

My 5 minute talk

about 2 of my ongoing projects

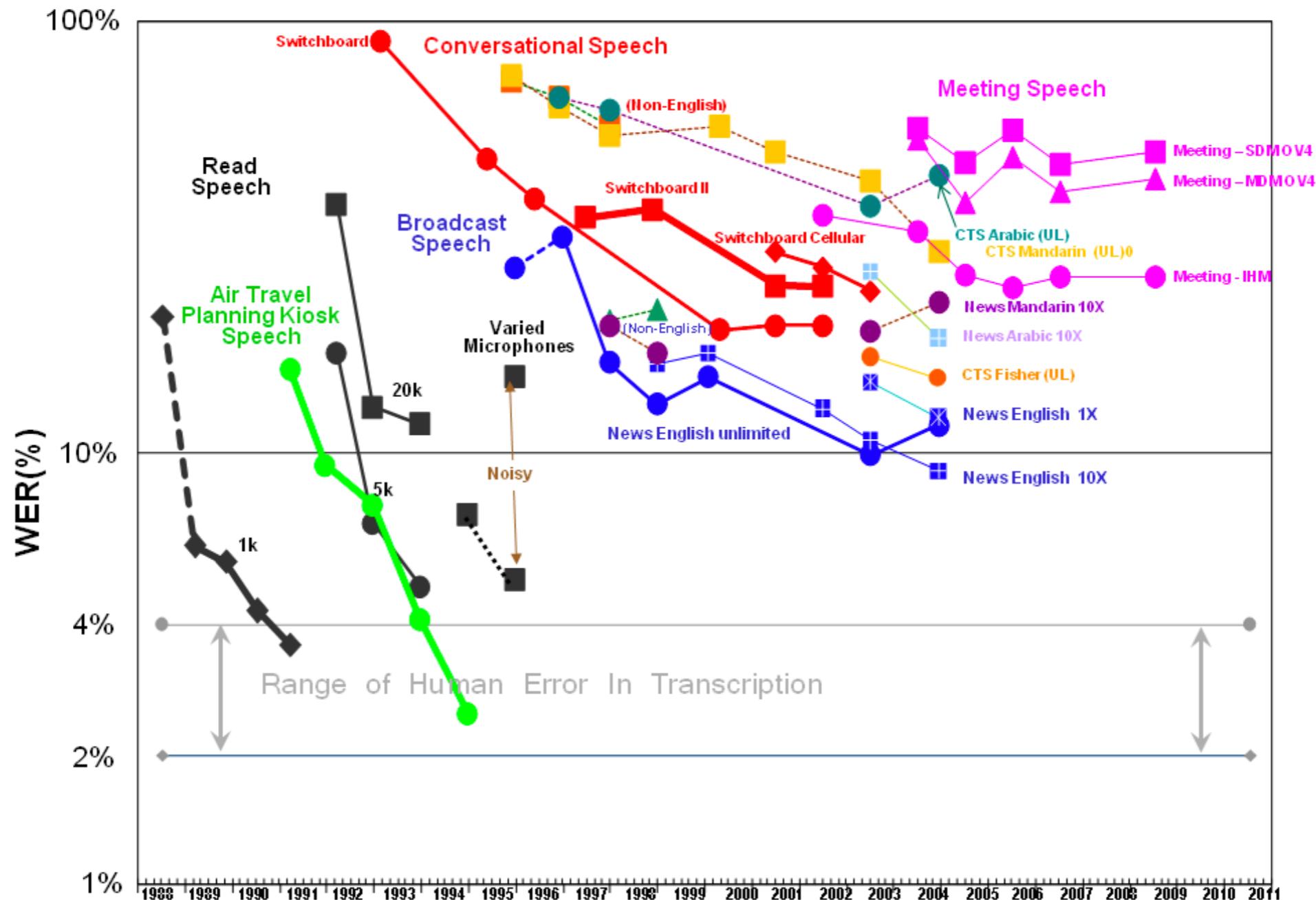
Desh Raj

Noise-aware Training of Acoustic Models

Noise-aware Training of Acoustic Models

History

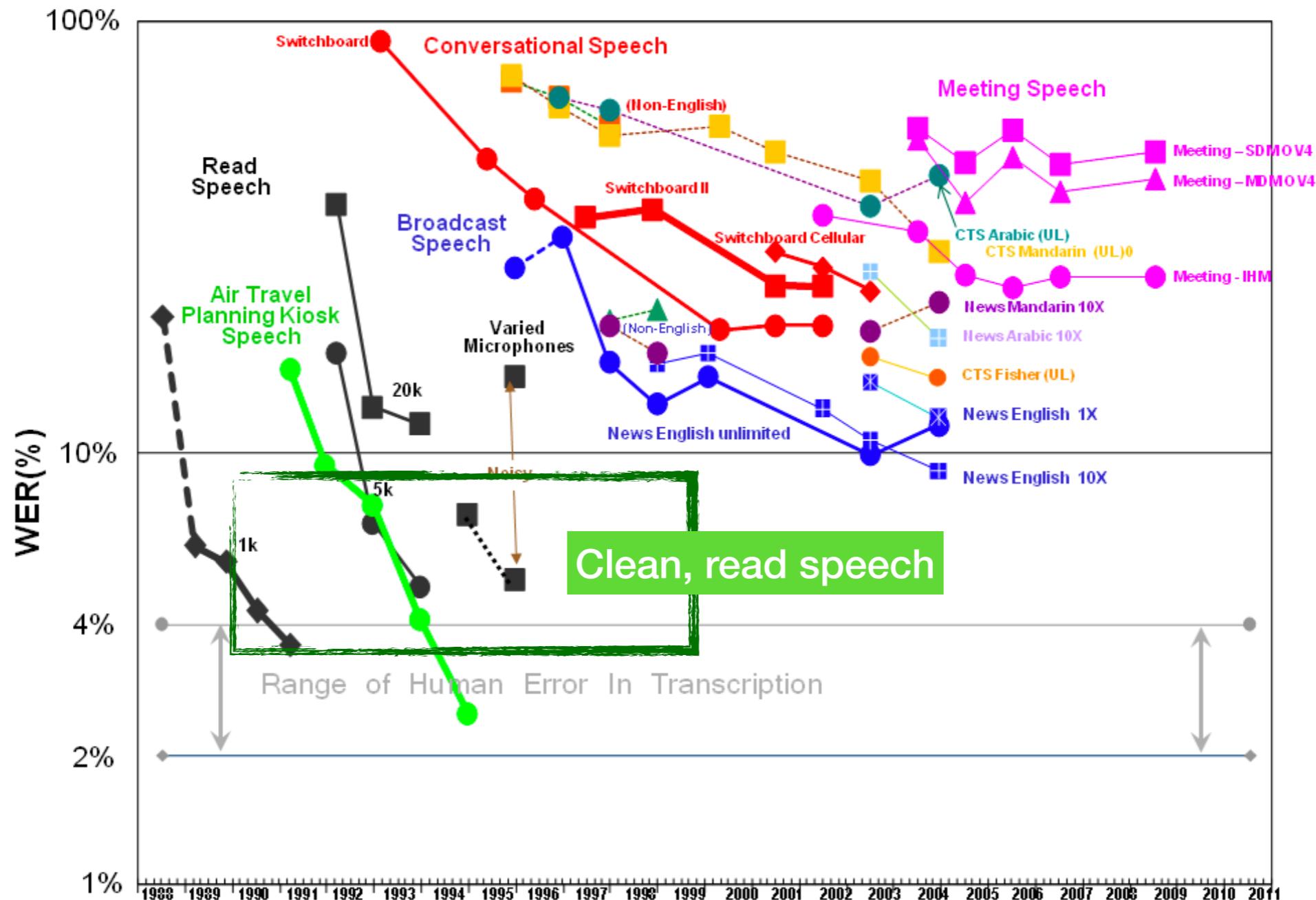
NIST STT Benchmark Test History – May. '09



Noise-aware Training of Acoustic Models

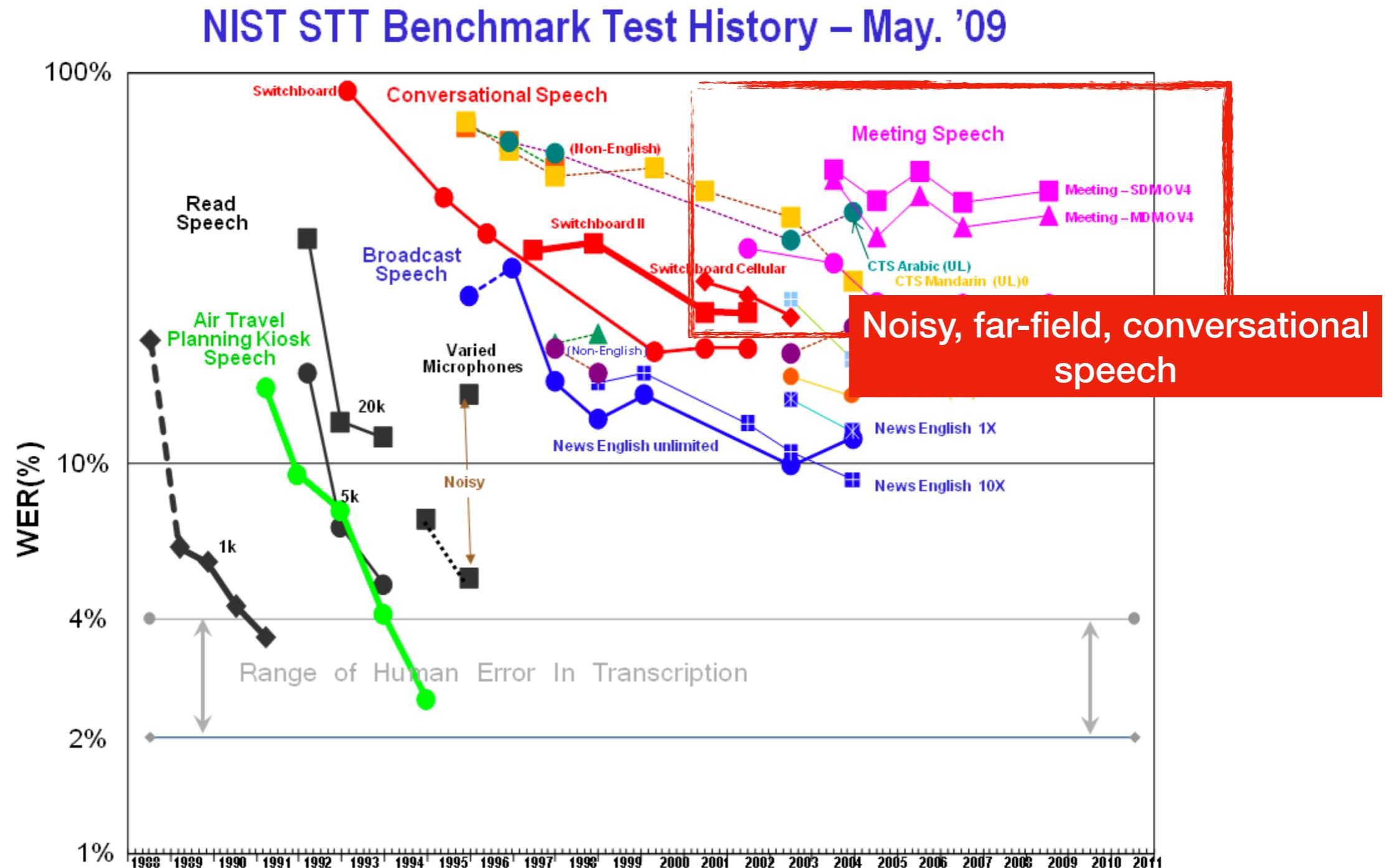
History

NIST STT Benchmark Test History – May. '09



Noise-aware Training of Acoustic Models

History



Noise-aware Training of Acoustic Models

The Problem

- Speaker effects
- Background noise
- Reverberation

Noise-aware Training of Acoustic Models

The Problem

- Speaker effects
- Background noise
- Reverberation



Noise-aware Training of Acoustic Models

Well, actually...

VERY DEEP CONVOLUTIONAL NEURAL NETWORKS FOR ROBUST SPEECH RECOGNITION

Yanmin Qian^{1,2}, Philip C Woodland²

¹ Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai, China

²Cambridge University Engineering Department, Cambridge CB2 1PZ, UK

yanminqian@sjtu.edu.cn, pcw@eng.cam.ac.uk

Noise-aware Training of Acoustic Models

Well, actually...

VERY DEEP CONVOLUTIONAL NEURAL NETWORKS FOR ROBUST SPEECH RECOGNITION

¹ Department of

²C

Very Deep Self-Attention Networks for End-to-End Speech Recognition

Ngoc-Quan Pham¹, Thai-Son Nguyen¹, Jan Niehues¹, Markus Müller¹, Sebastian Stüker¹, Alex Waibel^{1,2}

¹Interactive Systems Lab, Karlsruhe Institute of Technology, Karlsruhe, Germany

²Carnegie Mellon University, Pittsburgh PA, USA

ngoc.pham@kit.edu, thai.nguyen@kit.edu

Noise-aware Training of Acoustic Models

Well, actually...

VERY DEEP CONVOLUTIONAL NEURAL NETWORKS FOR ROBUST SPEECH RECOGNITION

1

Untangling in Invariant Speech Recognition

Cory Stephenson
Intel AI Lab
cory.stephenson@intel.com

Jenelle Feather
MIT
jfeather@mit.edu

Suchismita Padhy
Intel AI Lab
suchismita.padhy@intel.com

Oguz Elibol
Intel AI Lab
oguz.h.elibol@intel.com

Hanlin Tang
Intel AI Lab
hanlin.tang@intel.com

Josh McDermott
MIT/ Center for Brains, Minds, and Machines
jhm@mit.edu

SueYeon Chung
Columbia University/ MIT
sueyeon@mit.edu

Speech Recognition

Sebastian Stüker¹, Alex

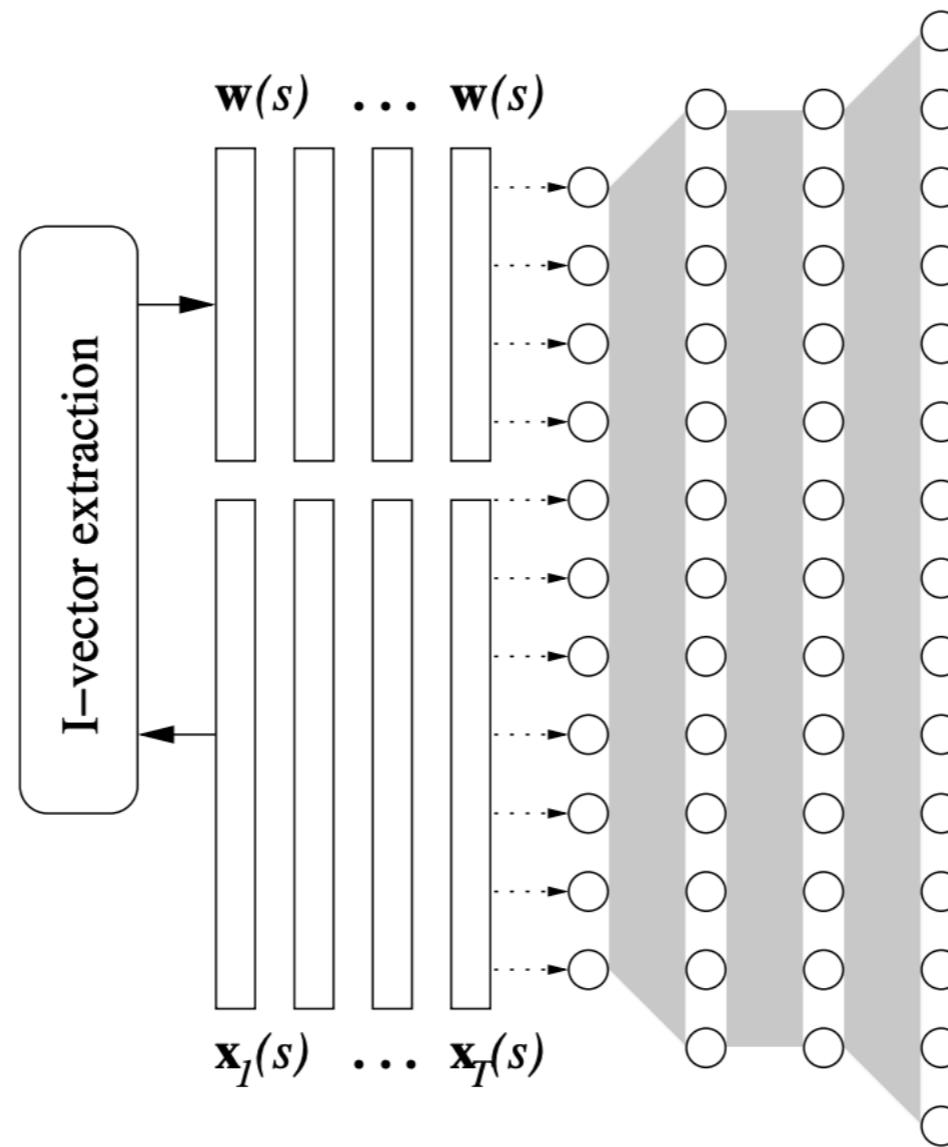
rlsruhe, Germany
A

“But we don’t have enough data or GPUs.”

- Every DL practitioner not working at Froogle

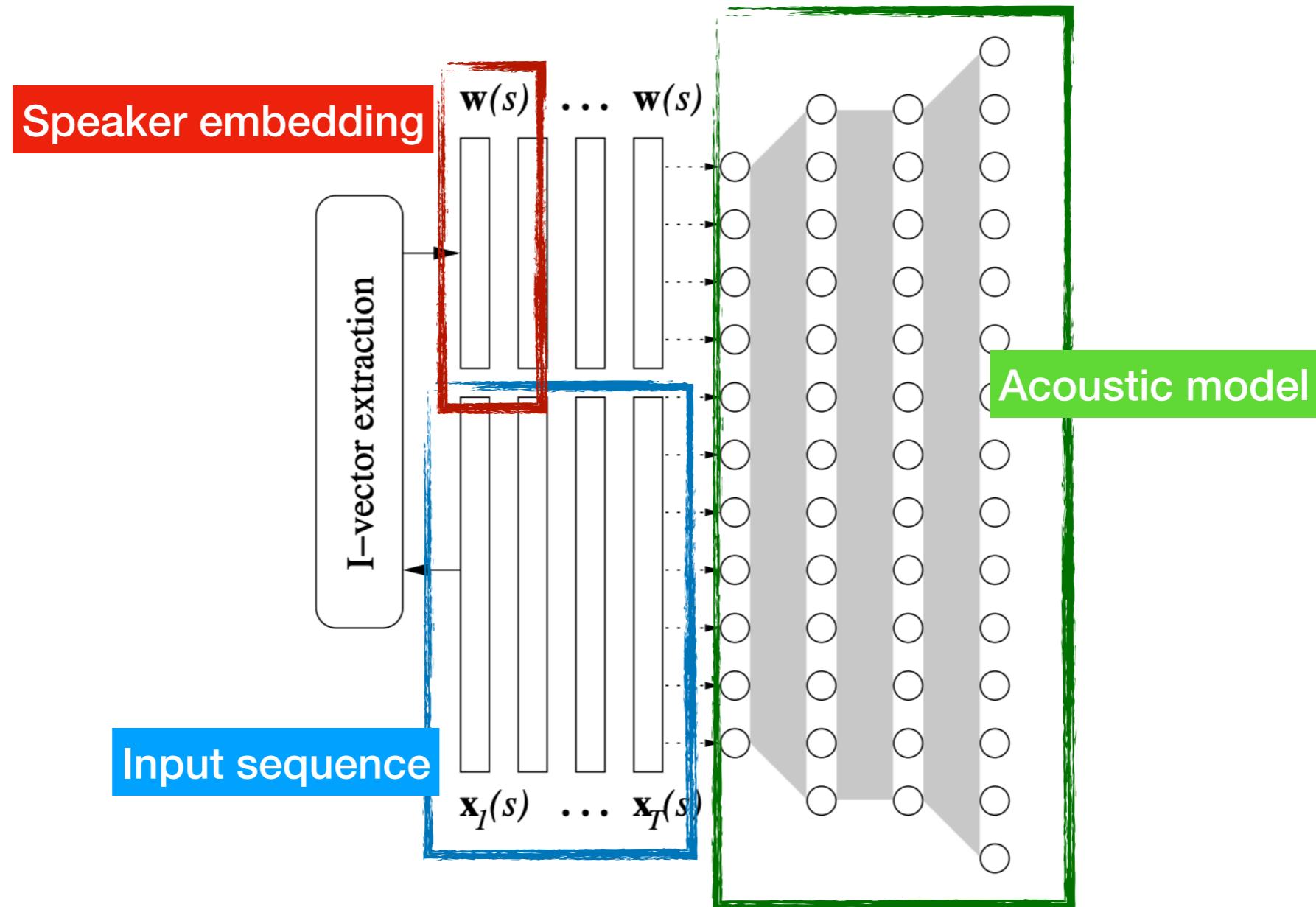
Noise-aware Training of Acoustic Models

An idea from speaker adaptation



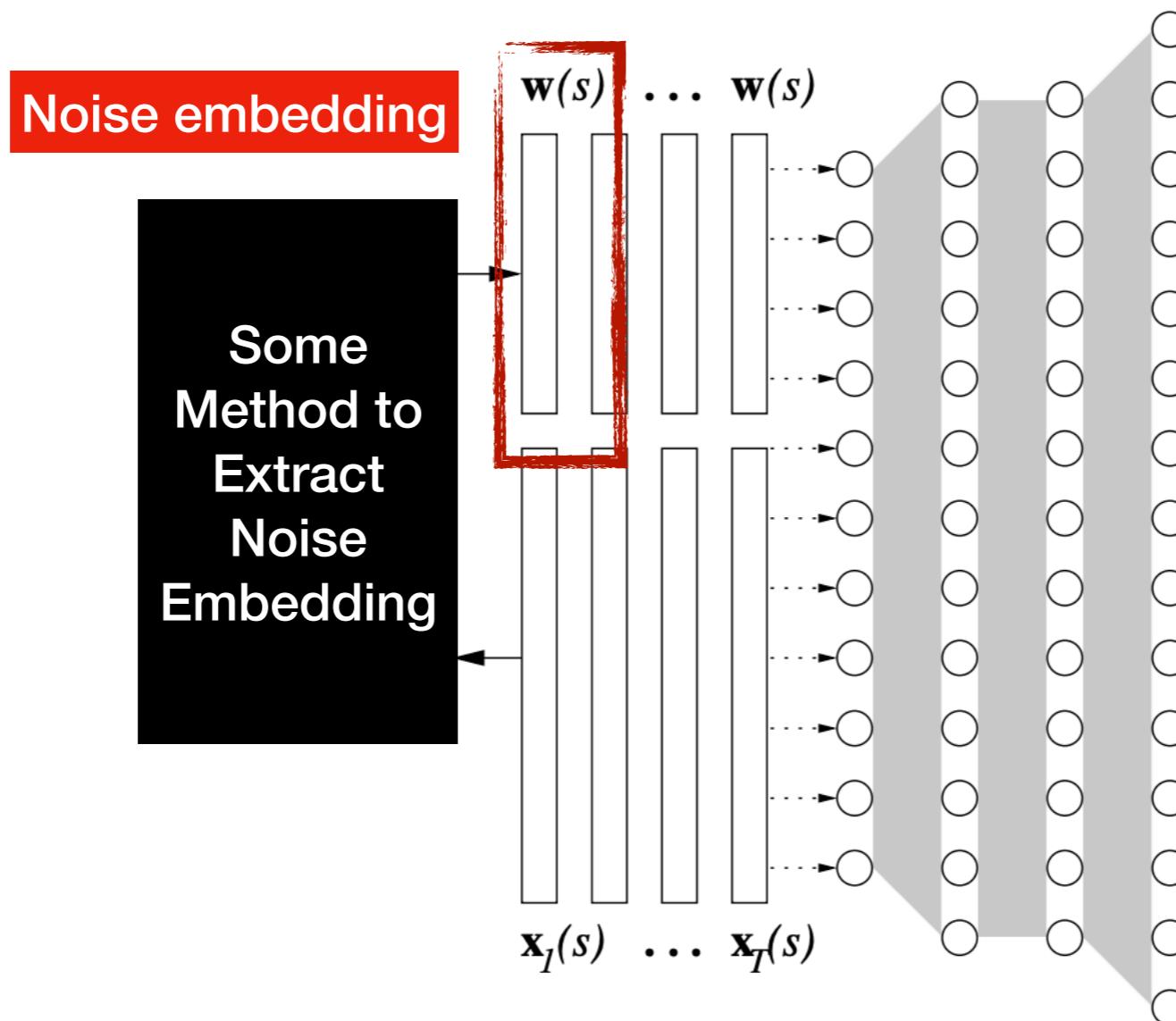
Noise-aware Training of Acoustic Models

An idea from speaker adaptation



Noise-aware Training of Acoustic Models

An idea from speaker adaptation



Noise-aware Training of Acoustic Models

Where to get noise embeddings?

Noise-aware Training of Acoustic Models

Where to get noise embeddings?

ASRU 2019

PROBING THE INFORMATION ENCODED IN X-VECTORS

Desh Raj, David Snyder, Daniel Povey, Sanjeev Khudanpur

Center for Language and Speech Processing & Human Language Technology Center of Excellence
The Johns Hopkins University, Baltimore, MD 21218, USA.

`draj@cs.jhu.edu, {david.ryan.snyder, dpovey}@gmail.com, khudanpur@jhu.edu`

From Google



Noise-aware Training of Acoustic Models

Where to get noise embeddings?

Discriminatively trained speaker embeddings

ASRU 2019

PROBING THE INFORMATION ENCODED IN X-VECTORS

Desh Raj, David Snyder, Daniel Povey, Sanjeev Khudanpur

Center for Language and Speech Processing & Human Language Technology Center of Excellence
The Johns Hopkins University, Baltimore, MD 21218, USA.

`draj@cs.jhu.edu, {david.ryan.snyder, dpovey}@gmail.com, khudanpur@jhu.edu`



Noise-aware Training of Acoustic Models

Where to get noise embeddings?

ASRU 2019

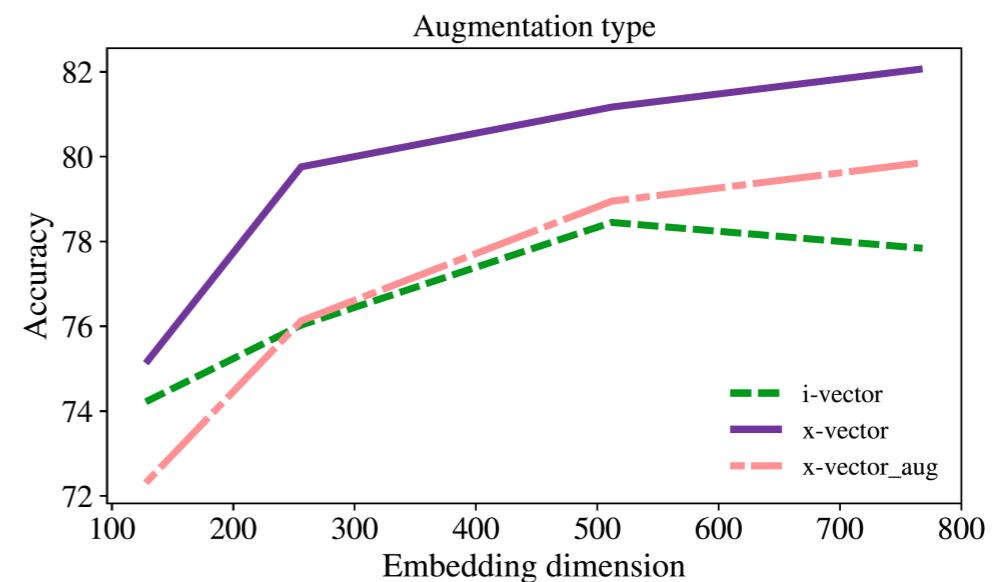
PROBING THE INFORMATION ENCODED IN X-VECTORS

Desh Raj, David Snyder, Daniel Povey, Sanjeev Khudanpur

Center for Language and Speech Processing & Human Language Technology Center of Excellence
The Johns Hopkins University, Baltimore, MD 21218, USA.

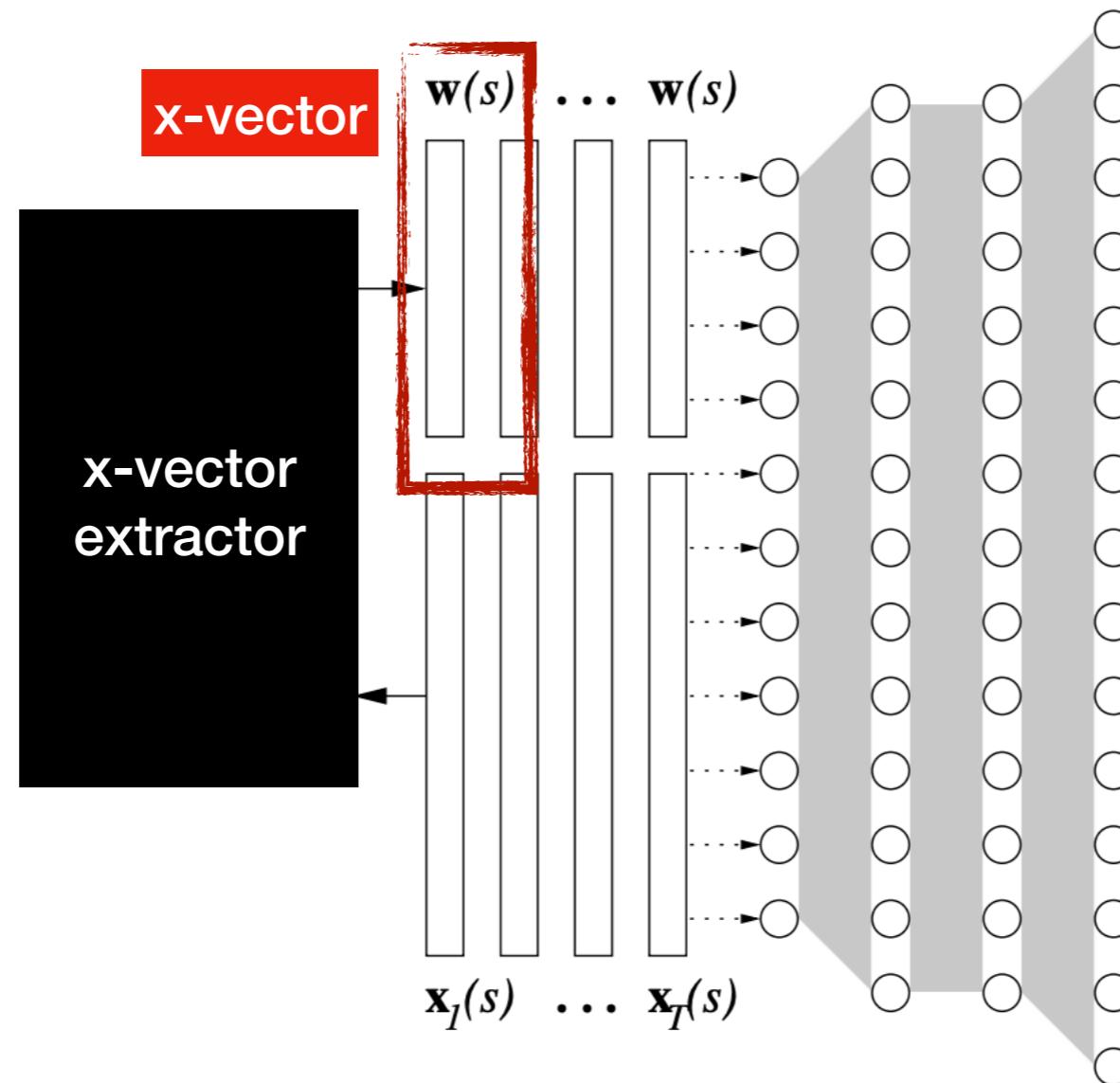
`draj@cs.jhu.edu, {david.ryan.snyder, dpovey}@gmail.com, khudanpur@jhu.edu`

But also contain noise information if
trained without augmentation



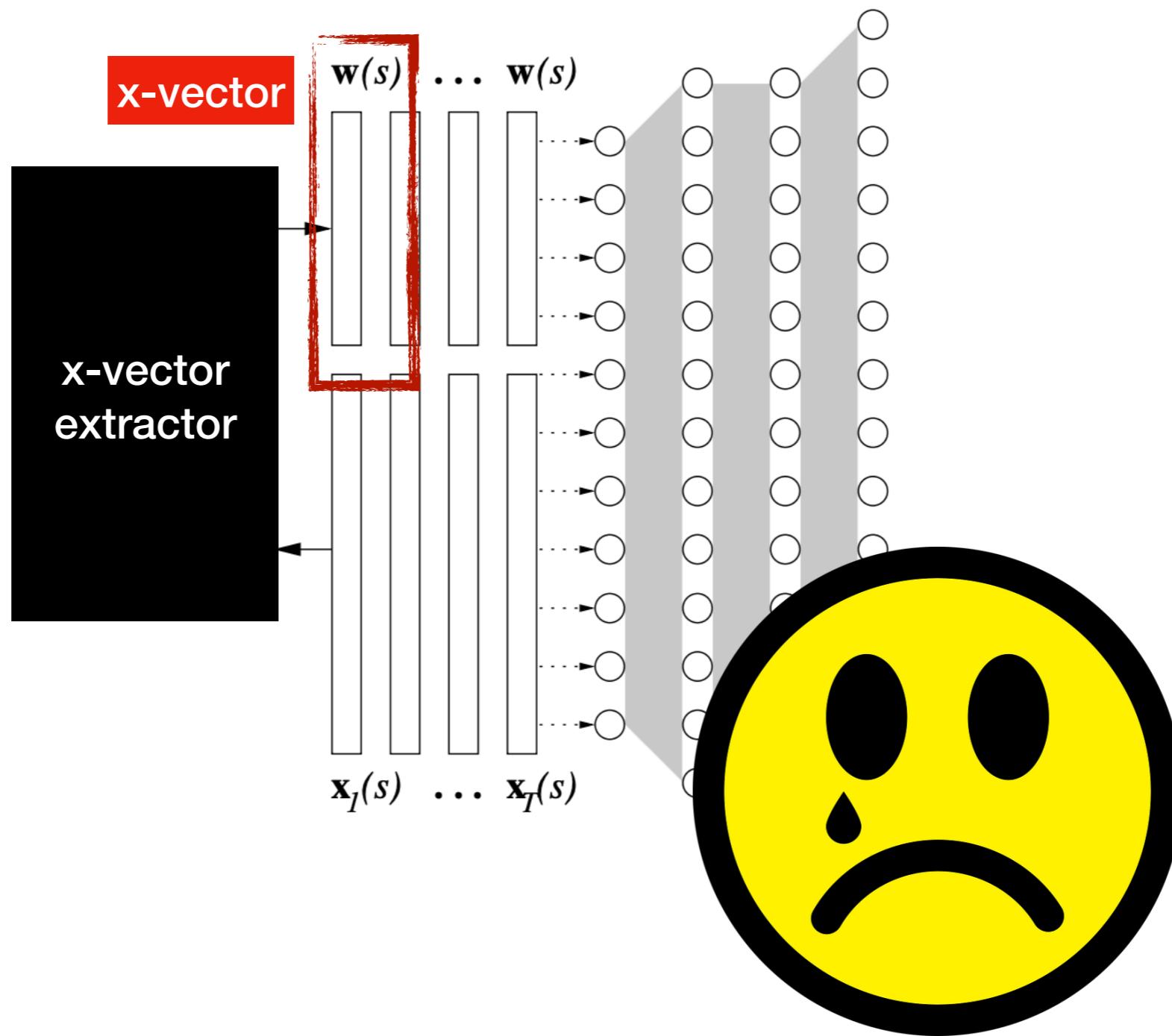
Noise-aware Training of Acoustic Models

Using x-vectors



Noise-aware Training of Acoustic Models

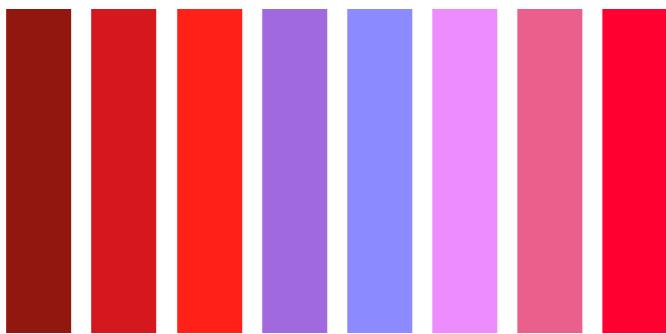
Using x-vectors



Noise-aware Training of Acoustic Models

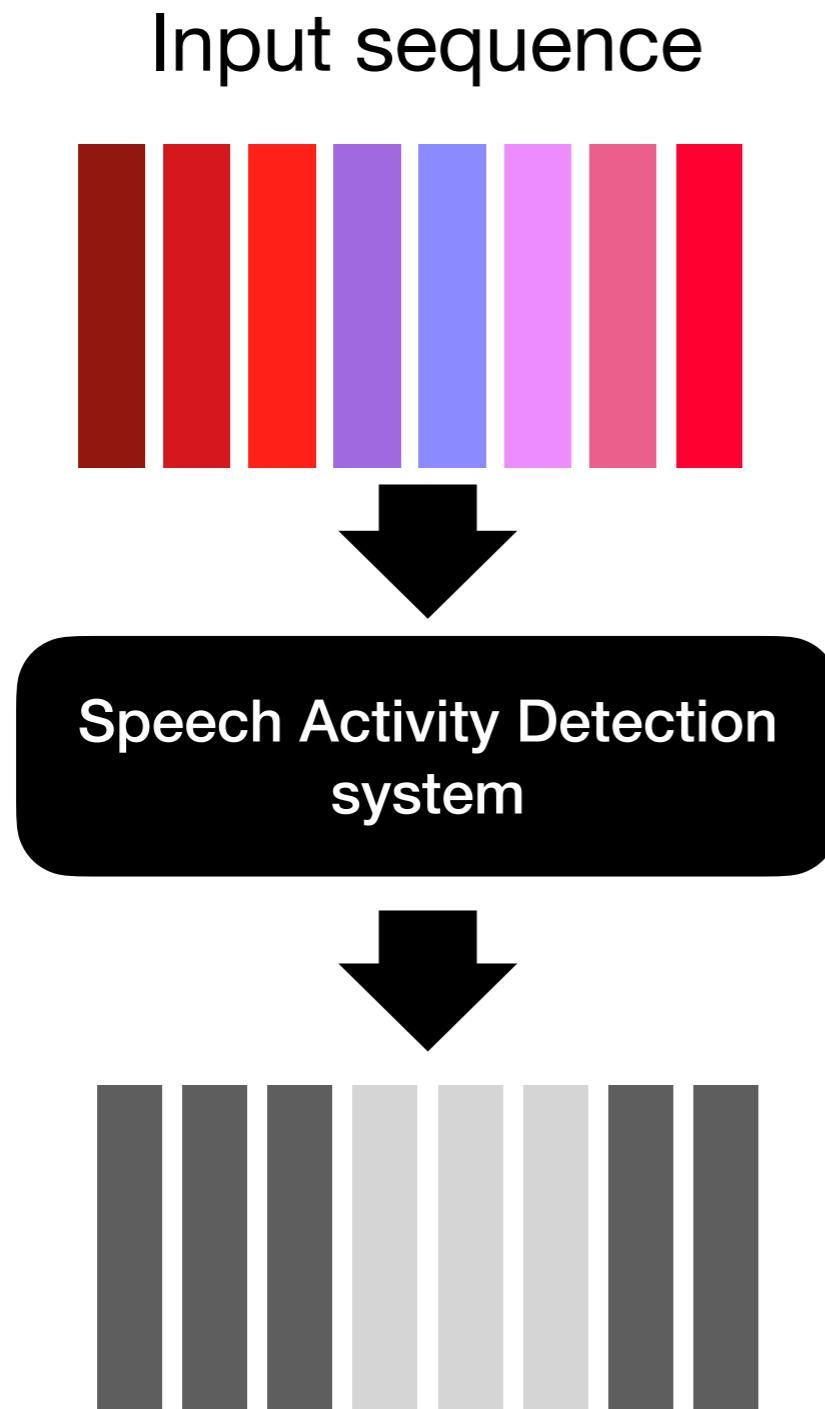
Speech/silence classification

Input sequence



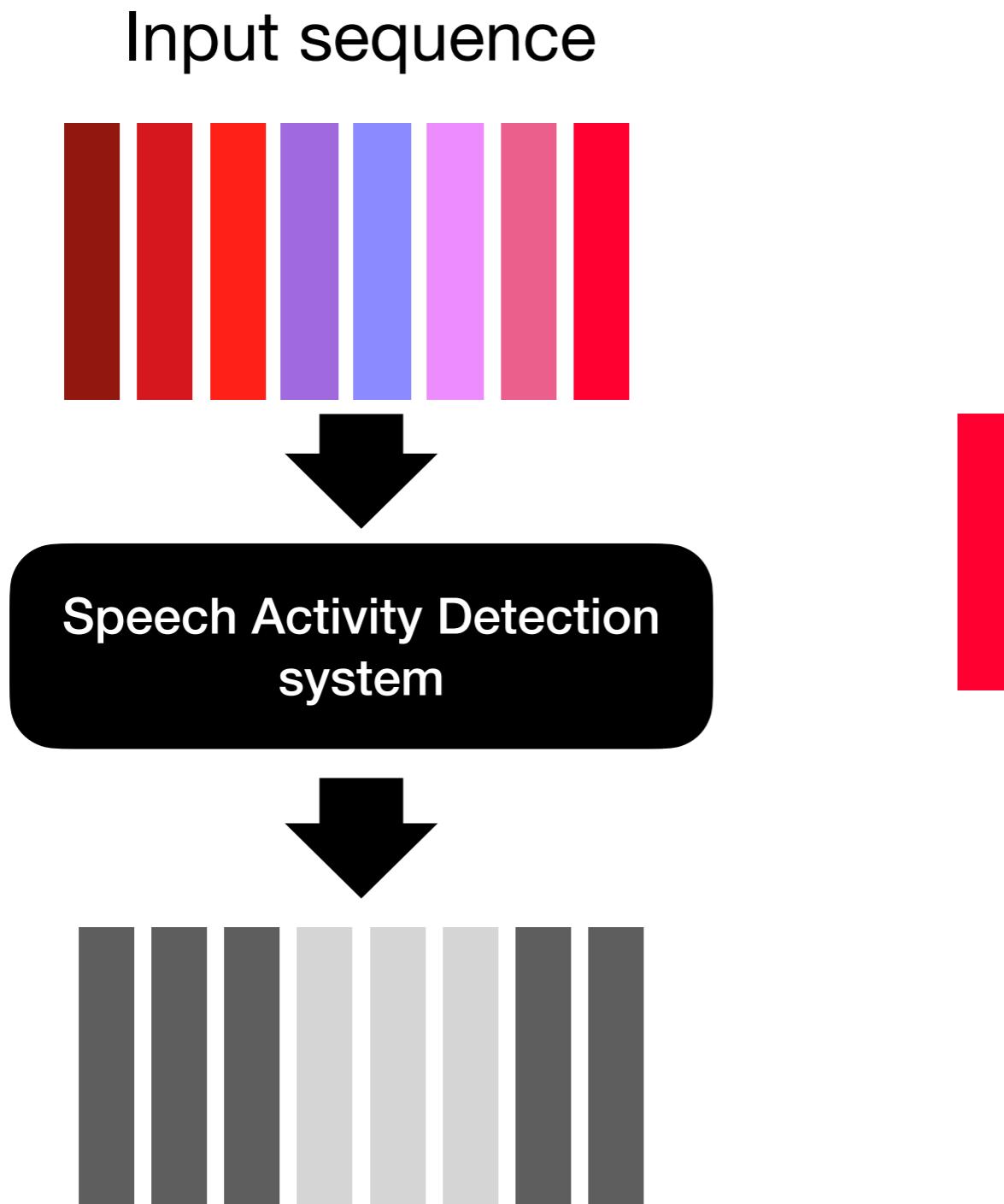
Noise-aware Training of Acoustic Models

Speech/silence classification



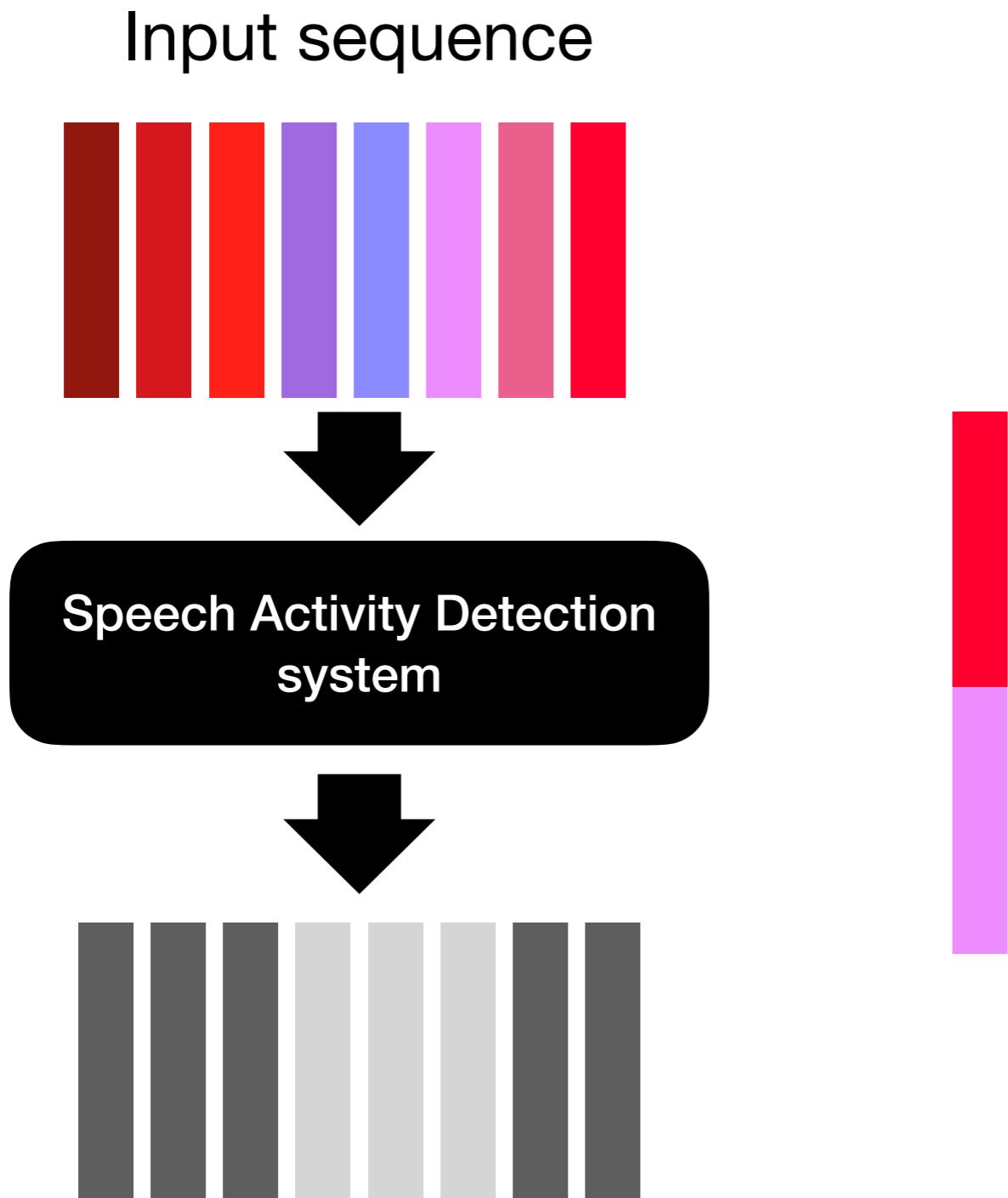
Noise-aware Training of Acoustic Models

Speech/silence classification



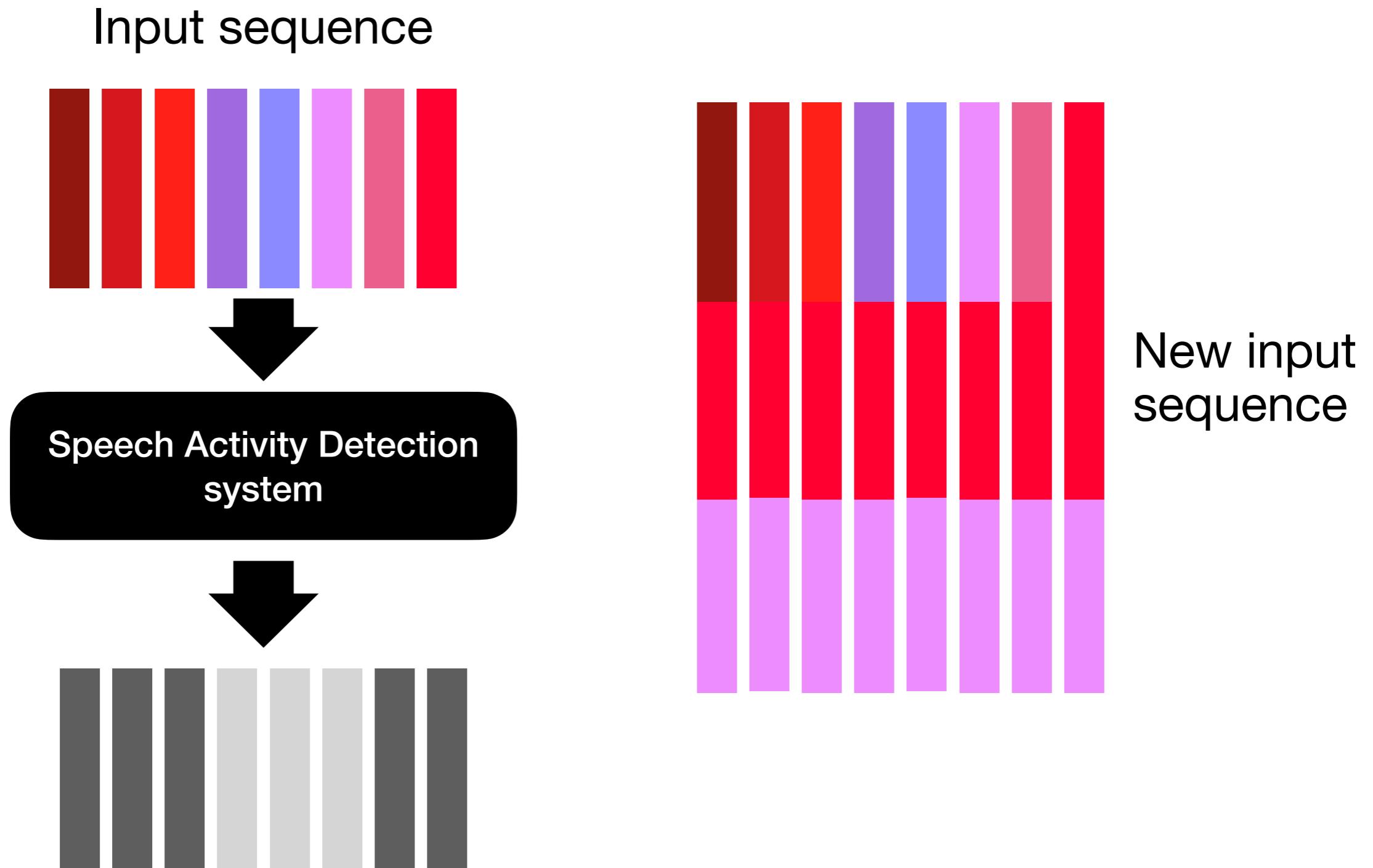
Noise-aware Training of Acoustic Models

Speech/silence classification



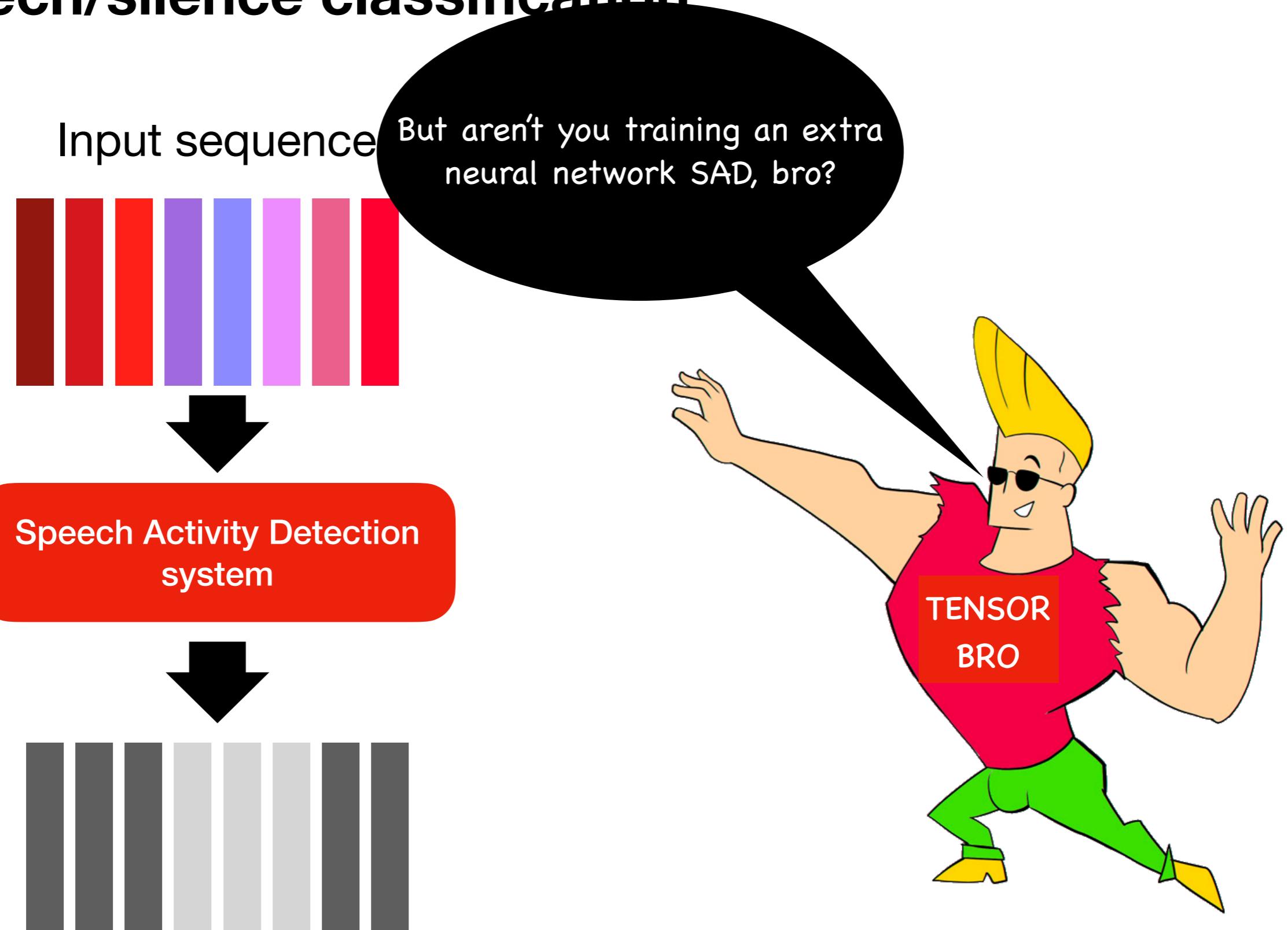
Noise-aware Training of Acoustic Models

Speech/silence classification



Noise-aware Training of Acoustic Models

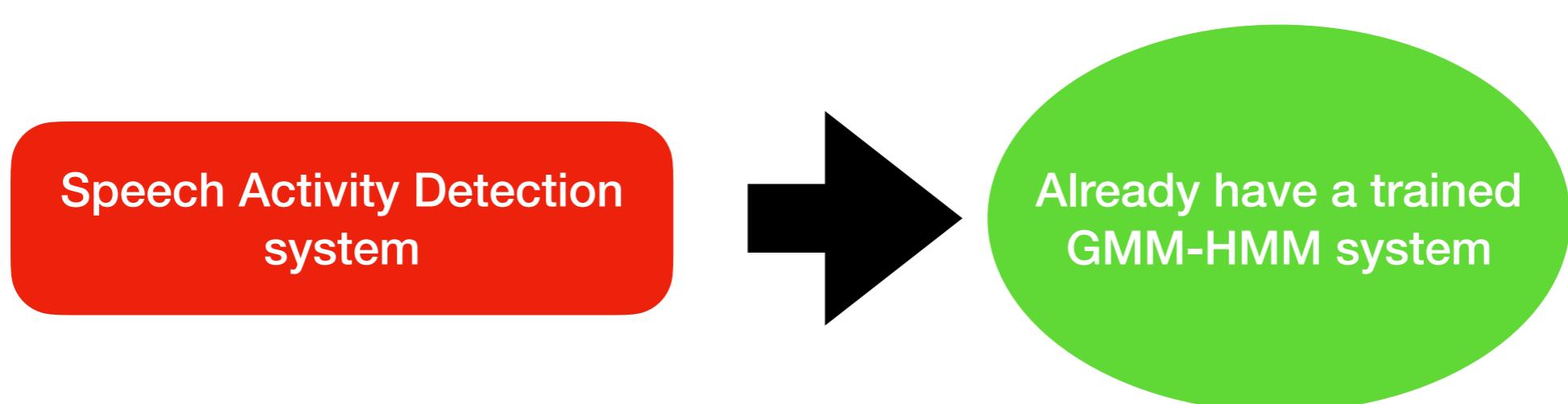
Speech/silence classification



Noise-aware Training of Acoustic Models

Speech/silence classification

No extra training needed!



Noise-aware Training of Acoustic Models

Extension to online ASR

How to estimate noise embedding in streaming ASR?

Noise-aware Training of Acoustic Models

Extension to online ASR

How to estimate noise embedding in streaming ASR?



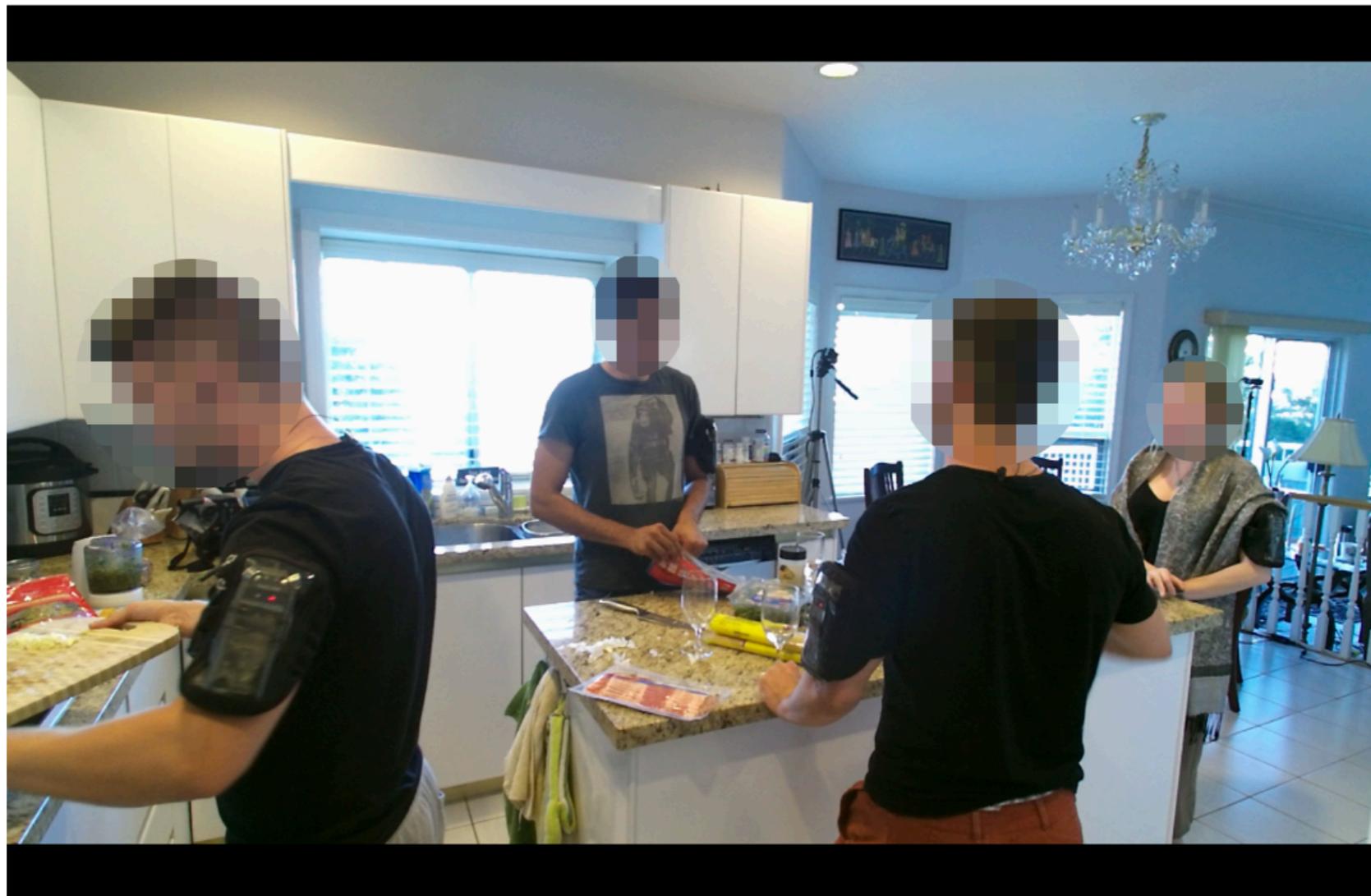
Bayesian Model

$$\mu_i = \begin{bmatrix} \mu_{s_i} \\ \mu_{n_i} \end{bmatrix}$$
$$\mathbf{K}_i = \begin{bmatrix} (1 + r_{s_i} N_{s_i}) \boldsymbol{\Lambda}_s & -\boldsymbol{\Lambda}_s \mathbf{B} \\ -\mathbf{B}^T \boldsymbol{\Lambda}_s & (1 + r_{n_i} N_{n_i}) \boldsymbol{\Lambda}_n + \mathbf{B}^T \boldsymbol{\Lambda}_s \mathbf{B} \end{bmatrix}$$
$$\mathbf{Q}_i = \begin{bmatrix} \boldsymbol{\Lambda}_s (\mathbf{a} + r_{s_i} \mathbf{F}_{s_i}) \\ \boldsymbol{\Lambda}_n (\mu_n + r_{n_i} \mathbf{F}_{n_i}) + \mathbf{B}^T \boldsymbol{\Lambda}_s \mathbf{a} \end{bmatrix}.$$

The CHiME-6 Challenge

The CHiME-6 Challenge

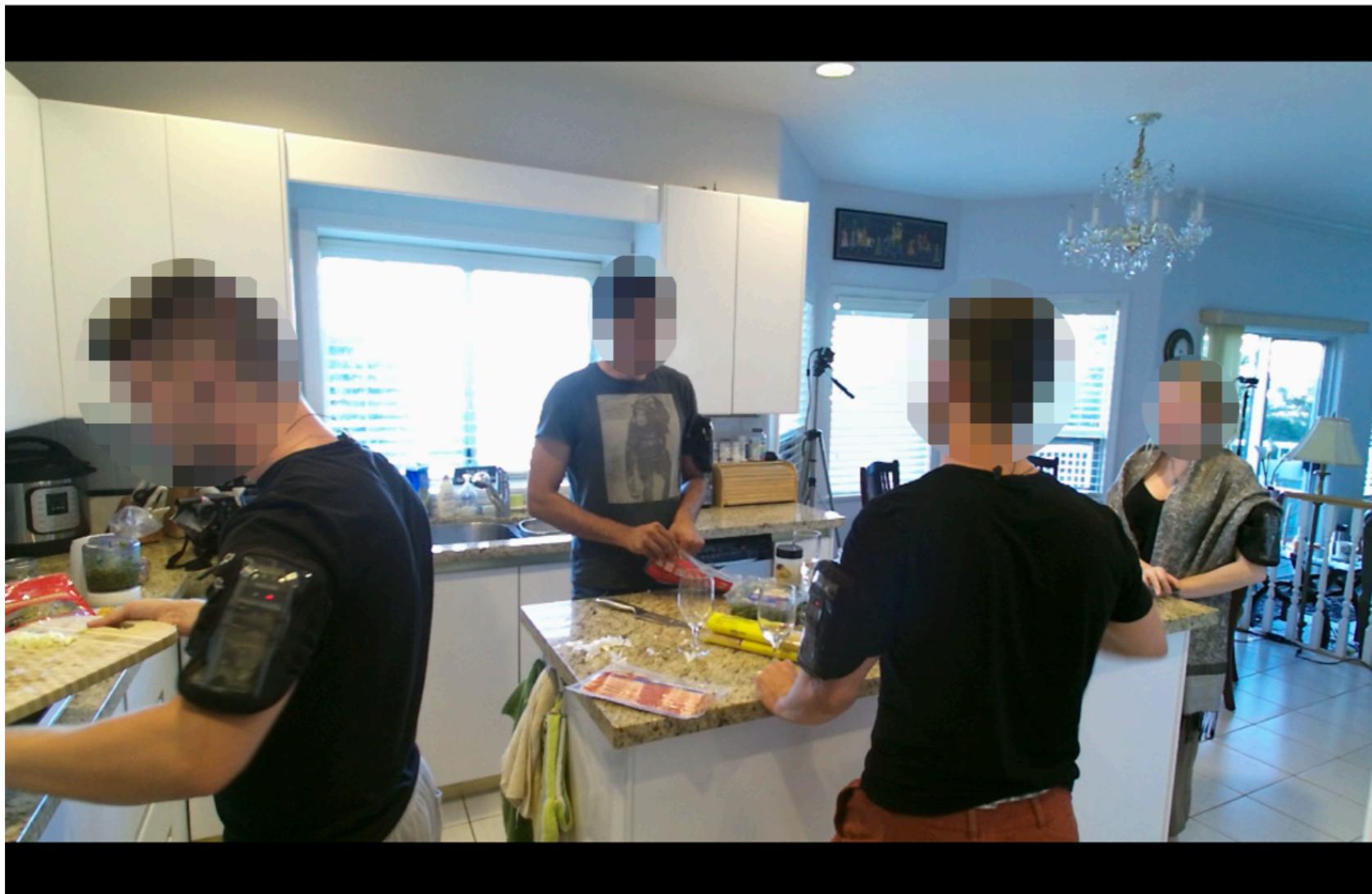
What is it?



<https://chimechallenge.github.io/chime6/overview.html>

The CHiME-6 Challenge

What is it?



Far-field

Noisy

Overlapping

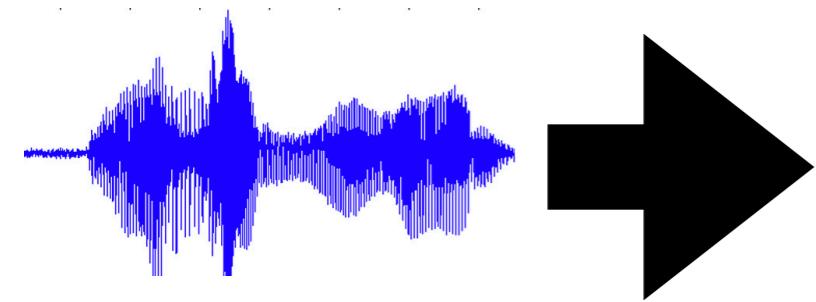
Real conversations

The CHiME-6 Challenge

How to solve it?

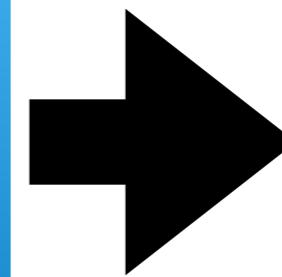
The CHiME-6 Challenge

How to solve it?



END-TO-END NEURAL
NETWORK

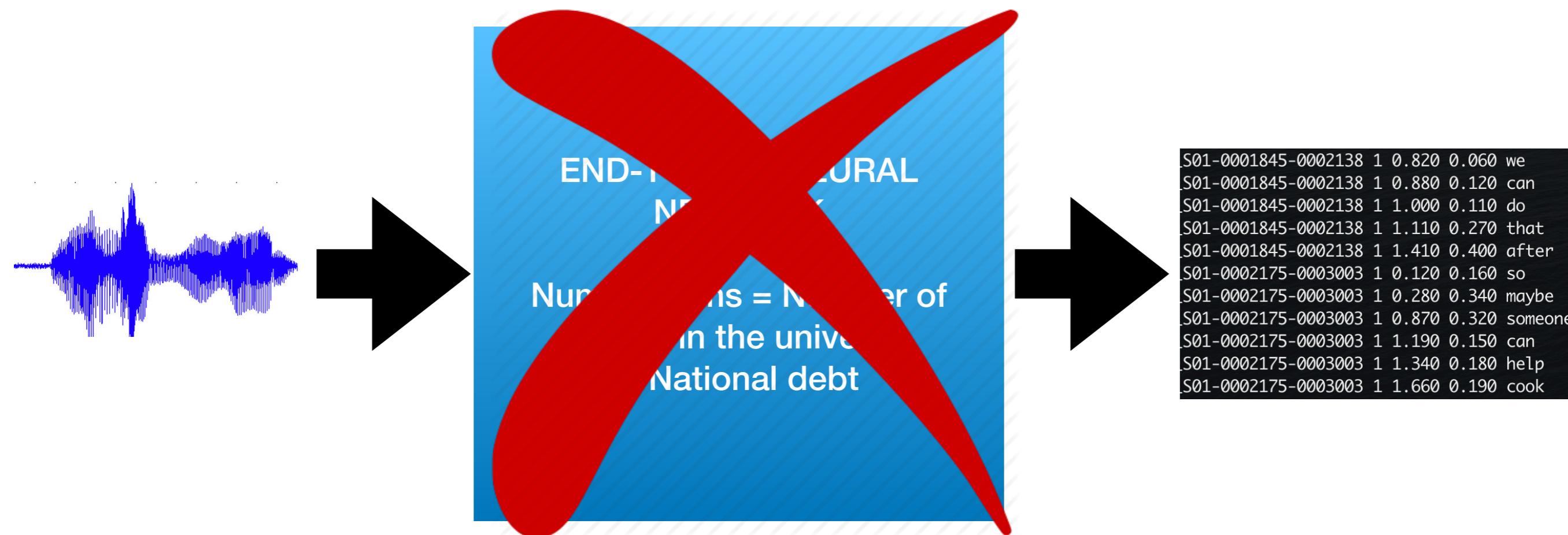
Num. Params = Number of
atoms in the universe
+ National debt

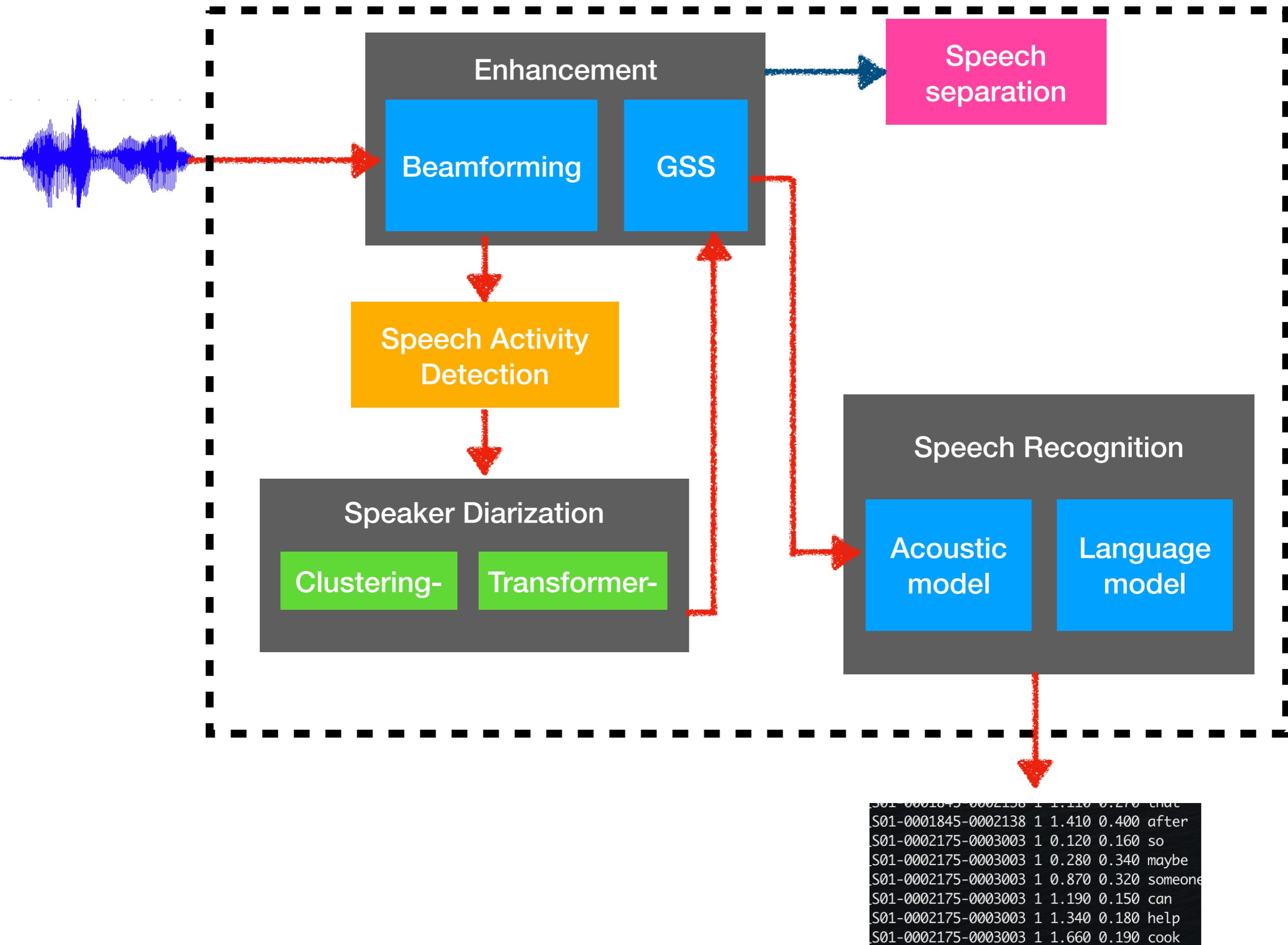


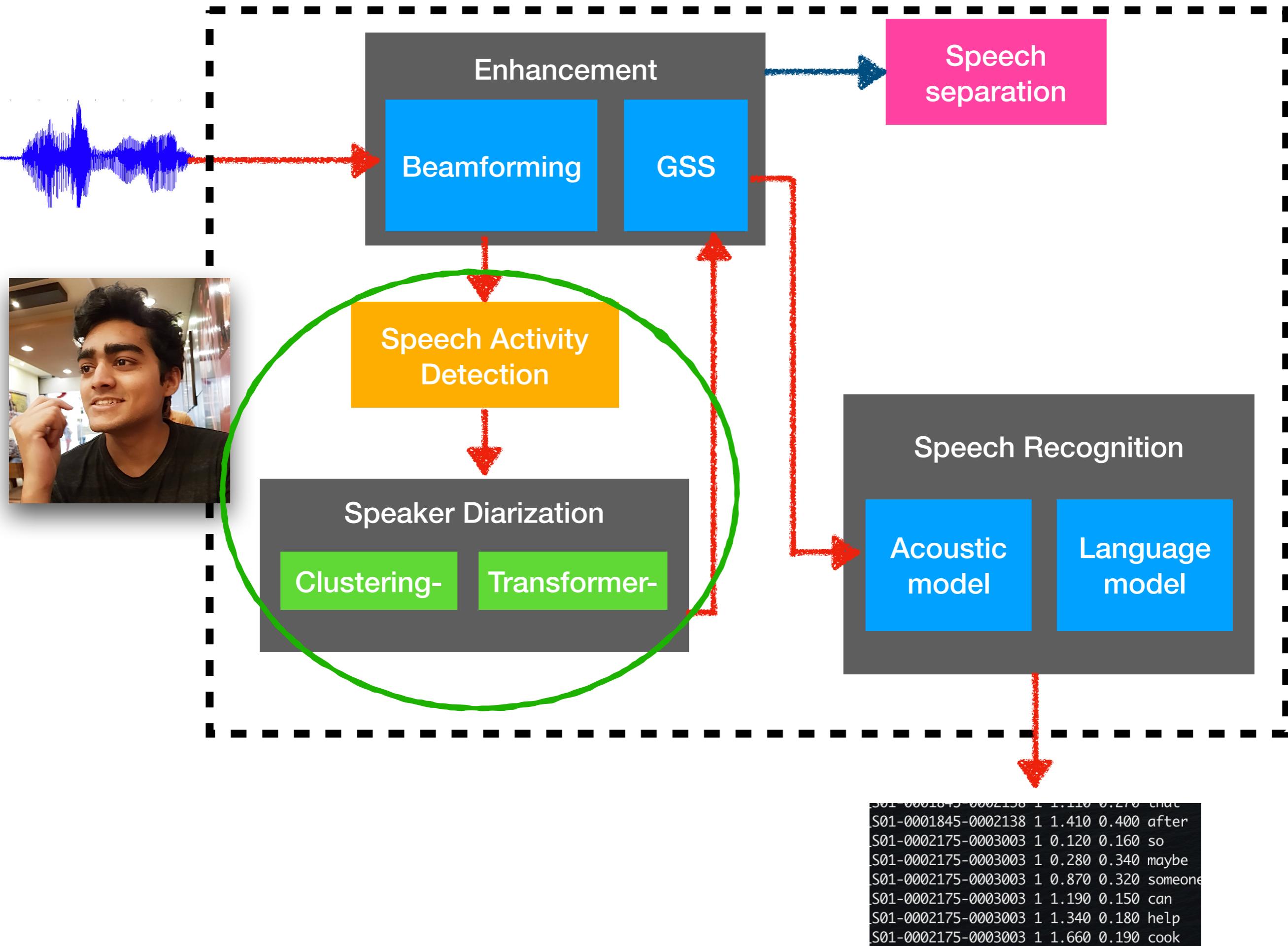
```
S01-0001845-0002138 1 0.820 0.060 we
S01-0001845-0002138 1 0.880 0.120 can
S01-0001845-0002138 1 1.000 0.110 do
S01-0001845-0002138 1 1.110 0.270 that
S01-0001845-0002138 1 1.410 0.400 after
S01-0002175-0003003 1 0.120 0.160 so
S01-0002175-0003003 1 0.280 0.340 maybe
S01-0002175-0003003 1 0.870 0.320 someone
S01-0002175-0003003 1 1.190 0.150 can
S01-0002175-0003003 1 1.340 0.180 help
S01-0002175-0003003 1 1.660 0.190 cook
```

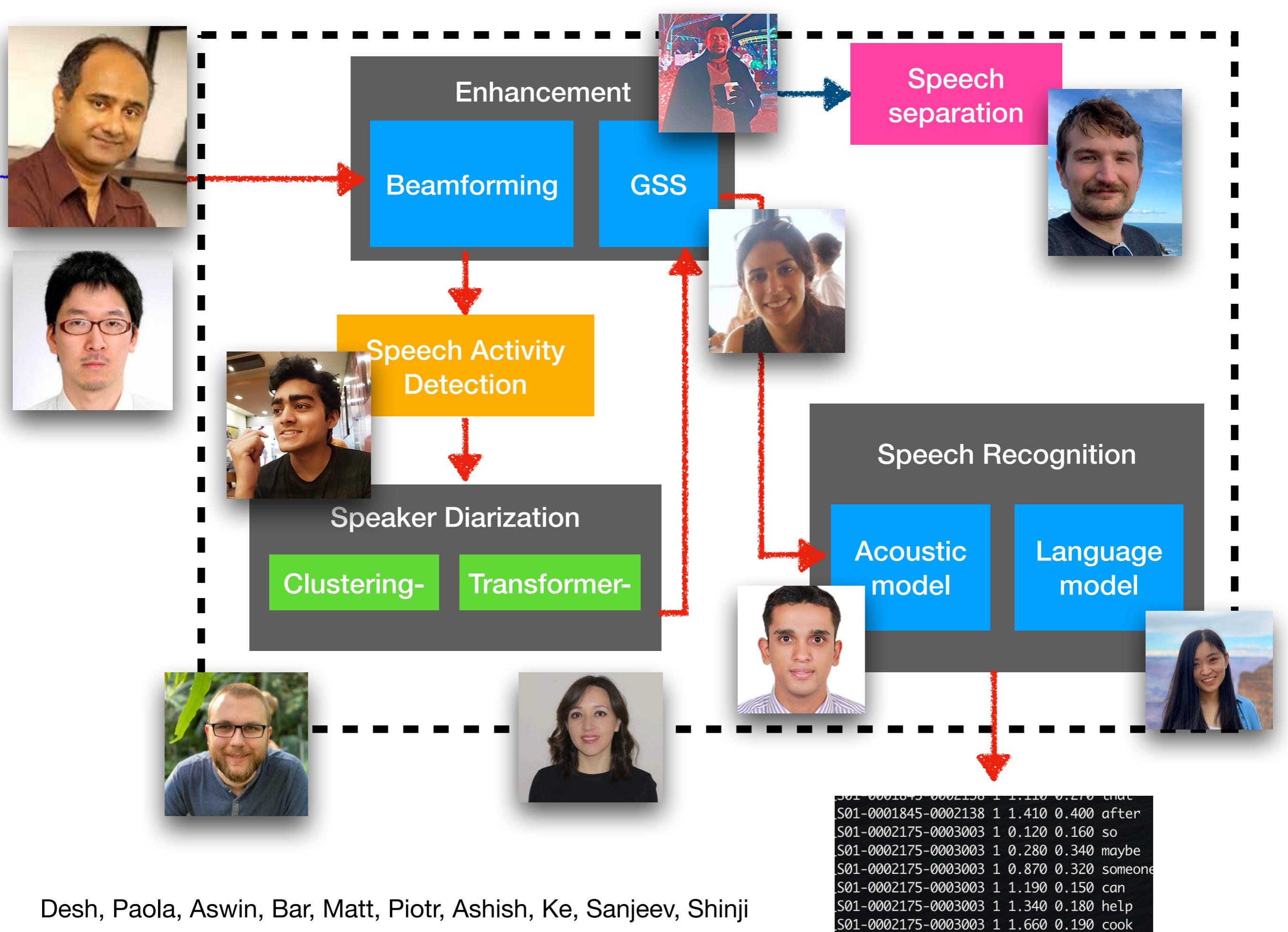
The CHiME-6 Challenge

How to solve it?









Desh, Paola, Aswin, Bar, Matt, Piotr, Ashish, Ke, Sanjeev, Shinji

Talk to me about...

(Inspired from Suzanna's talk)

Bouldering
Kaldi
Hybrid
Deep learning
End-to-end ASR^{but}
Transformers
CHiME-6
Diarization
Sequence modeling