

Speech Denoising via Nonnegative Matrix Factorization

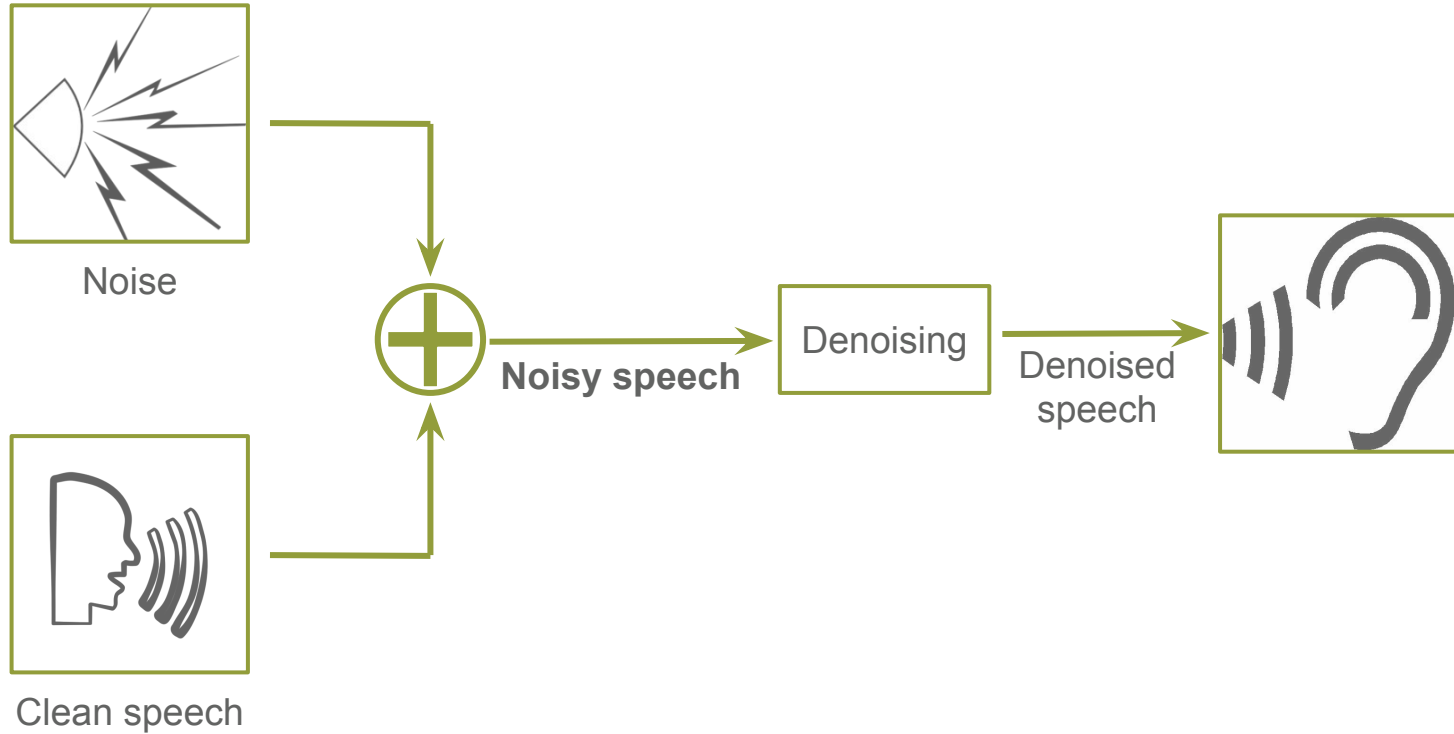


Your Answer

Oleg Alenkin
Artyom Chashchin
Vladimir Chernykh
German Novikov

December 2016

Problem



Applications

Automatic Speech
Recognition



Telephone conversations



Hearing aids

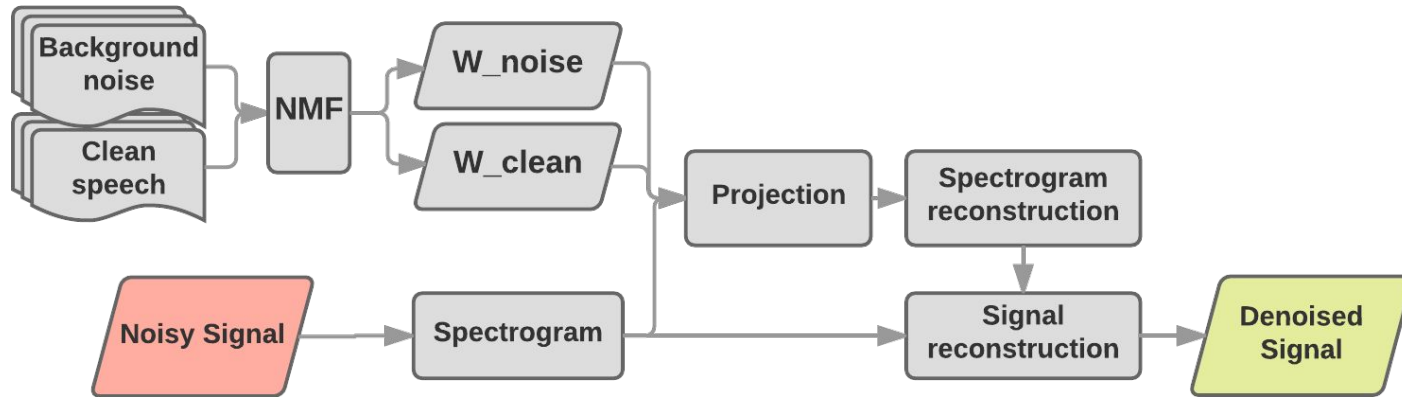


Data

- CHiME Speech Separation and Recognition Challenge
http://spandh.dcs.shef.ac.uk/chime_challenge/chime_download.html
Recordings of WSJ utterances + 8 hours of noise
- Berlin Database of Emotional Speech
<http://www.emodb.bilderbar.info/download/>
Clean utterances
- Aurora noising
<http://aurora.hsnr.de/download.html>

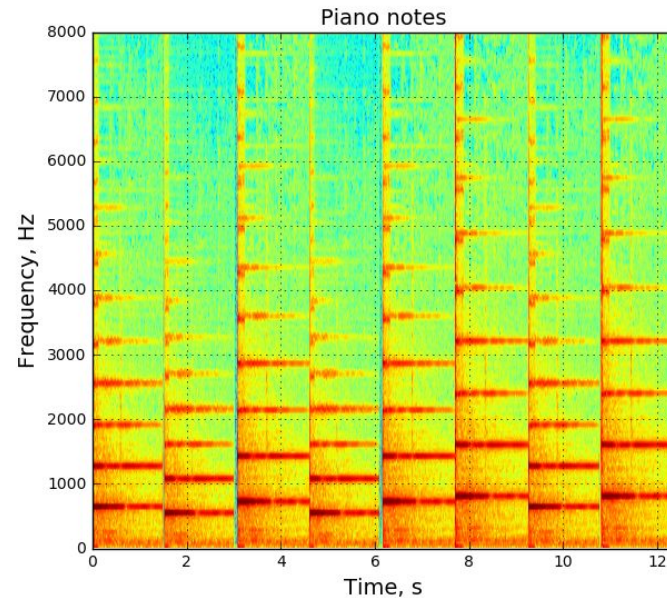
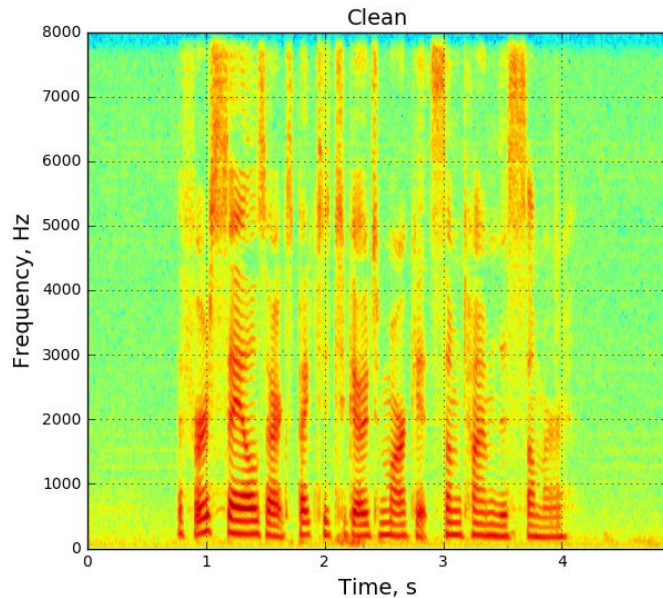
Solution

- Learn frequency patterns from speech and noise via NMF
- Decompose new signal on joint “dictionary” of patterns
- Take only projection corresponding to “clean speech”



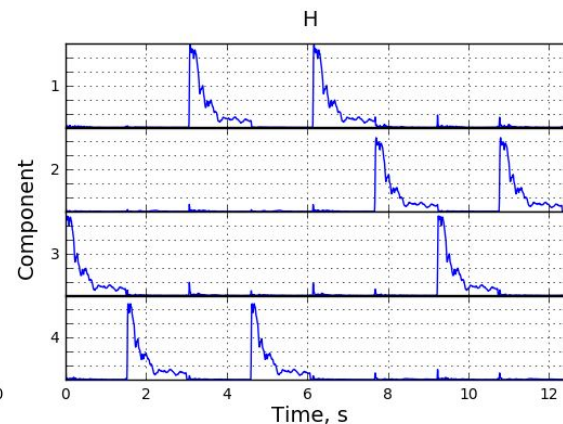
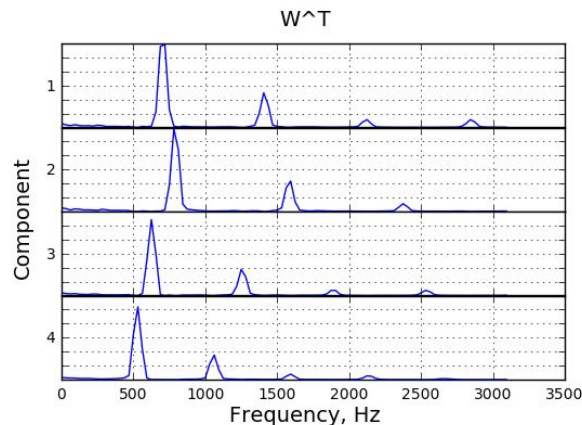
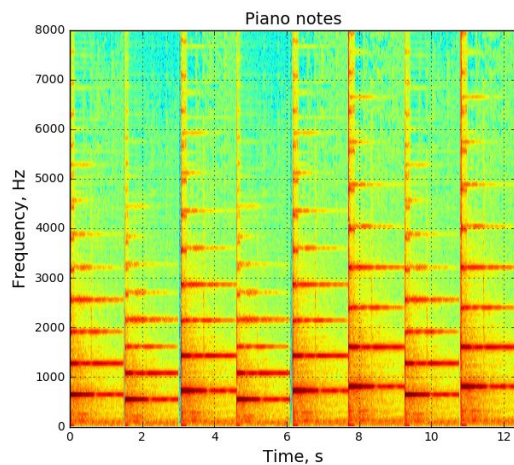
Spectrogram

- Represent signal in Time-Frequency domain
- Built via Short-Term Fourier Transform - FFT with sliding overlapped window
- STFT - Complex spectrogram $S \Rightarrow$ Amplitudes V



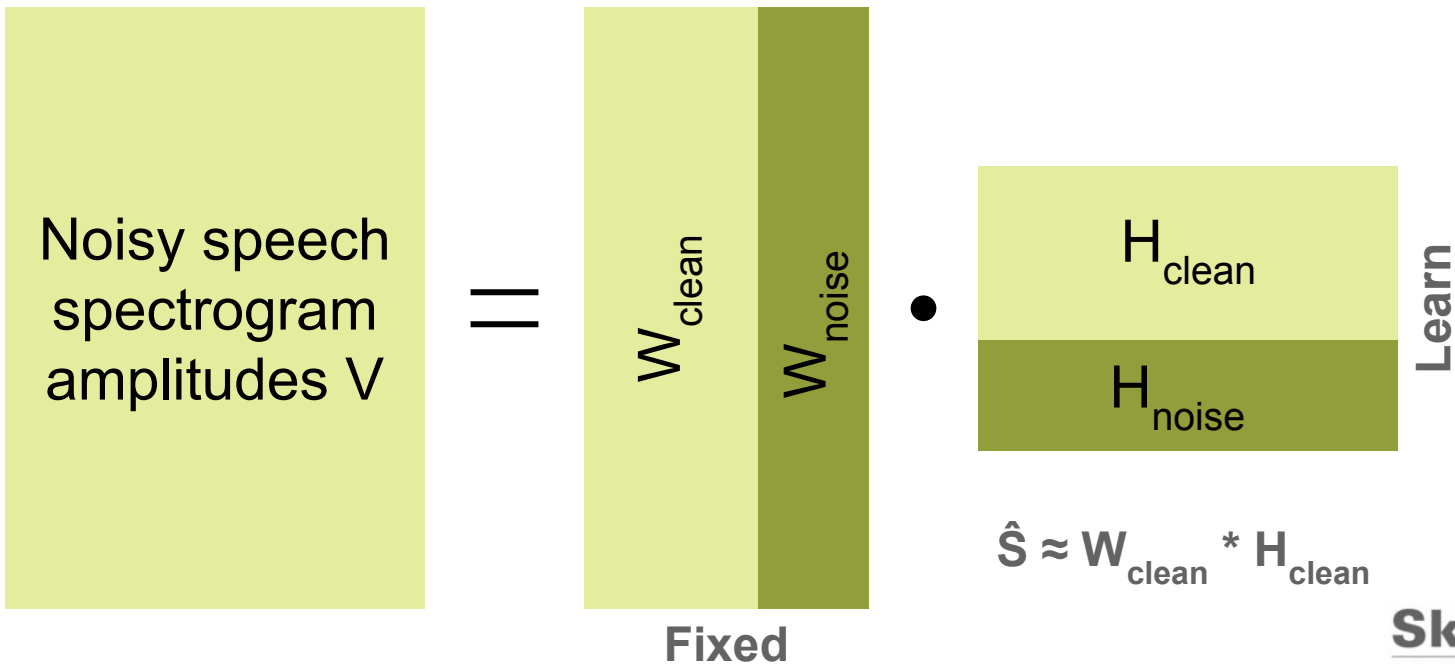
NMF

- Factorization $\mathbf{V} \approx \mathbf{W} * \mathbf{H}$, where \mathbf{V} , \mathbf{W} , \mathbf{H} - real nonnegative
- Interpretation: \mathbf{W} - frequency patterns, \mathbf{H} - time-activation matrix
- Hidden dimension \approx number of phonemes ≈ 40
- Learn speech and noise “building blocks”



Projection

- Join “dictionaries” - concatenate matrices W_{noise} and W_{clean}
- Project signal onto them



NMF: how to compute

Optimization problem:

$$(W^*, H^*) = \operatorname{argmin}_{W \geq 0, H \geq 0} D(V, WH)$$

$$D(P, Q) = \sum_{i=1}^m \sum_{j=1}^n d(p_{ij}, q_{ij})$$

The most popular metrics:

$$d(p, q) = (p - q)^2 \quad \text{Frobenius norm}$$

$$d(p, q) = p \ln\left(\frac{p}{q}\right) - p + q \quad \text{KL divergence}$$

Multiplicative Update
Method:

$$[\nabla_H]_{kj} = \frac{\partial D(V, WH)}{\partial h_{kj}} = [\nabla_H^+]_{kj} - [\nabla_H^-]_{kj}$$

$$h_{kj} \leftarrow h_{kj} - \frac{h_{kj}}{[\nabla_H^+]_{kj}} ([\nabla_H^+]_{kj} - [\nabla_H^-]_{kj})$$

NMF: how to compute

Alternating Nonnegative Least Squares:

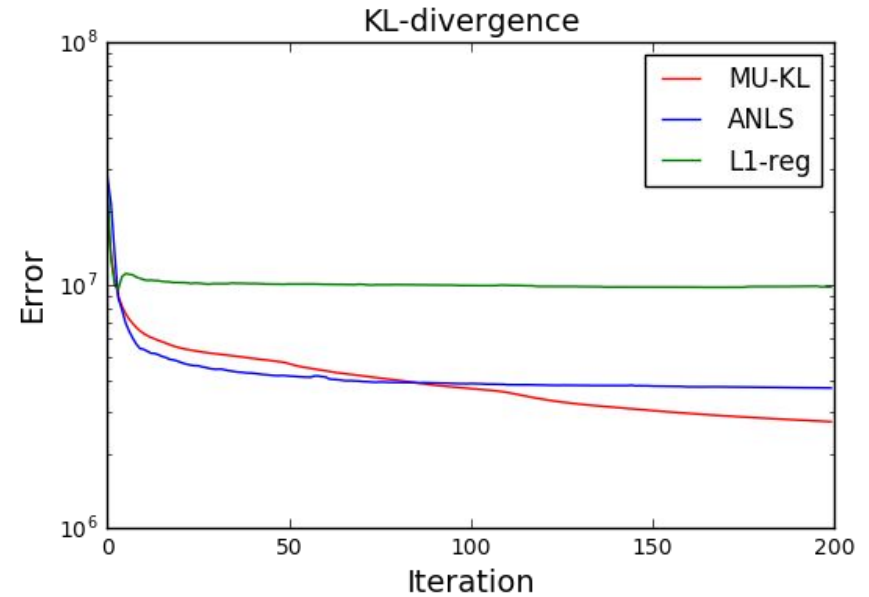
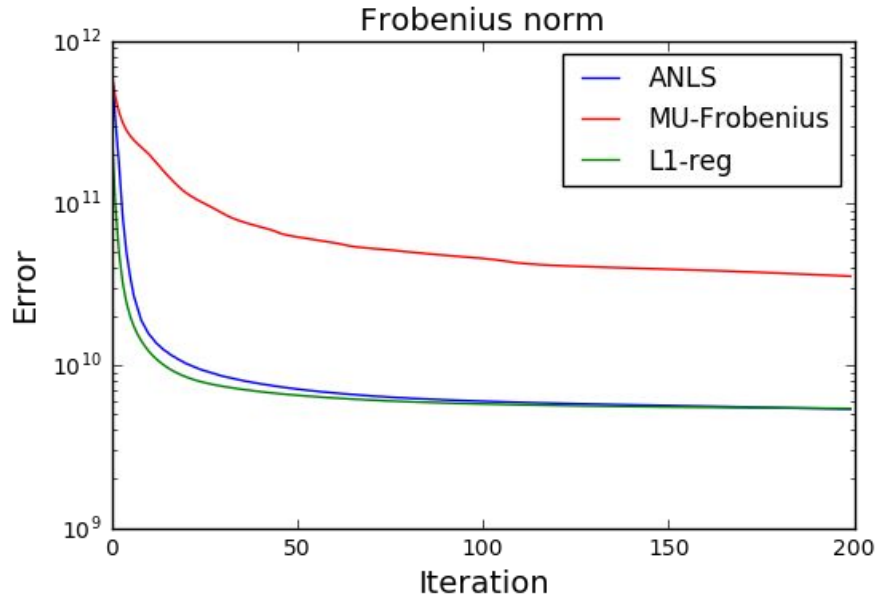
1) Initialize $W_{ia}^1 \geq 0, H_{bj}^1 \geq 0, \forall a, i, b, j.$

2) For $k=1,2,\dots$

$$W^{k+1} = \operatorname{argmin}_{W \geq 0} D(V, W H^k),$$

$$H^{k+1} = \operatorname{argmin}_{H \geq 0} D(V, W^{k+1} H).$$

Methods' Convergence



Quasi-Newton method

$$(W^*, H^*) = \operatorname{argmin}_{W>0, H>0} D_{KL}(V, WH)$$

$$W \leftarrow \max(\varepsilon, W - H_W^{-1} \nabla_W D_{KL})$$

$$\nabla_W D_{KL} = H^T J_{M \times K} - H^T (V \oslash (WH))$$

$$H_W = \operatorname{diag}\{h_{W,m}, m = 1, \dots, M\}$$

$$h_{W,m} = H \operatorname{diag}\{[V \oslash (Q \otimes Q)]_{m,:}\} H^T$$

$$H \leftarrow \max(\varepsilon, H - H_H^{-1} \nabla_H D_{KL})$$

$$\nabla_H D_{KL} = J_{M \times K} W^T - (V \oslash (WH)) W^T$$

$$H_H = \operatorname{diag}\{h_{H,k}, k = 1, \dots, K\}$$

$$h_{H,k} = W^T \operatorname{diag}\{[V \oslash (Q \otimes Q)]_{:,k}\} W$$

Problem: hard to compute inverse Hessian

Signal reconstruction

- Naive method with zero phase:

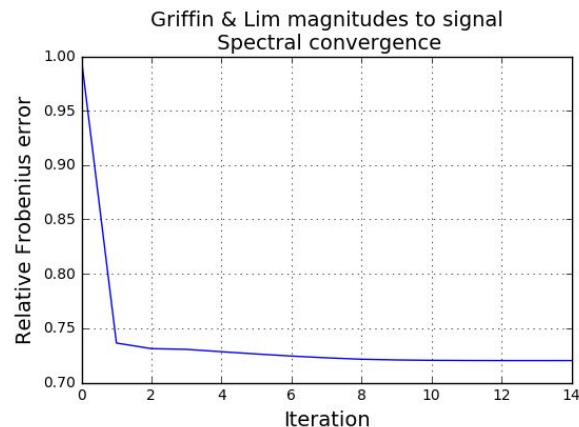
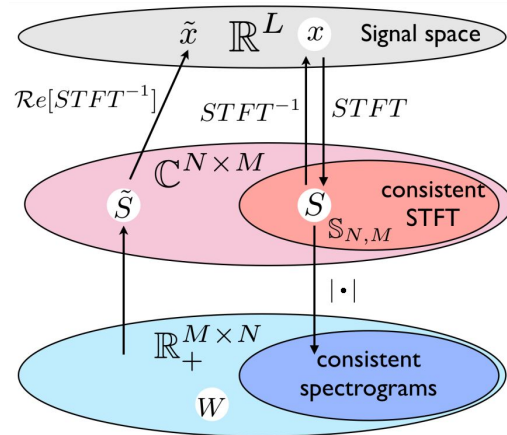
$$\hat{X} = STFT^{-1}(\hat{S})$$

- Noisy signal phases:

$$\hat{X} = STFT^{-1}(\hat{S} \times \exp(i\angle STFT(X_{noisy})))$$

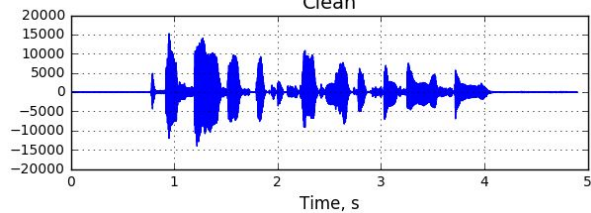
- Griffin & Lim iterative method:

$$\hat{X}_n = STFT^{-1}(\hat{S} \times \exp(i\angle STFT(\hat{X}_{n-1})))$$

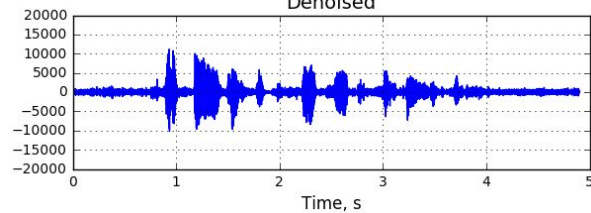


Demo

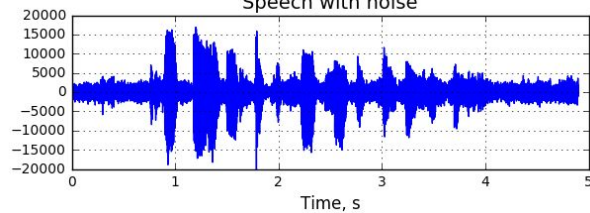
Clean



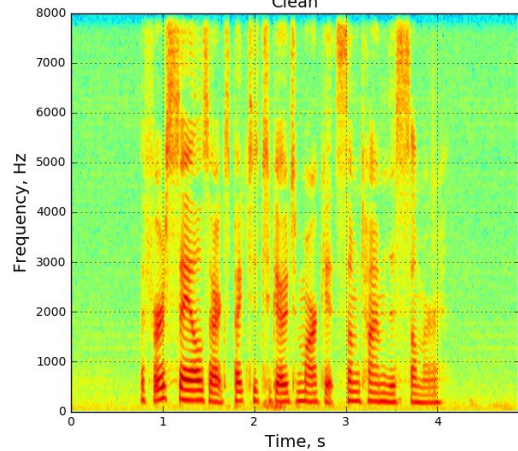
Denoised



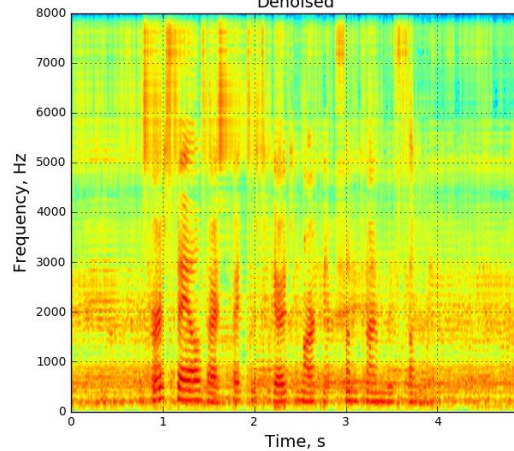
Speech with noise



Clean



Denoised



Speech with noise

