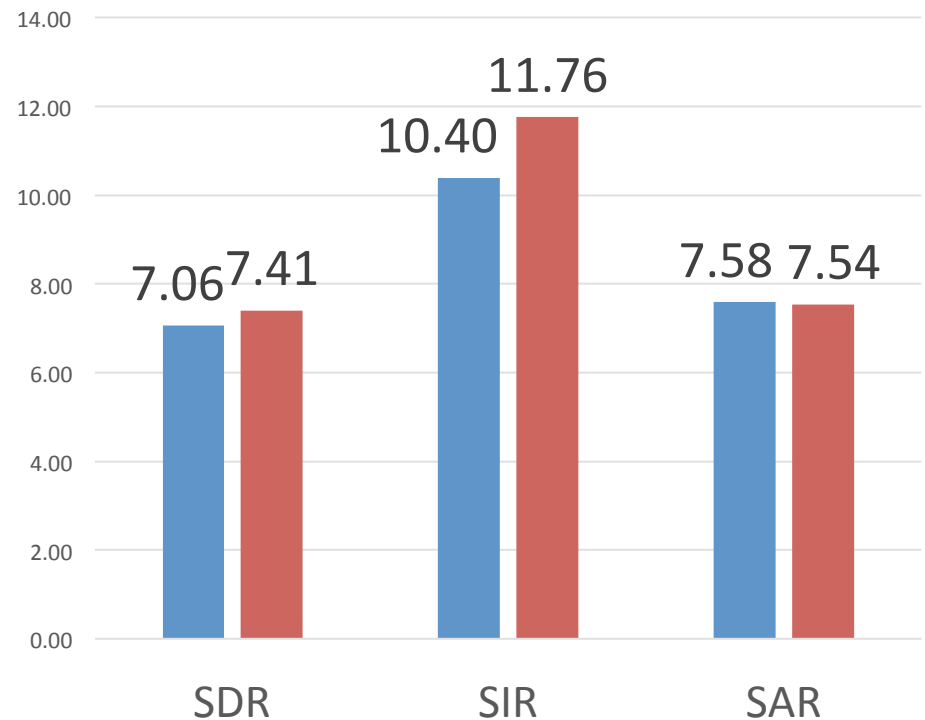# Experimental Results – Discriminative Training

- Discriminative training provides extra regularization
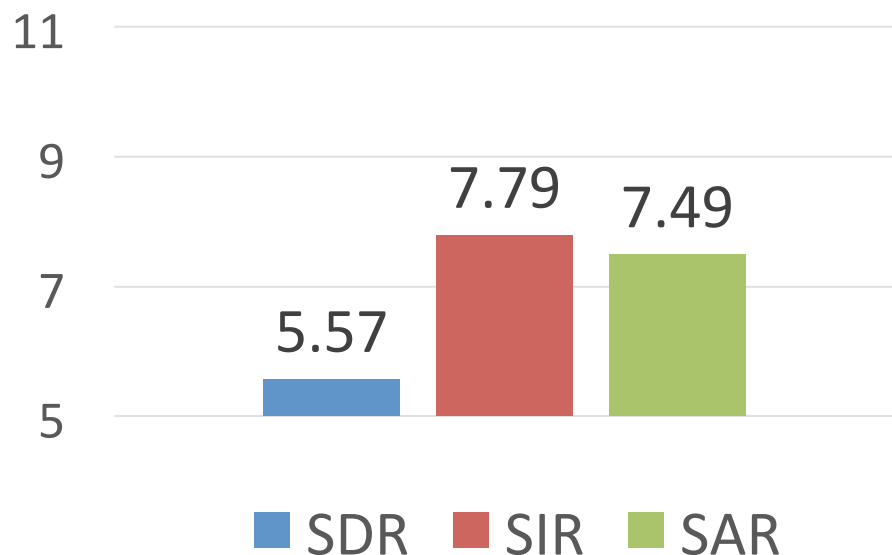


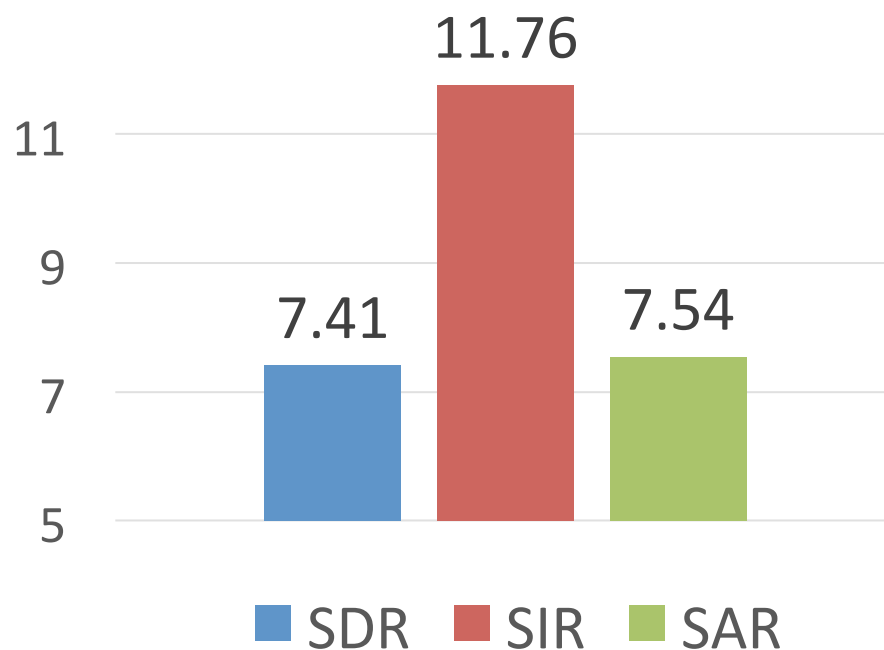win=1, logmel with joint mask training

win=3, logmel with joint mask training

# Summary

- Comparison between NMF and DRNN with log-mel, joint mask training, and discriminative training objective



NMF with soft masking

SDR: 5.57, SIR: 7.79, SAR: 7.49

DRNN with soft masking

SDR: 7.41, SIR: 11.76, SAR: 7.54

ILLINOIS

# DEMO



TIMIT mixture

ILLINOIS

# Conclusion

- Propose using deep learning models for monaural speech separation.
  - Propose the joint optimization of a soft masking function and deep learning models
  - Discriminative training criterion to further improve the SIR
- Overall, our proposed models achieve 3.8~4.9 dB SIR gain compared to the NMF baseline
- Future work
  - Explore longer temporal information with neural networks
  - Apply many other applications such as robust ASR

ILLINOIS

# Thank you!

https://sites.google.com/site/deeplearningsourceseparation

# MIR1K GNSDR

Train on two singers and test on other 17 singers



MIR 1K GNSDR

- RPCA [ICASSP12]: 3.15
- RPCAh [MM 12]: 3.84
- MLRR [ISMIR13]: 3.85
- SBS [ISMIR12]: 4.96
- DRNN: 7.45