



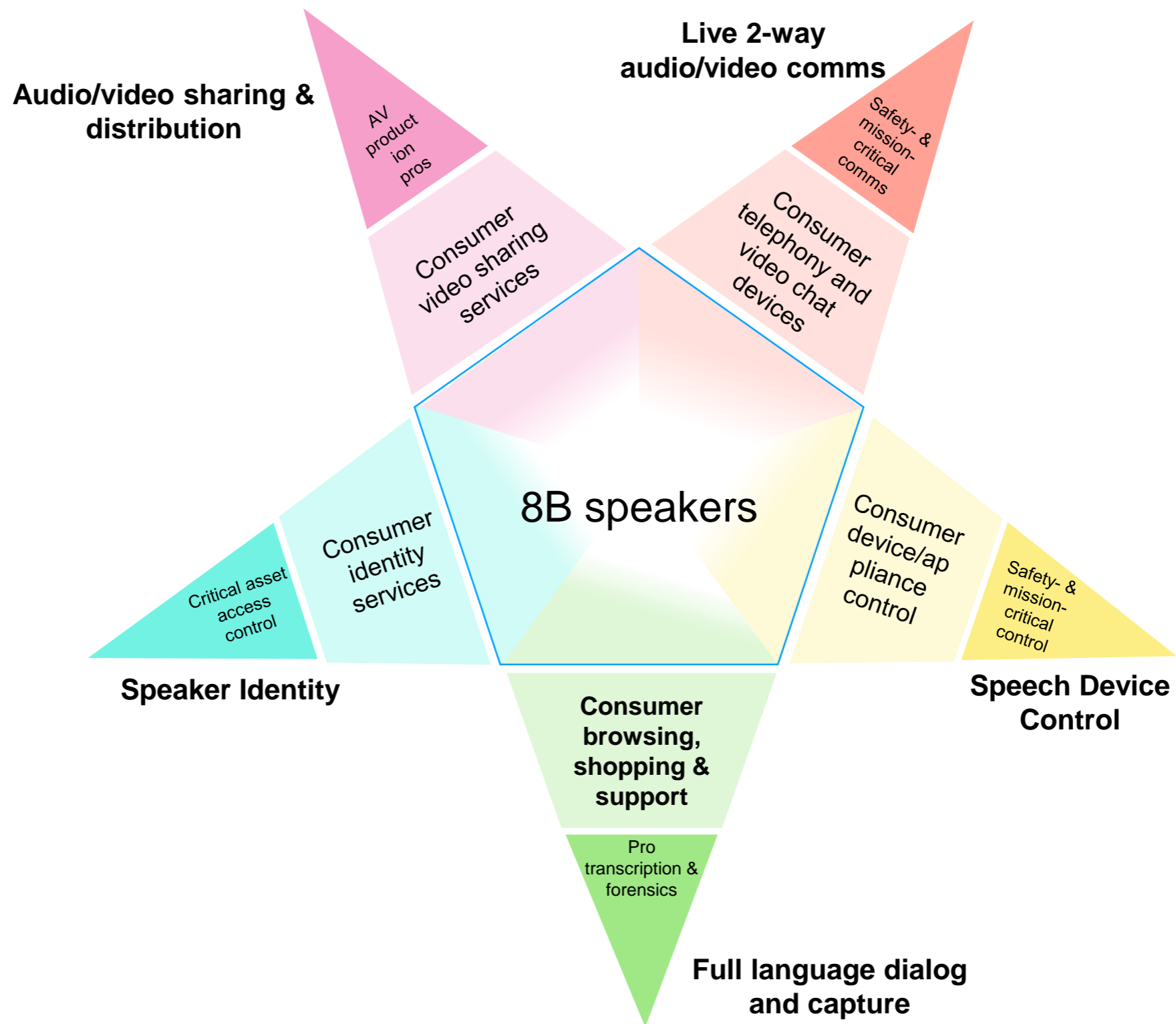
**babblelabs**

**Samer Hijazi**

**CTO**

**BabbleLabs Inc.**

# The New World of Speech Technology



# The Opportunity

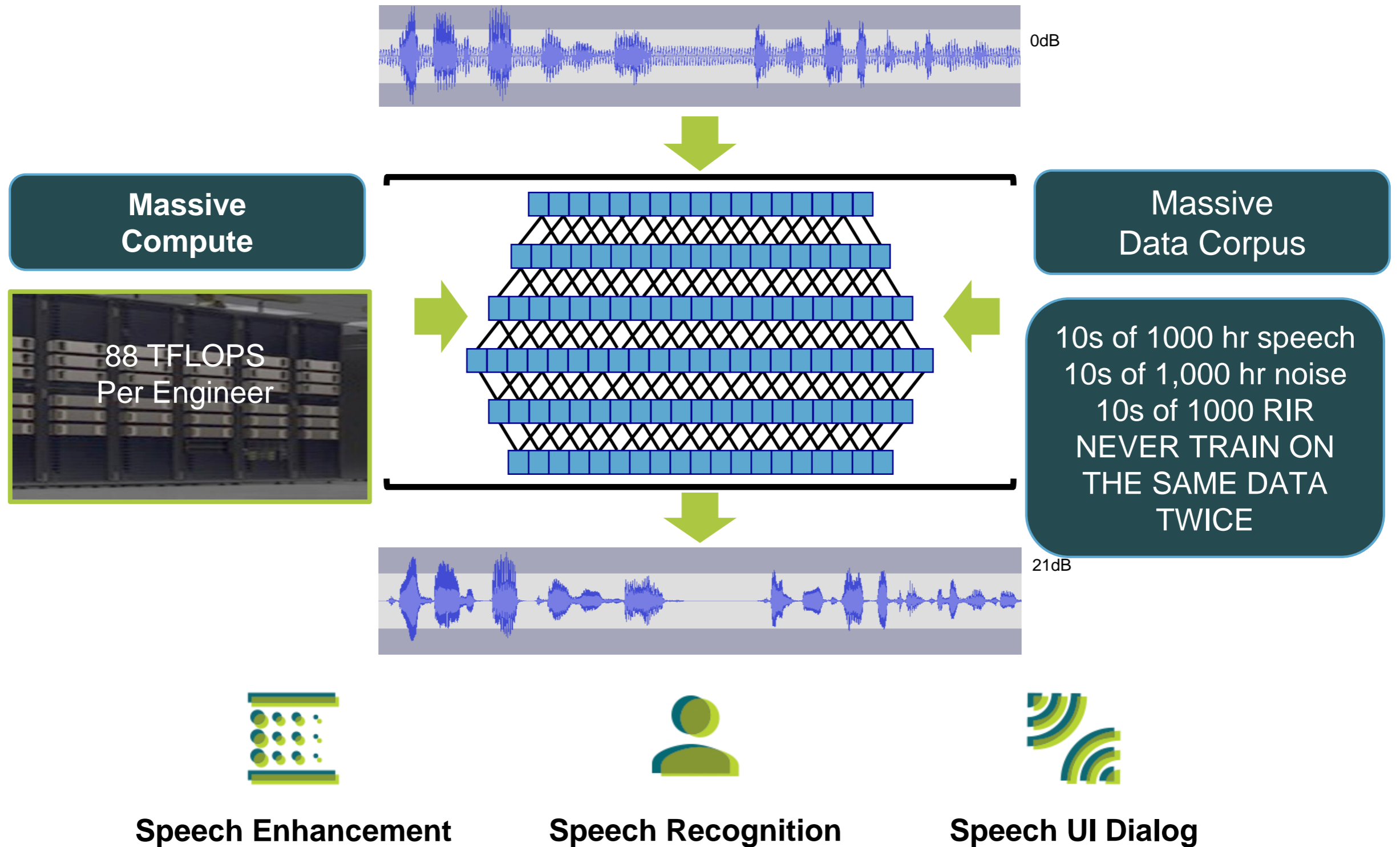
- 22B microphones by 2020
- 7B phones + radios + TVs delivering voice
- YouTube uploads: 13B minutes per year
- 200T minutes per year of device interaction
- 1Q words per year in voice calls

**Speech Market Growth: 38.3%**

-Statista 2018: speech recognition technology market 2016-2014

# AI meets speech

*more sophisticated models, more data, more training*

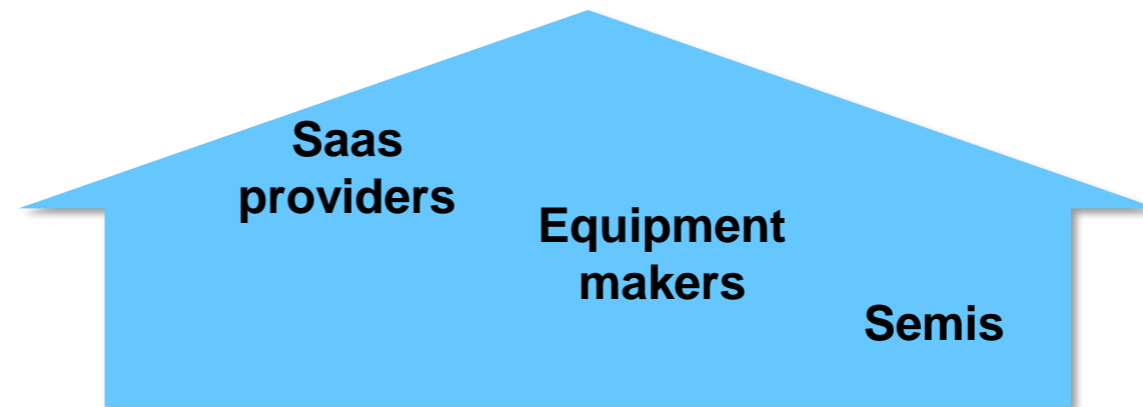


# Technology → Product → Customers → End Users

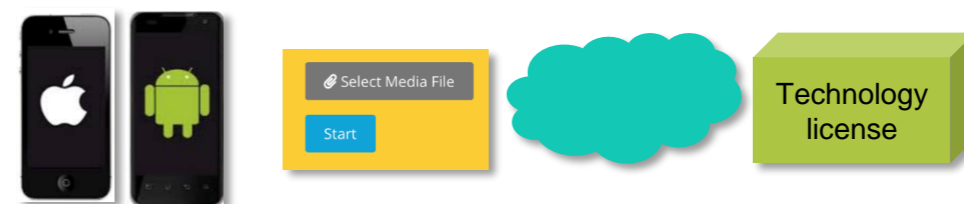
End user:  
*their problem*

**Consumer audio/video sharer:**  
*Recording in the real world*

Customers



Product: Platform-optimized solutions

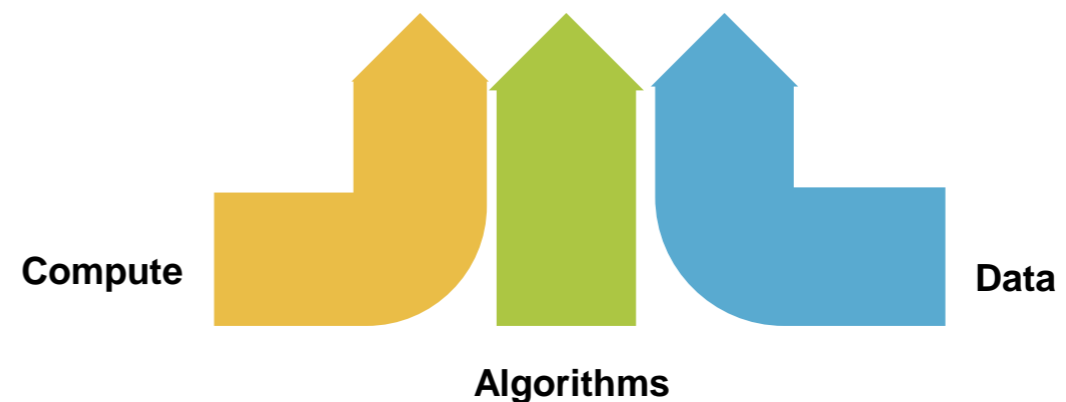


Product: Deep learning speech software



**Speech Enhancement**


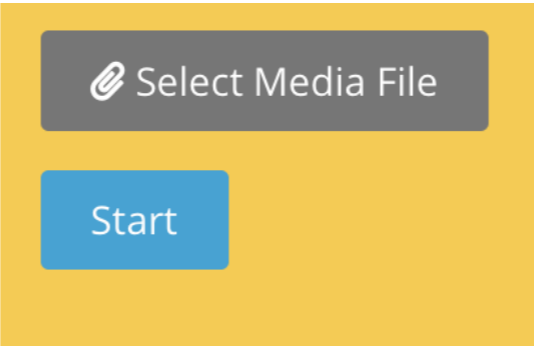
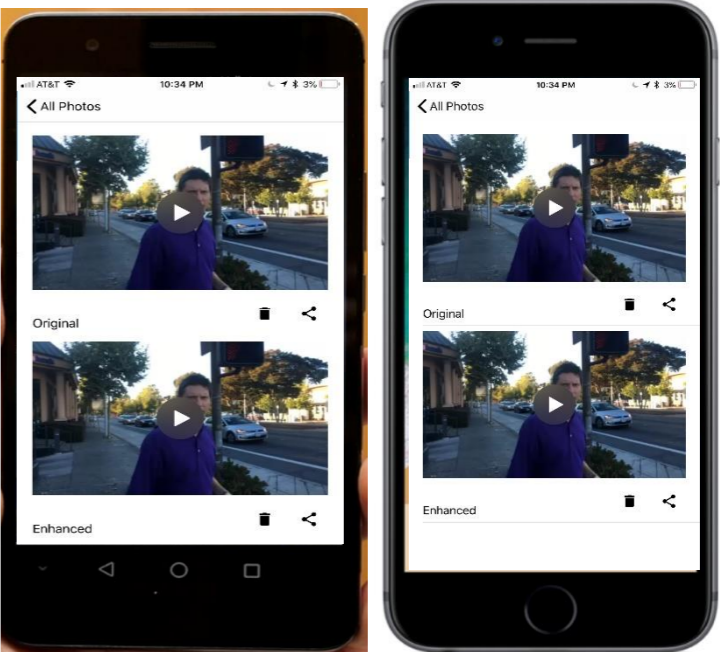
Technology: Unique data-sets and training



# Clear Speech Everywhere

*In production for real-world video sharing, production, streaming, and audio*

**Foundation to  
product release  
in 28 weeks!**

Common product delivered across platforms	
<i>Primary products</i>	<div style="display: flex; justify-content: space-around;"><div style="text-align: center;"><h3>Cloud API</h3><pre>host=api.babblelabs.com token=\$(cat token.txt) #token.txt created by login.sh email=&lt;yourEmailAddr&gt; api="audio" # audio or video contentType="audio/wav" # a normal supported mime-type of audio or video for infile in *.wav # Replace this with a different pattern, or a list of file names (of the same type) do ( set -x time curl --progress-bar \ --data-binary @\$infile \ --dump-header \$infile.hdr \ -H "Content-Type":"\$contentType" \ -o \$infile.denoised \ -H "Authorization: Bearer \$token" \ "\$https://\$host/audioEnhancer/api /\$api/stream/\$email" ) &amp; done</pre></div><div style="text-align: center;"><h3>On-Device License</h3></div></div>
<i>For visibility and demonstration</i>	<div style="display: flex; justify-content: space-around;"><div style="text-align: center;"><h3>Web UI</h3></div><div style="text-align: center;"><h3>Android &amp; iPhone Apps</h3></div></div>



**babble**labs

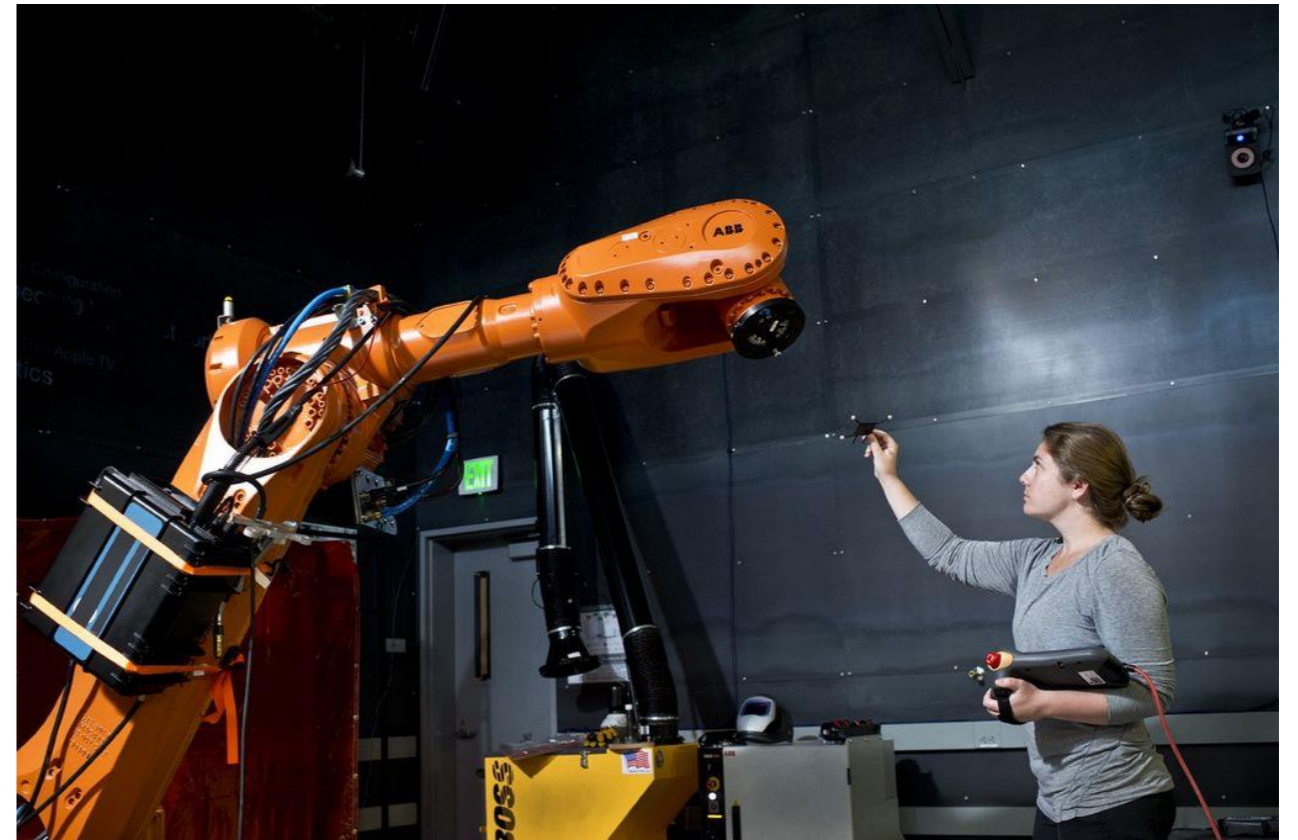
What is Speech Enhancement

# Human – Human Interface Challenges





# Human – Machine Interface Challenges



# BabbleLabs Answer to these Challenges: Clear Cloud™



**Noisy**



**Enhanced**

# Outline

**A bit about noisy speech**

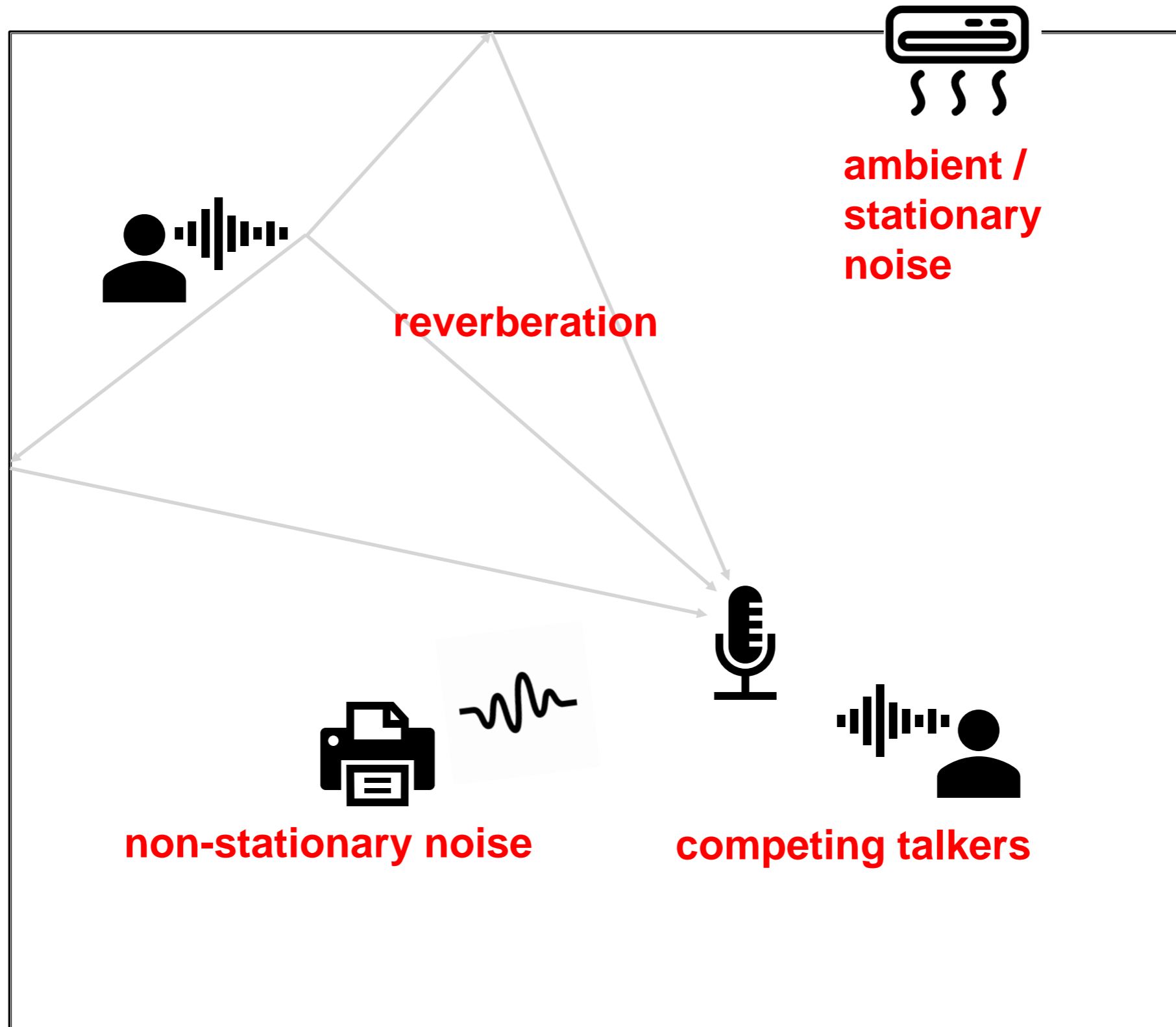
**Traditional speech enhancement**

**Deep neural network approaches**

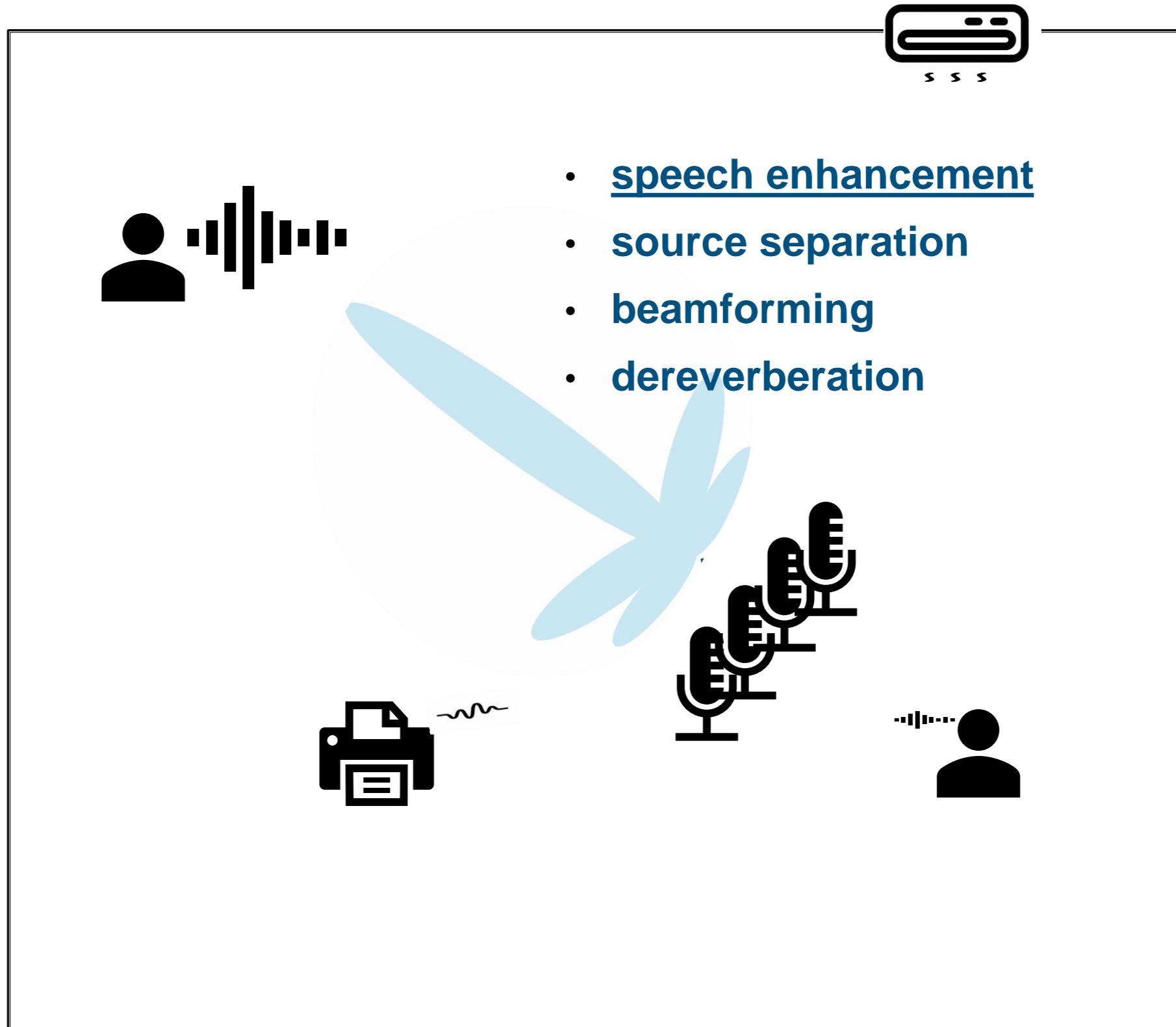
**Closing thoughts**



# Acoustic Impairments Model



# Solutions to Acoustic Impairments



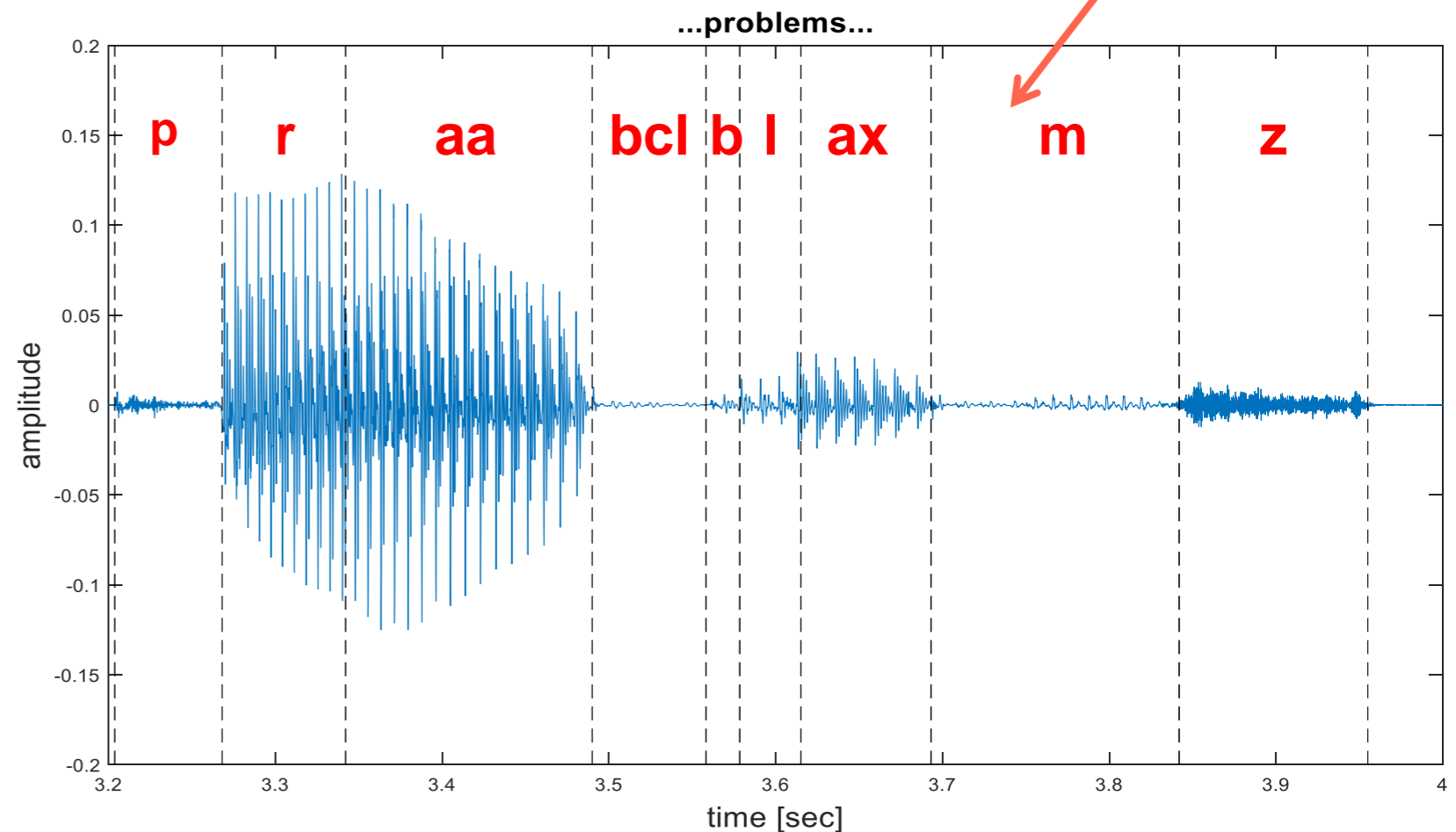
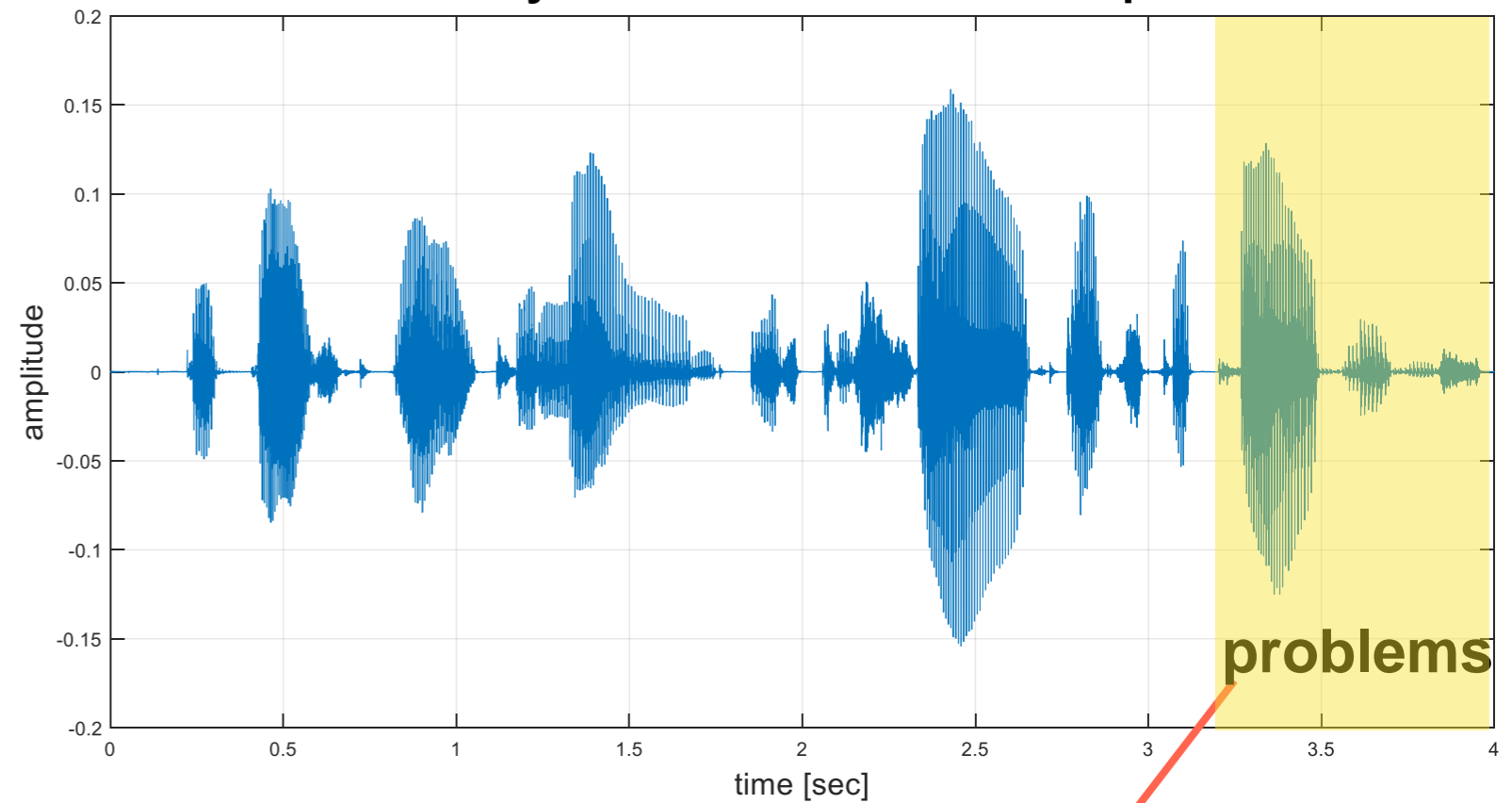
# Classes of Speech Sounds



## ARPABET Phonetic Symbols

Phone	Phoneme
Vowel	iy, ih, eh, ey, ae, <b>aa</b> , aw, ay, ah, ao, oy, ow, uh, uw, ux, er, <b>ax</b> , ix, axr, axh
Semivowel	<b>l</b> , <b>r</b> , w, y, hh, hv, el
Affricate	jh, ch
Stops	<b>b</b> , d, g, <b>p</b> , t, k, dx, q
Nasal	<b>m</b> , n, ng, em, en, eng, nx
Fricative	s, sh, <b>z</b> , zh, f, th, v, dh

The best way to learn is to solve extra problems.



# Speech Spectrum

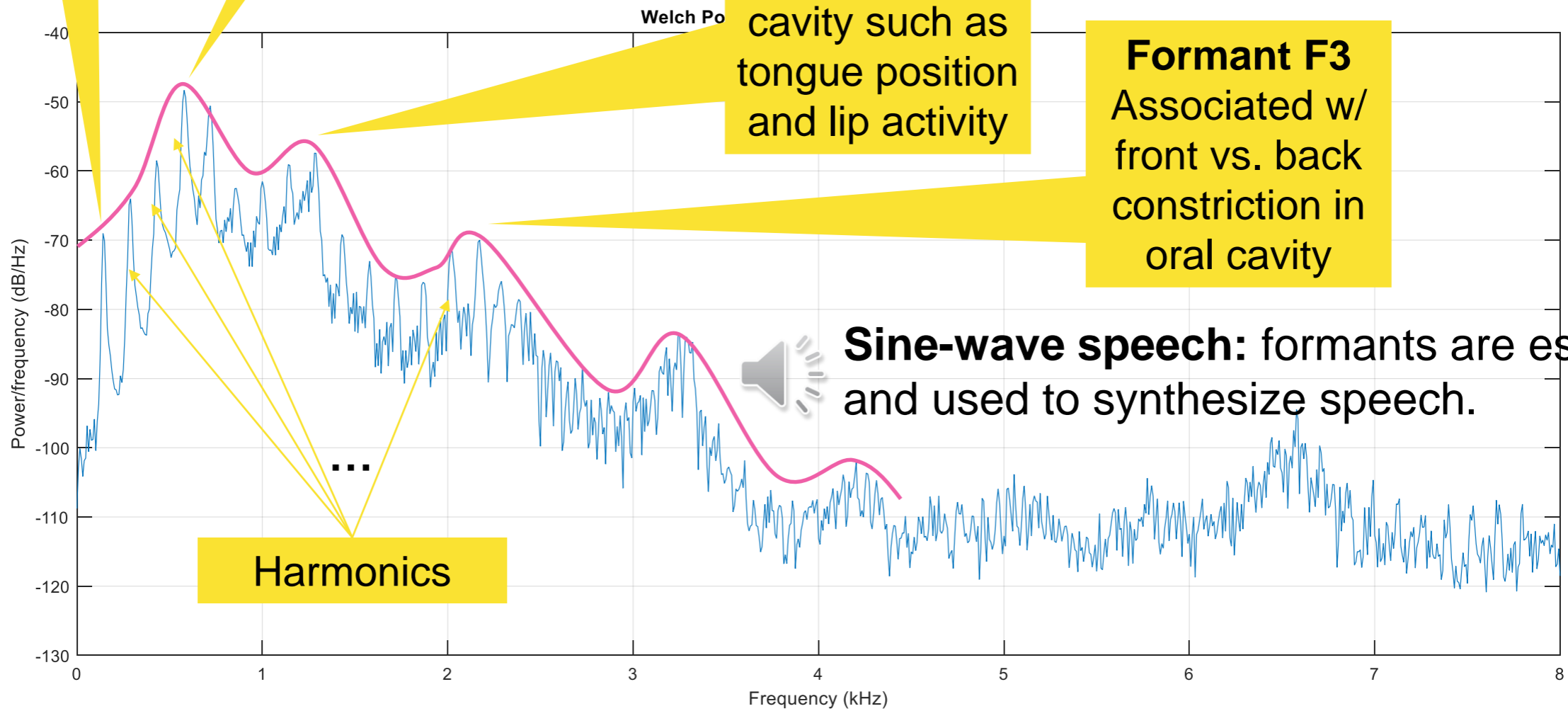
Spectral Characteristics	Frequency Range [Hz]
Fundamental Frequency, F0	Females/Children: 200 to 400 Males: 60 to 150
Harmonics	Up to 20K
Hearing Range	20 to 20K
Typical Audio Sampling Rates	In KHz: 8 (Telephony), 11.025, 22.05 (MP3s), 32 (Cassette), 44.1 (CD), 48 (DVD)

**Fundamental Frequency F0**  
(e.g. ~150Hz)

**Formant F1**  
Associated w/  
size of mouth  
opening;  
proportional to  
frequency  
e.g. AA ~580Hz

**Formant F2**  
Associated w/  
changes in oral  
cavity such as  
tongue position  
and lip activity

**Formant F3**  
Associated w/  
front vs. back  
constriction in  
oral cavity



**Sine-wave speech:** formants are estimated and used to synthesize speech.

# Noise and Speech Levels

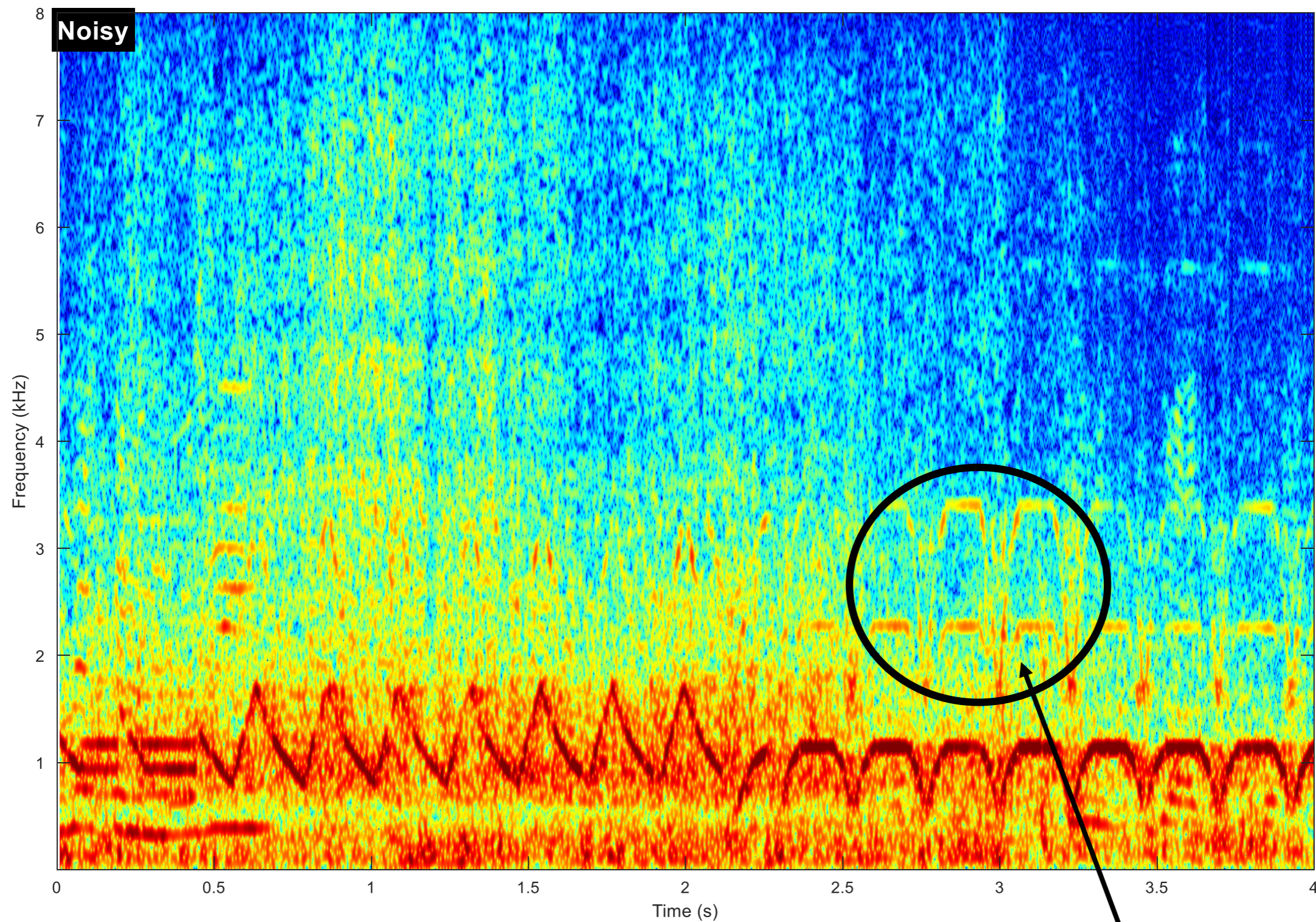
Level [dB]	Classroom, Hospital Home, Store	Trains, Airplanes	Restaurants
Speech SPL	60 to 70	60 to 70	60 to 70
Noise SPL	50 to 55	70 to 75	59 to 80
SNR	+5 to +20	-15 to 0	-20 to +11

*SPL: Sound Pressure Level relative to threshold of human hearing  
(20 micro-Pascals (force per square meter) ~mosquito flying 3m away)*

**Typical target range for speech  
enhancement: -5 to 15dB**

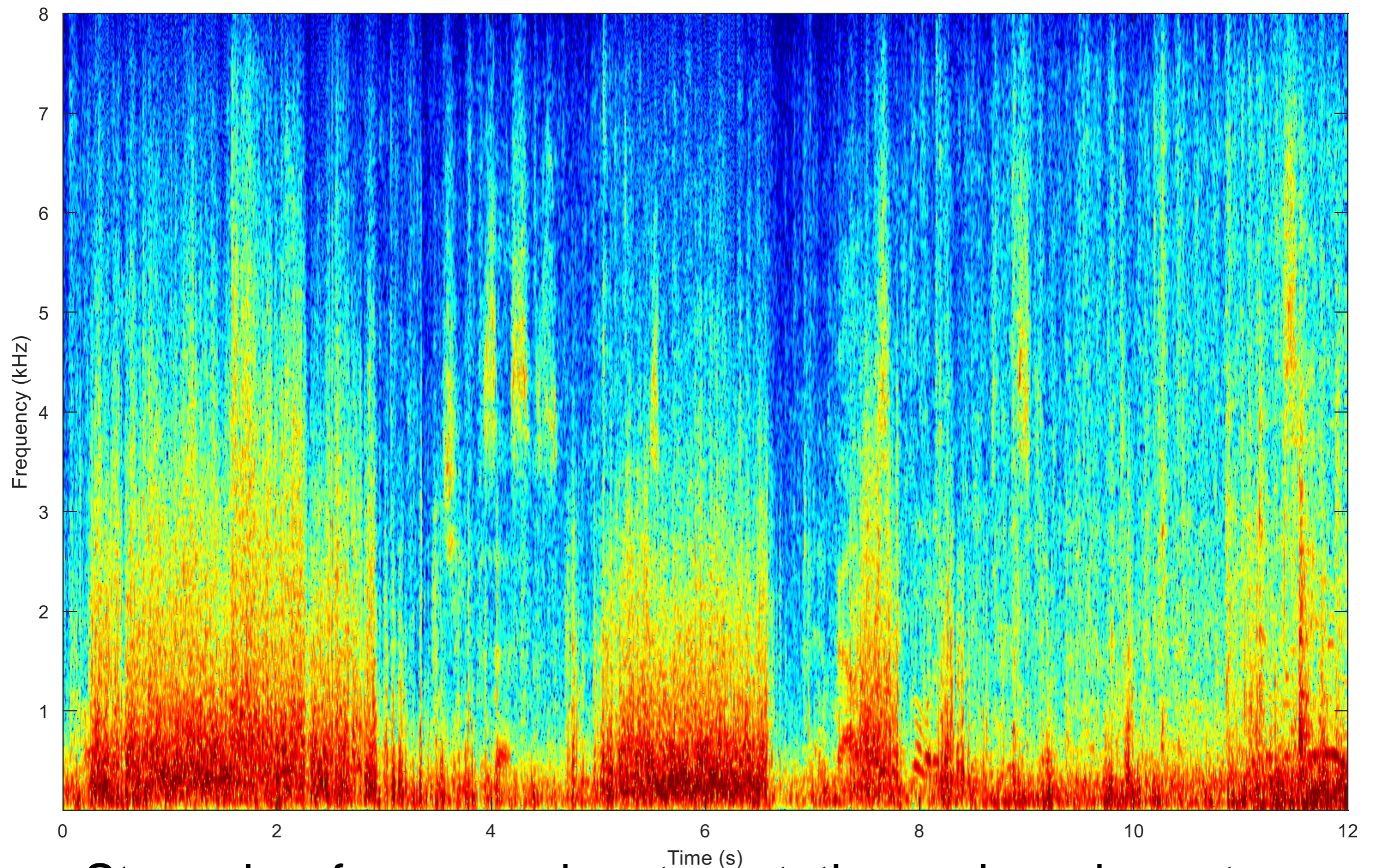


# Sirens



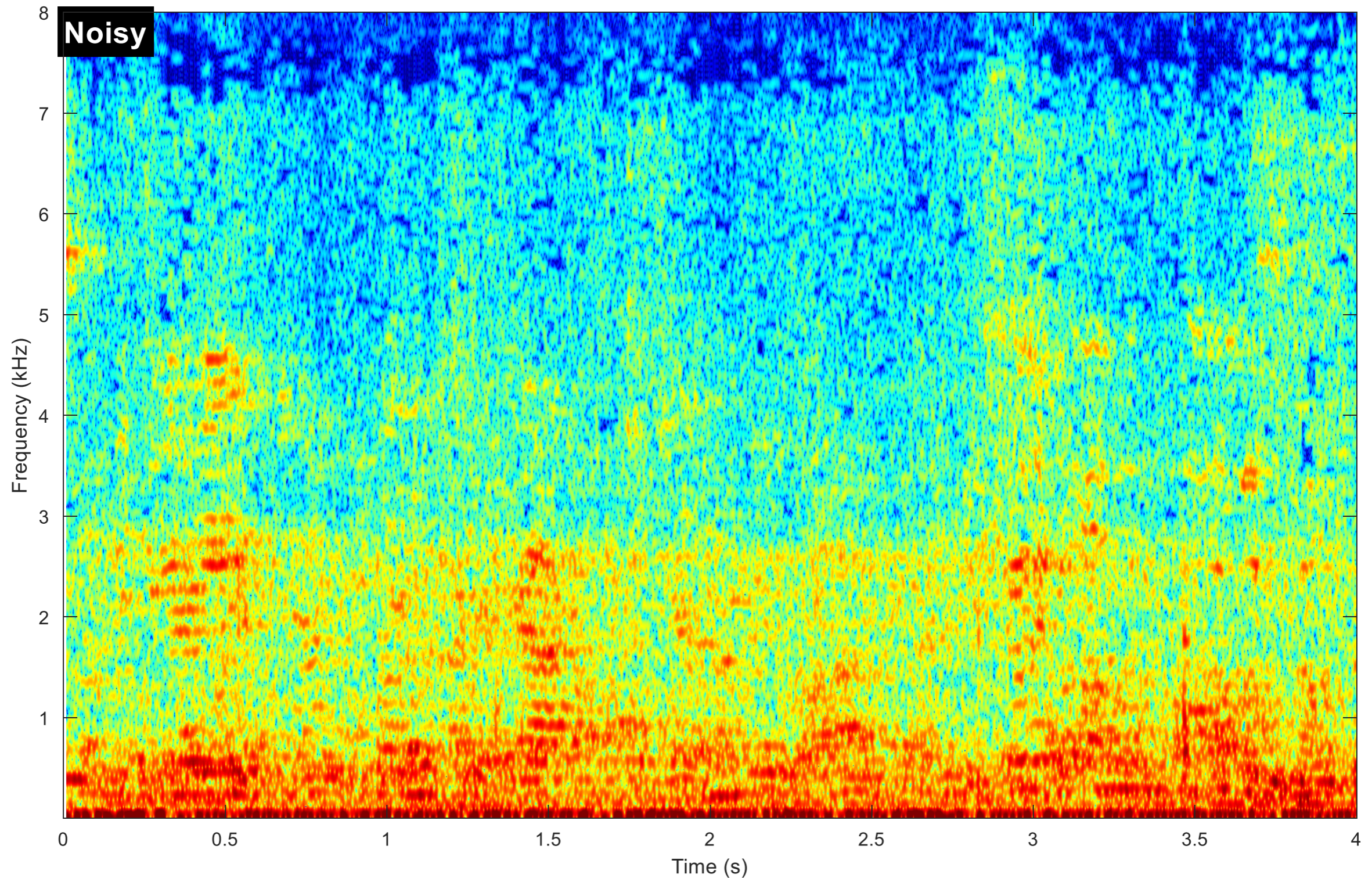
Strong, structured frequency modulated tones & overtones

# Wind Noise



Strong low frequency bursts + stationary broad spectrum

# Crowd



Broad, non-stationary spectrum in speech range

# Evaluating Performance of Speech Enhancers

- **Quality** measures assess **how** a speaker produces an utterance.
  - Is the utterance “natural”, “raspy”, “hoarse”, “scratchy”?
  - Does it sound good or bad?
- **Intelligibility** measures **what** a speaker said.
  - What did you understand?
  - What is the word error rate?

# Subjective Measures of Quality

## ITU-T P.835 Standard for Speech Enhancement Quality Assessment

Rating	Signal Distortion (SIG)	Background Distortion (BAK)	Overall Quality (OVL) Based on Mean Opinion Score Rating Scale (MOS)
5	Very natural, no degradation	Not noticeable	Excellent: Imperceptible
4	Fairly natural, little degradation	Somewhat noticeable	Good: Just perceptible, but not annoying
3	Somewhat natural, somewhat degraded	Noticeable but not intrusive	Fair: Perceptible and slightly annoying
2	Fairly unnatural, fairly degraded	Fairly conspicuous, somewhat intrusive	Poor: Annoying, but not objectionable
1	Very unnatural, very degraded	Very conspicuous, very intrusive	Bad: Very annoying and objectionable

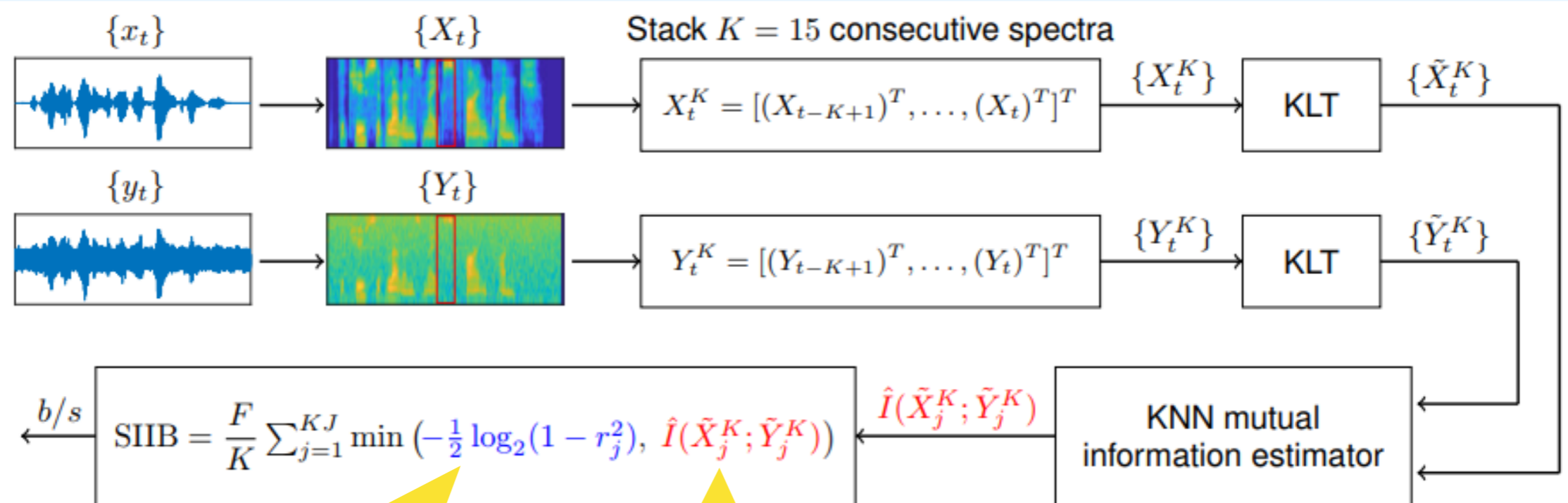
# Objective Measures of Quality and Intelligibility

Quality	Intelligibility
<p>Segmental SNR (SNRseg)                      Frequency Weighted Segmental SNR (fwSNRseg)                      Weighted Spectral Slope (WSS)                      Log-likelihood Ratio (LLR)                      Itakura-Saito (IS)                      Cepstral Distance (CEP)                      Hearing Aid Speech Quality Index (HASQI)                      Perceptual Evaluation of Video Quality (PEVQ)                      Perceptual Evaluation of Audio Quality (PEAQ)                      Perceptual Evaluation of Speech Quality (PESQ)  <b>Perceptual Objective Listening Quality Analysis (POLQA)</b>                      Composite Metrics</p>	<p>Normalized Covariance Metrics (NCM)                      Speech Intelligibility Index (SII)                      High-energy Glimpse Proportion Metric                      Coherence and Speech Intelligibility Index (CSII)                      Quasi-stationary Speech Transmission Index (QSTI)                      Short-time Objective Intelligibility Measure (STOI)                      Extended STOI Measure (ESTOI)                      Hearing-Aid Speech Perception Index (HASPI)                      K-Nearest Neighbor Mutual Information Intelligibility Measure (MIKNN)                      Speech Intelligibility Prediction based on a Mutual Information Lower Bound (SIMI)  <b>Speech Intelligibility in Bits (SIIB)</b>                      Speech-based Envelop Power Spectrum Model with Short-Time correlation (sEPSM)                      Automatic Speech Recognition (ASR)</p>

Effectiveness of metrics is evaluated by measuring correlation of metric predictions against subjective test data

# Speech Intelligibility in Bits (SIIB)

- Measures amount of information between speaker and listener.
- Linguistic models for “clean” speech communication measure 50-100 bps typical information rate.



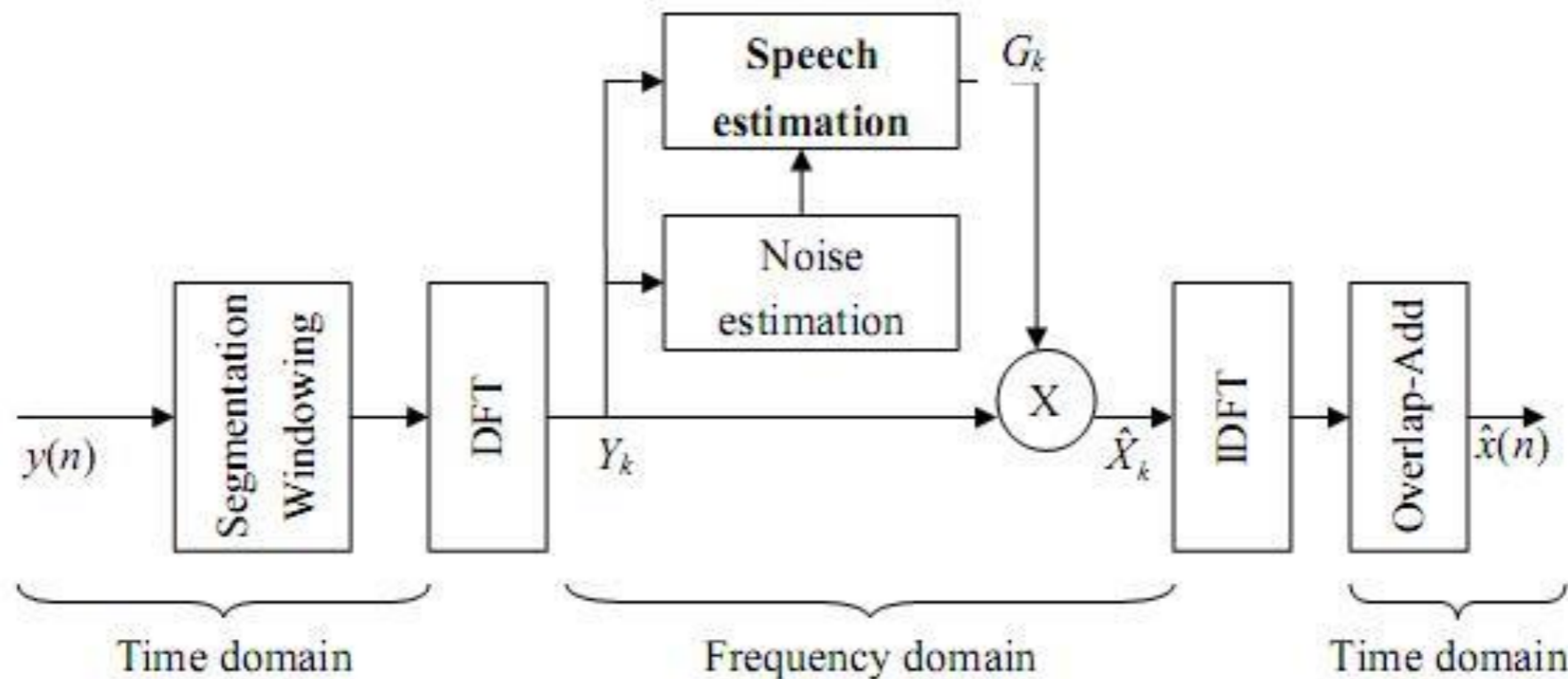
Mutual Info between  
“text” message and  
clean speech

Mutual Info  
between clean  
and noise  
speech

From: S. Van Kuyk; W. B. Kleijn; R. C. Hendriks; “An instrumental intelligibility metric based on information theory,” in *IEEE Signal Processing Letters*, 2018

# Traditional Methods of Speech Enhancement

- Most commonly employ a short-time Fourier transform based analysis-modification-synthesis framework
- Frequency dependent noise suppression function
- Noises suppression based on estimates of speech and noise statistics

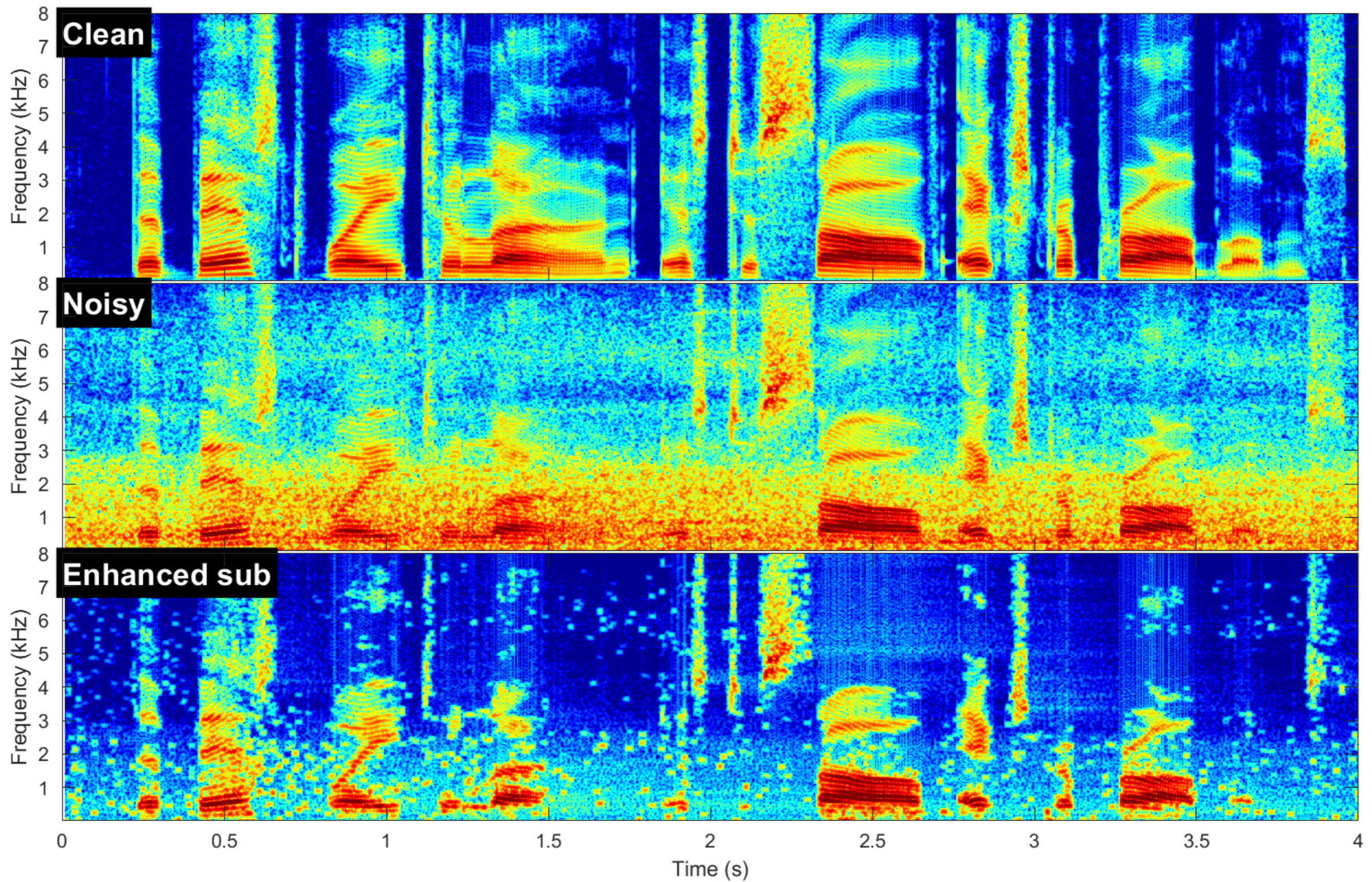




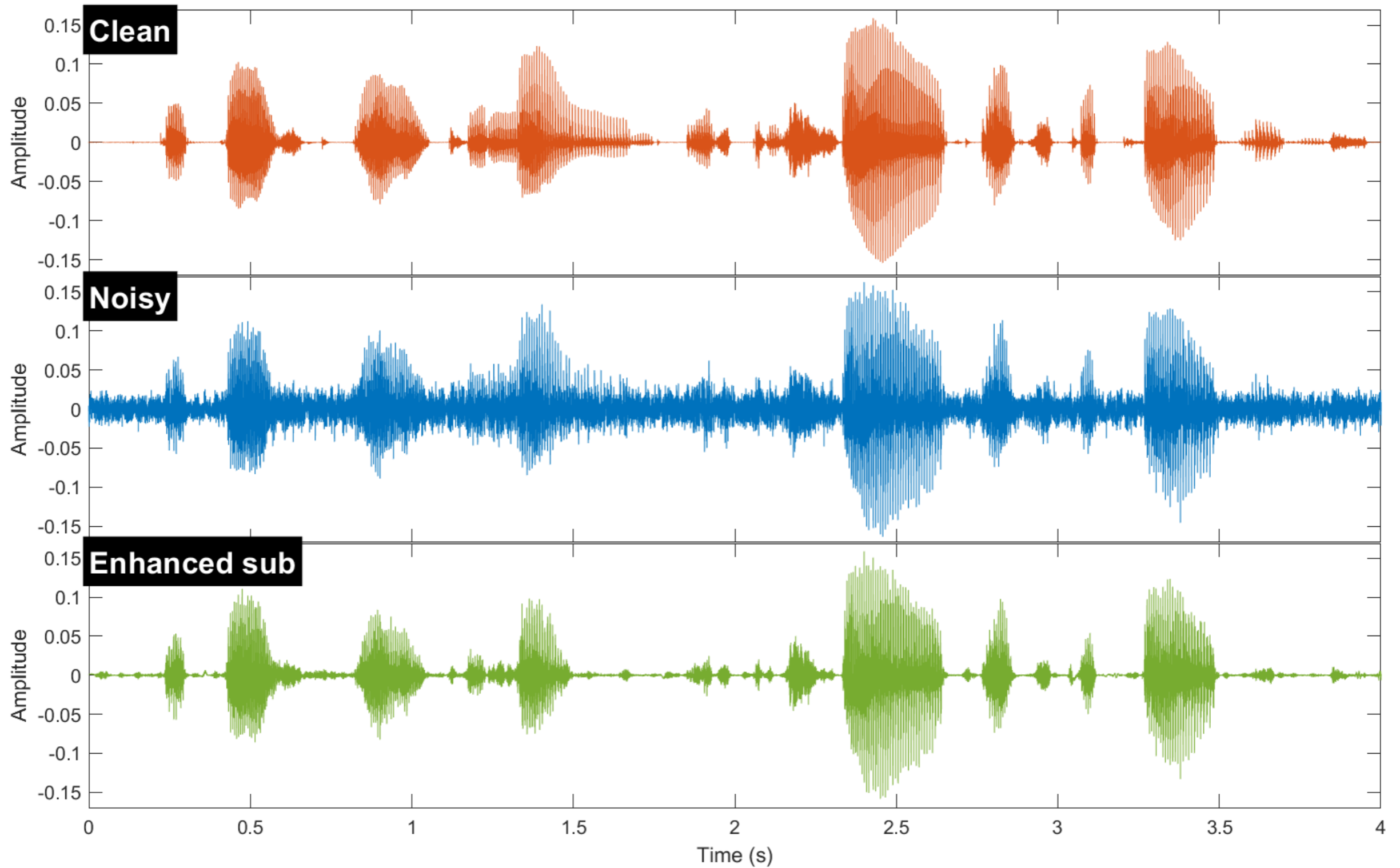
# Traditional Methods: Spectral Subtraction

$\underbrace{R(\omega)}_{\substack{\text{noisy} \\ \text{speech}}} = \underbrace{S(\omega)}_{\substack{\text{clean} \\ \text{speech}}} + \underbrace{D(\omega)}_{\text{noise}}$	Noisy speech model
$ \hat{D} ^2 = E\{ R ^2\} = E\{ \hat{D} ^2\}$ <p>when <math>S = 0</math></p>	<b>Noise magnitude estimate</b> measured during period of speech inactivity using Voice Activity Detector
$ R ^2 =  S ^2 +  D ^2 + \underbrace{2\text{Re}\{SD^*\}}_{\substack{\text{ignore} \\ \text{this term!!}}}$	Noisy speech magnitude Cross term is ignored because clean speech and noise are uncorrelated
$ \hat{S} ^2 =  R ^2 -  \hat{D} ^2$	<b>Clean speech magnitude estimate</b>
$\hat{S}(\omega) = \underbrace{ \hat{S}(\omega) }_{\substack{\text{clean} \\ \text{magnitude} \\ \text{estimate}}} \exp\left\{j \underbrace{\Phi_r(\omega)}_{\substack{\text{noisy} \\ \text{phase}}}\right\}$	<b>Clean speech synthesized</b> from noisy phase and magnitude estimate Difference in noisy and clean phase not perceptible for SNRs > 8dB

# Spectral Subtraction: Spectrograms



# Spectral Subtraction: Waveforms



# Deep Neural Networks for Speech Enhancement

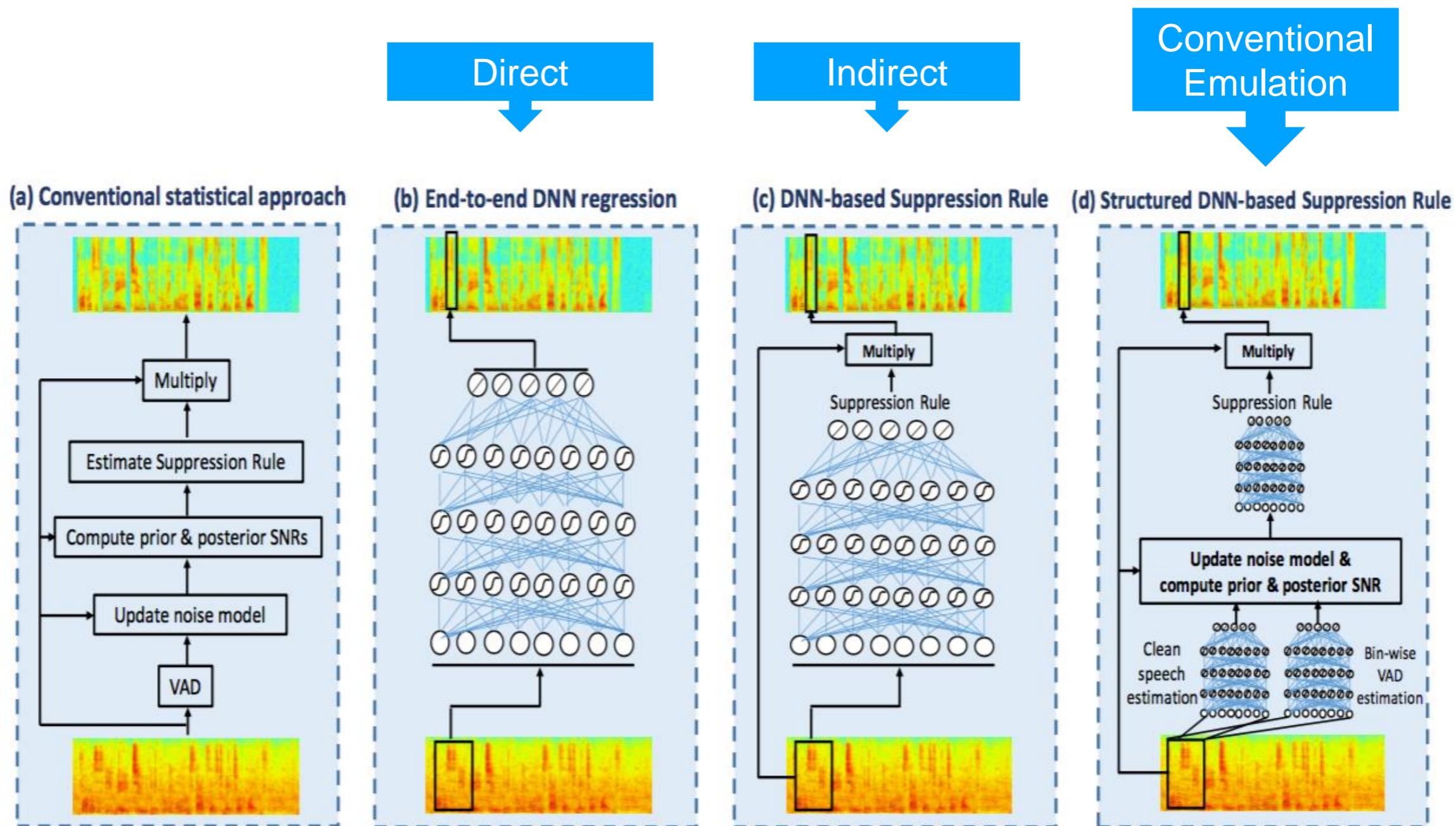


Figure 1: Different architectures for DNN-based speech enhancement: (a) Statistical noise suppression. (b) Enhancement by end-to-end DNN regression from noisy spectral features to clean features. (c) Estimating binwise suppression gain directly by a DNN. (d) Employing separate DNNs replacing the different components of conventional suppression gain estimation.

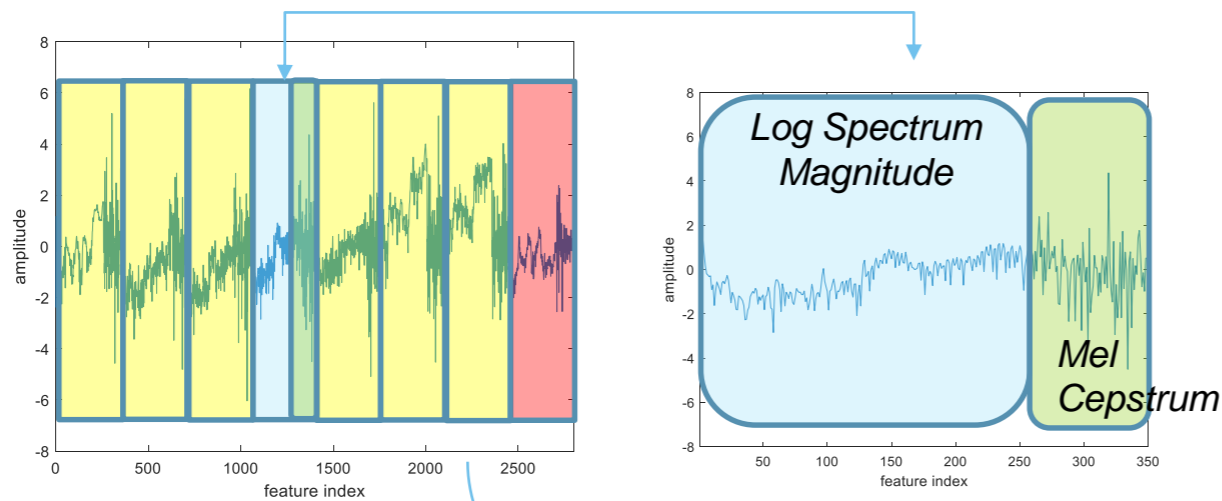
Mirsamadi, Seyedmahdad, and Ivan Tashev. "Causal Speech Enhancement Combining Data-Driven Learning and Suppression Rule Estimation." INTERSPEECH. 2016.

# Common Ideal Target Masks

target mask/filter	formula	optimality principle
IBM:	$a^{\text{ibm}} = \delta( s  >  n ),$	max SNR $a \in \{0, 1\}$
IRM:	$a^{\text{irm}} = \frac{ s }{ s  +  n },$	max SNR $\theta_s = \theta_n,$
“Wiener like”:	$a^{\text{wf}} = \frac{ s ^2}{ s ^2 +  n ^2},$	max SNR, expected power
ideal amplitude:	$a^{\text{iaf}} =  s / y ,$	exact $ \hat{s} ,$ max SNR $\theta_s = \theta_y$
phase-sensitive filter:	$a^{\text{psf}} = \frac{ s }{ y } \cos(\theta),$	max SNR given $a \in \mathbb{R}$
ideal complex filter:	$a^{\text{icf}} = s/y,$	max SNR given $a \in \mathbb{C}$

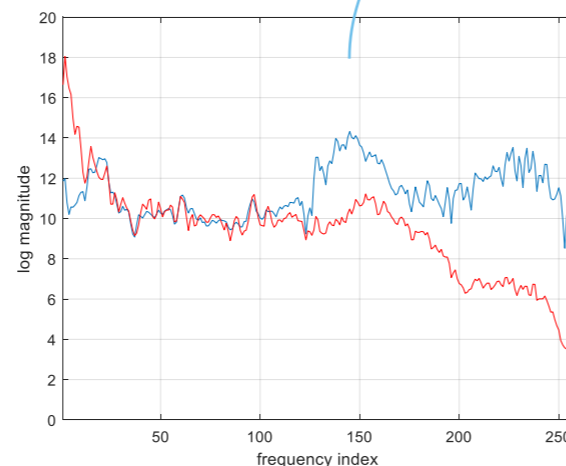
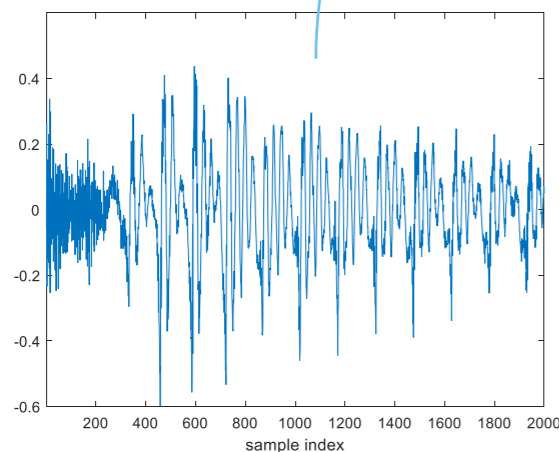
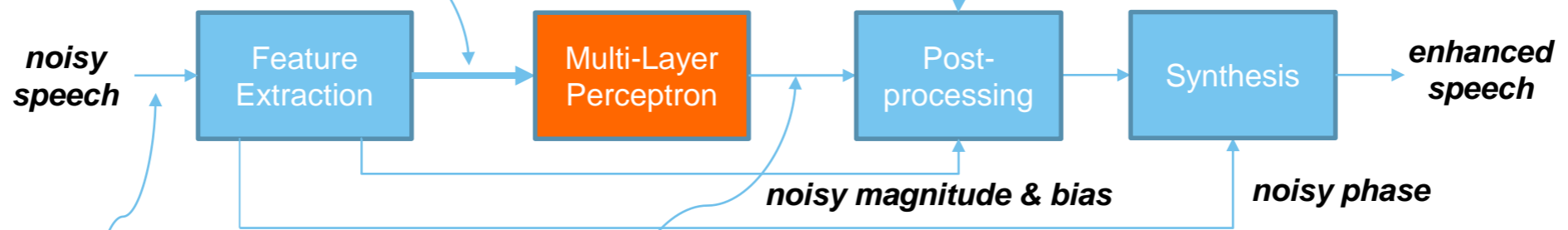
H. Erdogan, J. R. Hershey, S. Watanabe, and J. L. Roux, “Phase-sensitive and recognition-boosted speech separation using deep recurrent neural networks,” in *Proc. Int. Conf. Acoust., Speech, Signal Process.*, 2015, pp. 708–712

# GeorgiaTech System



**DNN Input Features:** 7x noisy speech frames + 1 frame noise only of concatenated Log Spectrum + Mel Cepstrum with Global mean removed & normalized by Global variance

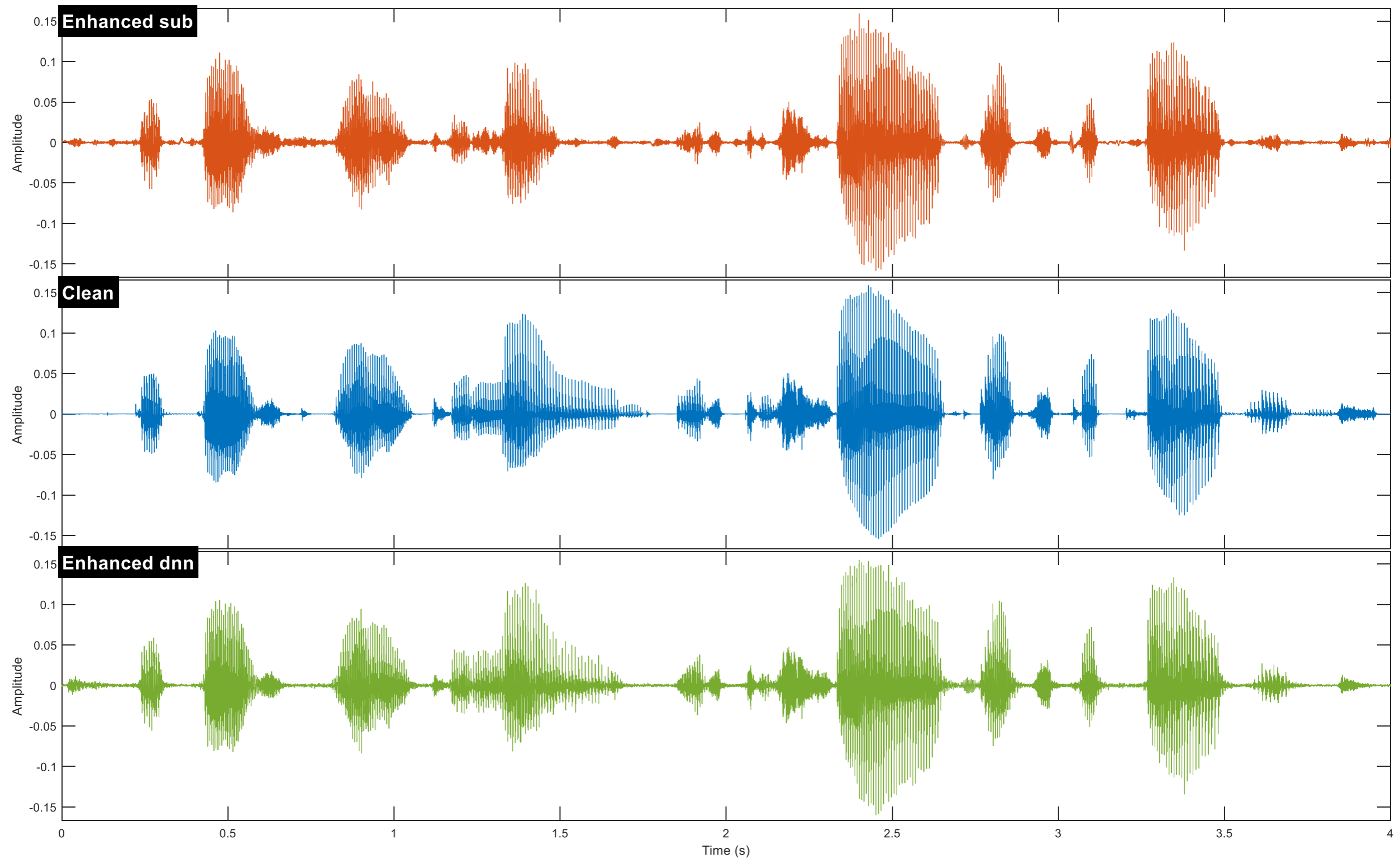
- Derive IRM from speech and noise spectrum estimates
- Mix DNN output with bias & noisy magnitude according to IRM



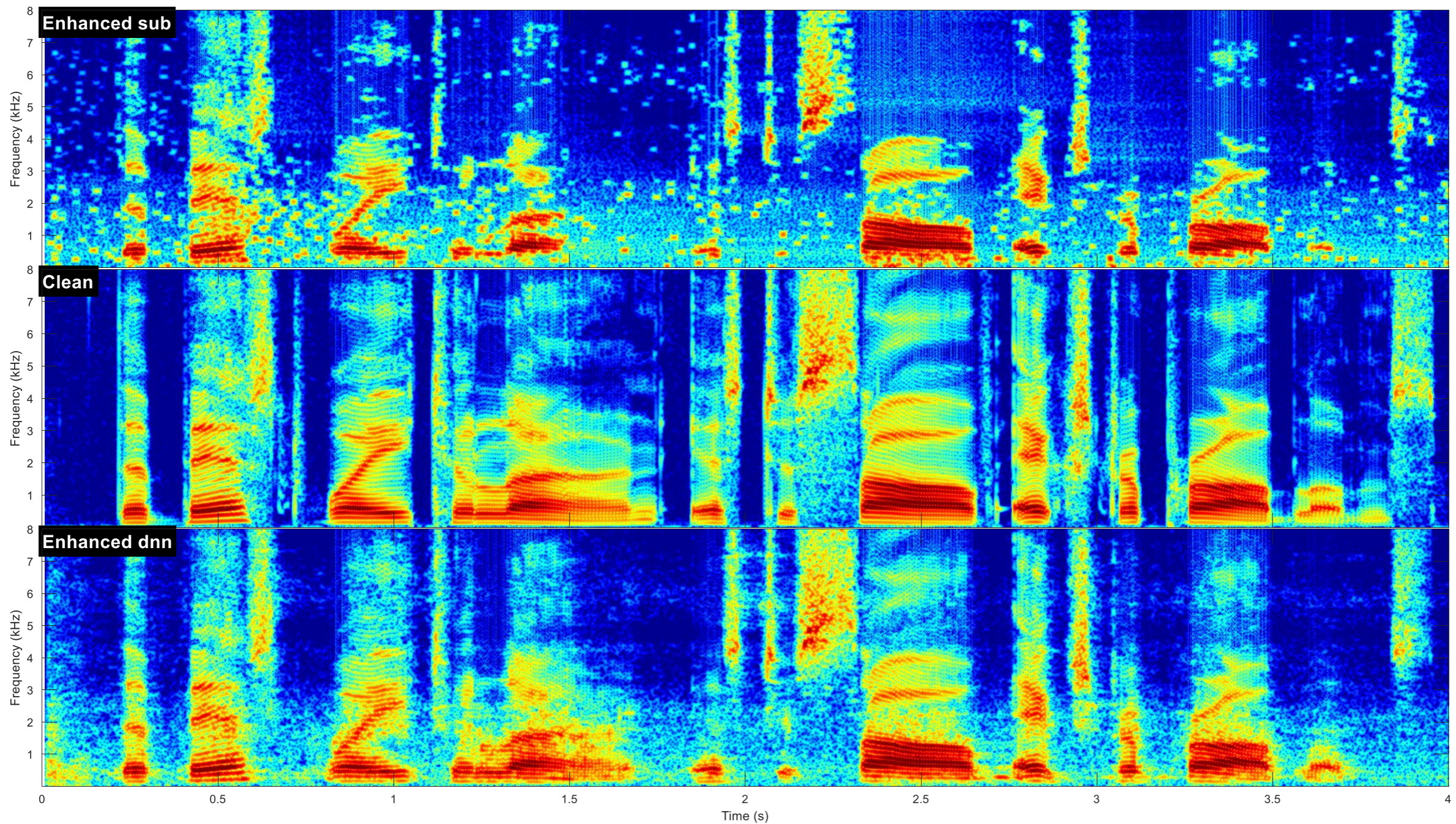
**Speech + Noise Log Spectrum**

Xu, Yong, et al. "A regression approach to speech enhancement based on deep neural networks." *IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP)* 23.1 (2015): 7-19.

# Spectral Subtractive vs. BabbleLabs DNN



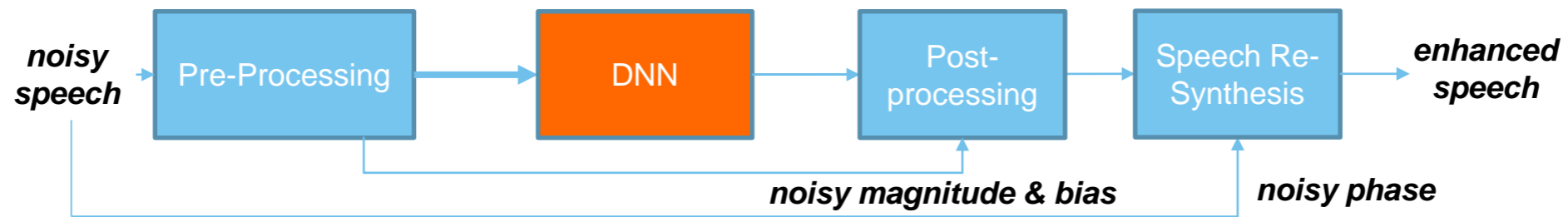
# Spectral Subtractive vs. BabbleLabs DNN



Metric	Noisy	Subtractive	BabbleLabs DNN
SNR [dB]	5.98	10.53	13.57
PESQ	1.18	1.42	2.01
ESTOI	0.44	0.49	0.71
SIIB_Gauss [bps]	65	56	93



# BabbleLabs Production Flow



- 90% of the code in the blue boxes
- 90% of the compute in the orange box
- Prototyping is in blocking format, while deployment is in streaming format.
- Using Matlab and the GPU coder, we were able to convert from reference to deployment code in 6 man-weeks.
- Currently we are porting the DNN using other open source tools.
  - Exploring the migration to GPU coder to unify the flow if possible.

# References

- Loizou, Philipos C. *Speech enhancement: theory and practice*. CRC press, 2007.
- Van Kuyk, Steven, W. Bastiaan Kleijn, and Richard C. Hendriks. "An evaluation of intrusive instrumental intelligibility metrics." arXiv preprint arXiv:1708.06027 (2017).
- Mirsamadi, Seyedmahdad, and Ivan Tashev. "Causal Speech Enhancement Combining Data-Driven Learning and Suppression Rule Estimation." INTERSPEECH. 2016.
- Xu, Yong, et al. "A regression approach to speech enhancement based on deep neural networks." *IEEE/ACM Transactions on Audio, Speech and Language Processing (TASLP)* 23.1 (2015): 7-19.
- H. Erdogan, J. R. Hershey, S. Watanabe, and J. L. Roux, "Phase-sensitive and recognition-boosted speech separation using deep recurrent neural networks," in Proc. Int. Conf. Acoust., Speech, Signal Process., 2015, pp. 708–712
- From <http://www.vision.huji.ac.il/visual-speech-enhancement/>
- <https://looking-to-listen.github.io/>



s p e a k   y o u r   m i n d