

Speaker Verification over Cellular Networks

Ran Gazit

Persay

Standard Speech Corpora

- **TIMIT** clean
- **YOHO** clean
- **KING** clean
- **SIVA** PSTN
- **PolyVar** PSTN
- **POLYCOST** ISDN
- **Switchboard I-II** PSTN
- **OGI Speaker Recognition** PSTN

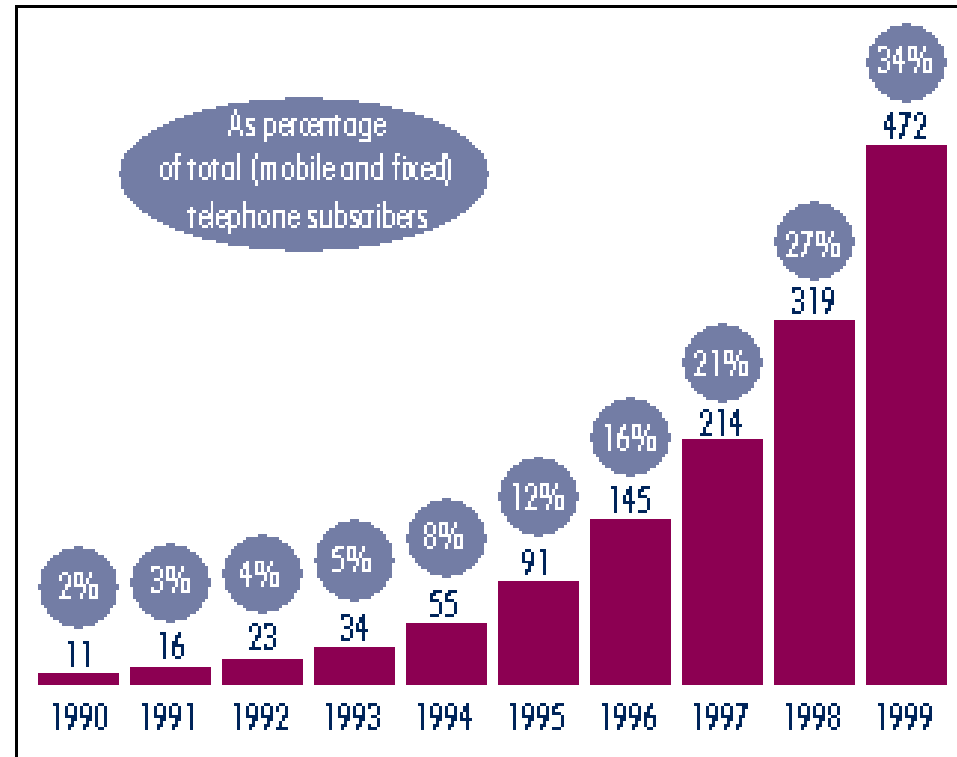
*Campbell and Reynolds, ICASSP 99
Melin, COST 250, 2000*

Cellular Penetration Rate

“It seems highly likely that the number of mobile cellular subscribers will surpass conventional fixed lines during the first decade...”

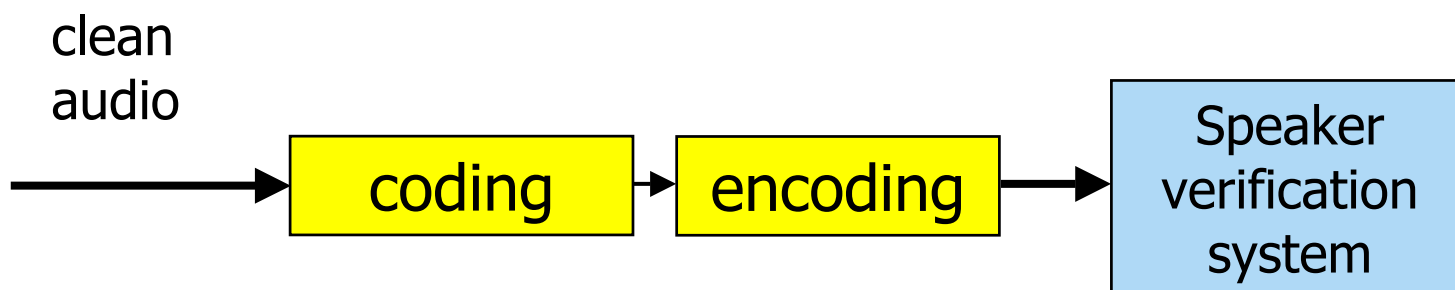
*ITU World
Telecommunication
Development Report,
1999*

**WORLDWIDE MOBILE CELLULAR SUBSCRIBERS
(MILLIONS)**



Source: ITU, 2000.

Previous Work

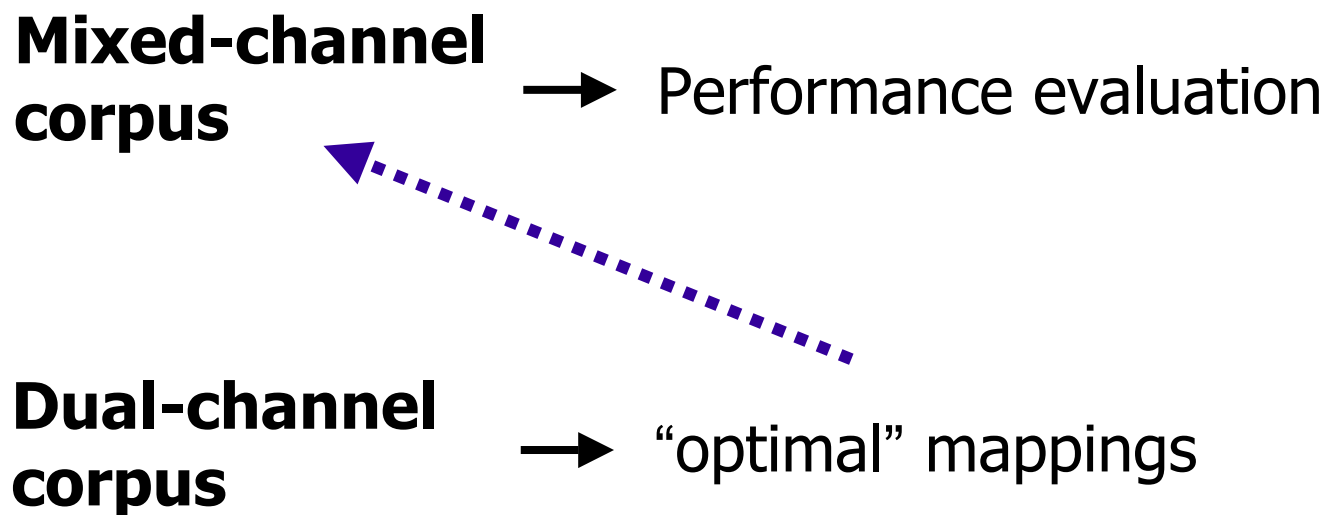


Kuitert and Boves, Eurospeech 97 (SESP - GSM)

Quatieri et. al., Eurospeech 99 (Swbrd - GSM, G.729, G.723.1)

Besacier et. al., ICASSP 2000 (TIMIT - GSM)

Method



Mixed-Channel Corpus

Target speakers

43F / 42M

10 1 min. Land-line

10 1 min. Cellular (various)

Background speakers

60F / 60M

1 min Land-line

1 min Cellular (various)

Experiment Setting

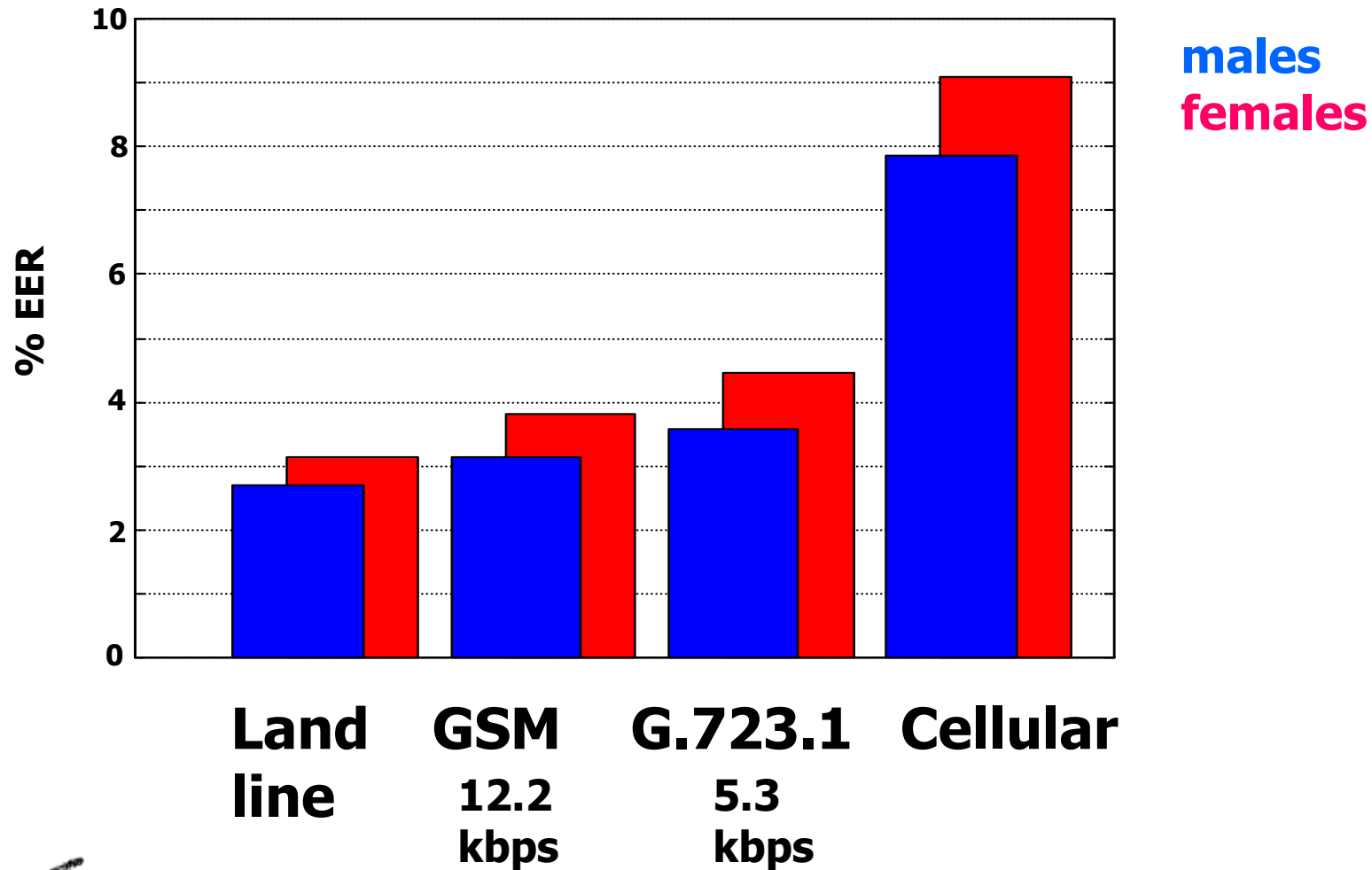
Working point

3 min. train
30 sec. test

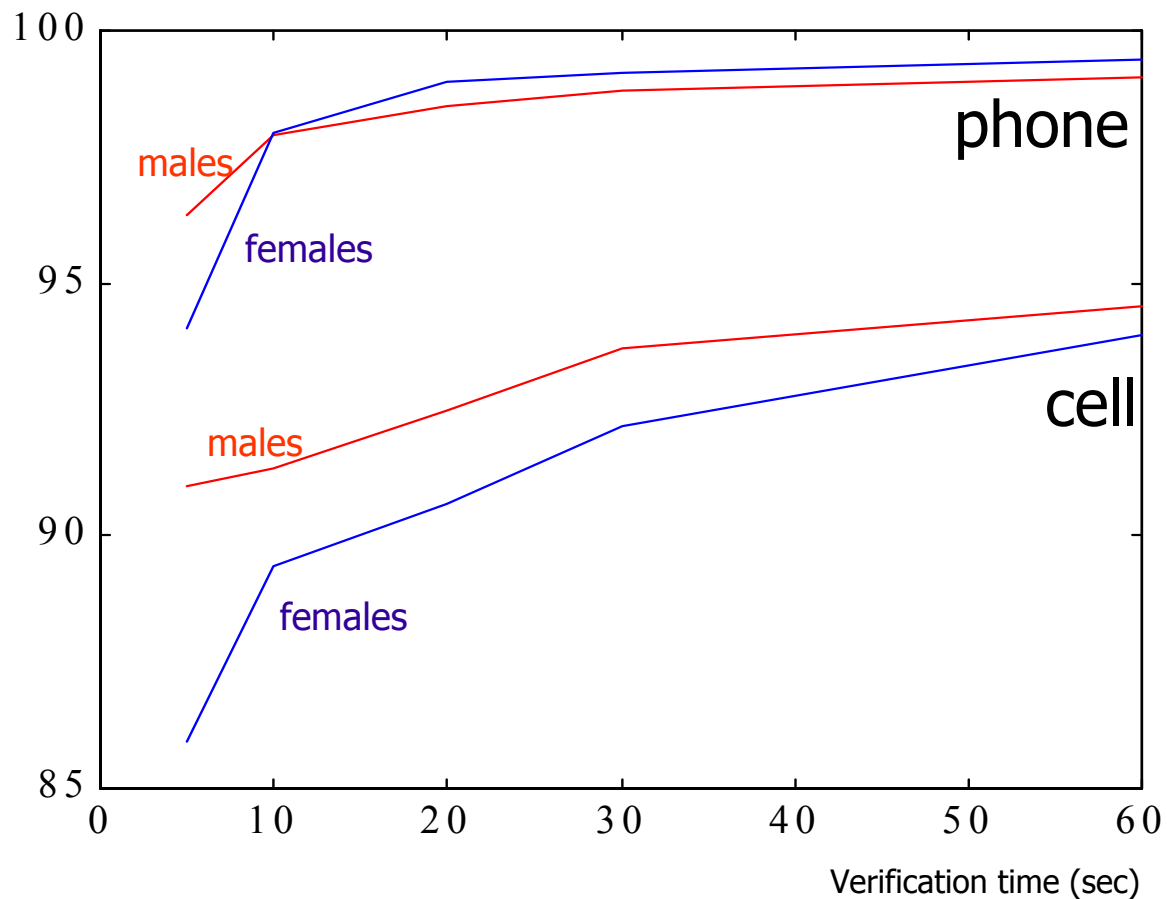
System

25 ms frames, 12.5 ms shift
Energy-based VOX
12 LPC-Ceps, CMS
12 delta-Ceps
30-mixture GMM
Gender, channel-dependent BM

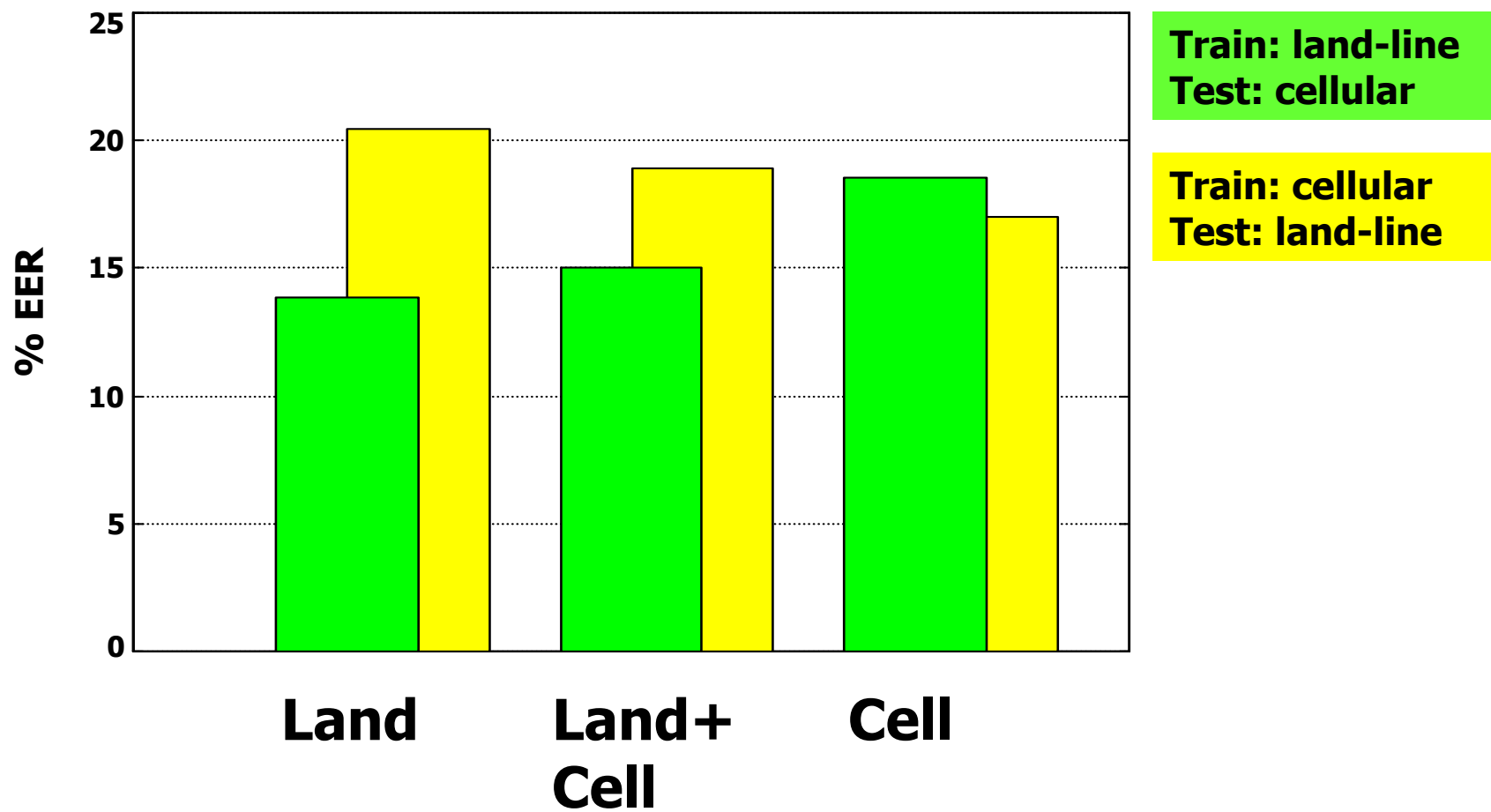
Matched-Conditions EER



Confidence Level



Mismatched-Conditions EER



Dual-Channel Corpus

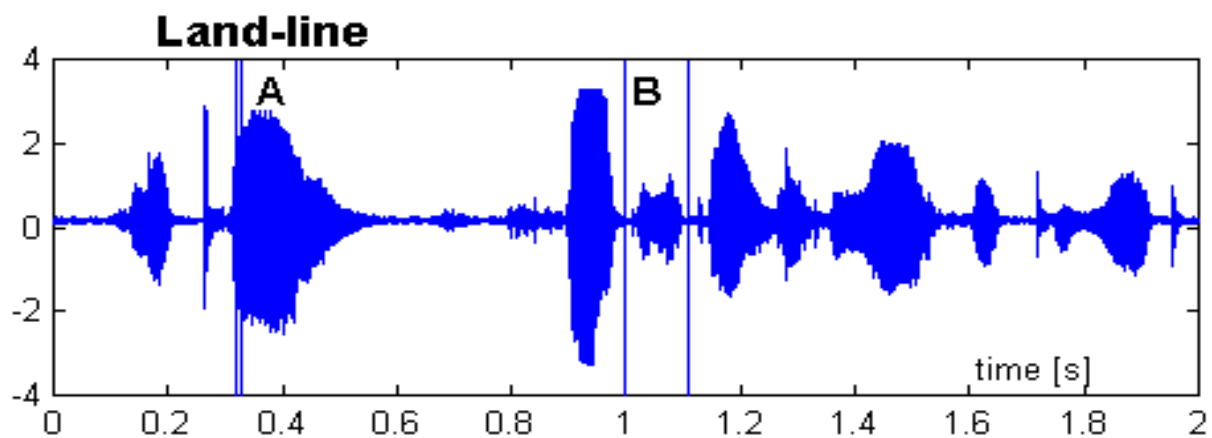
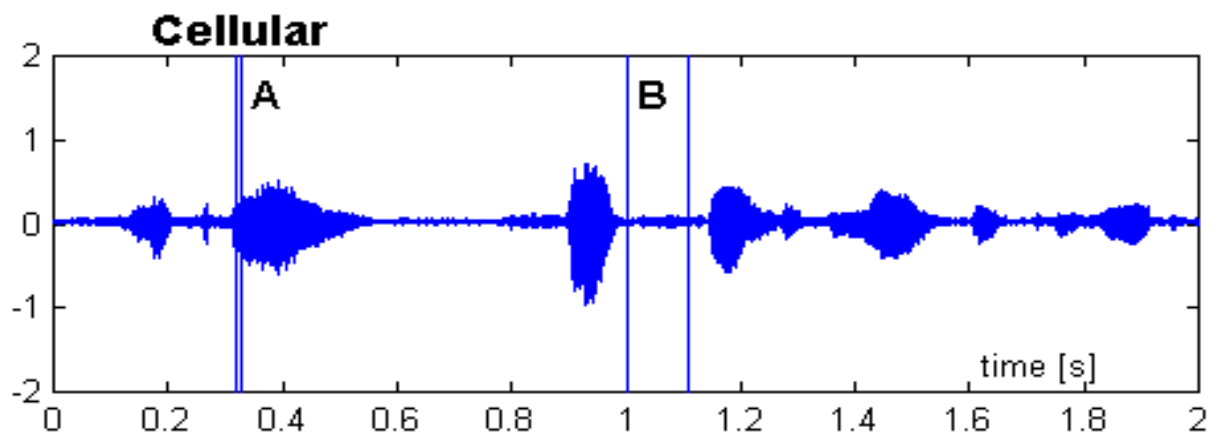
**Target
speakers**

10F / 10M

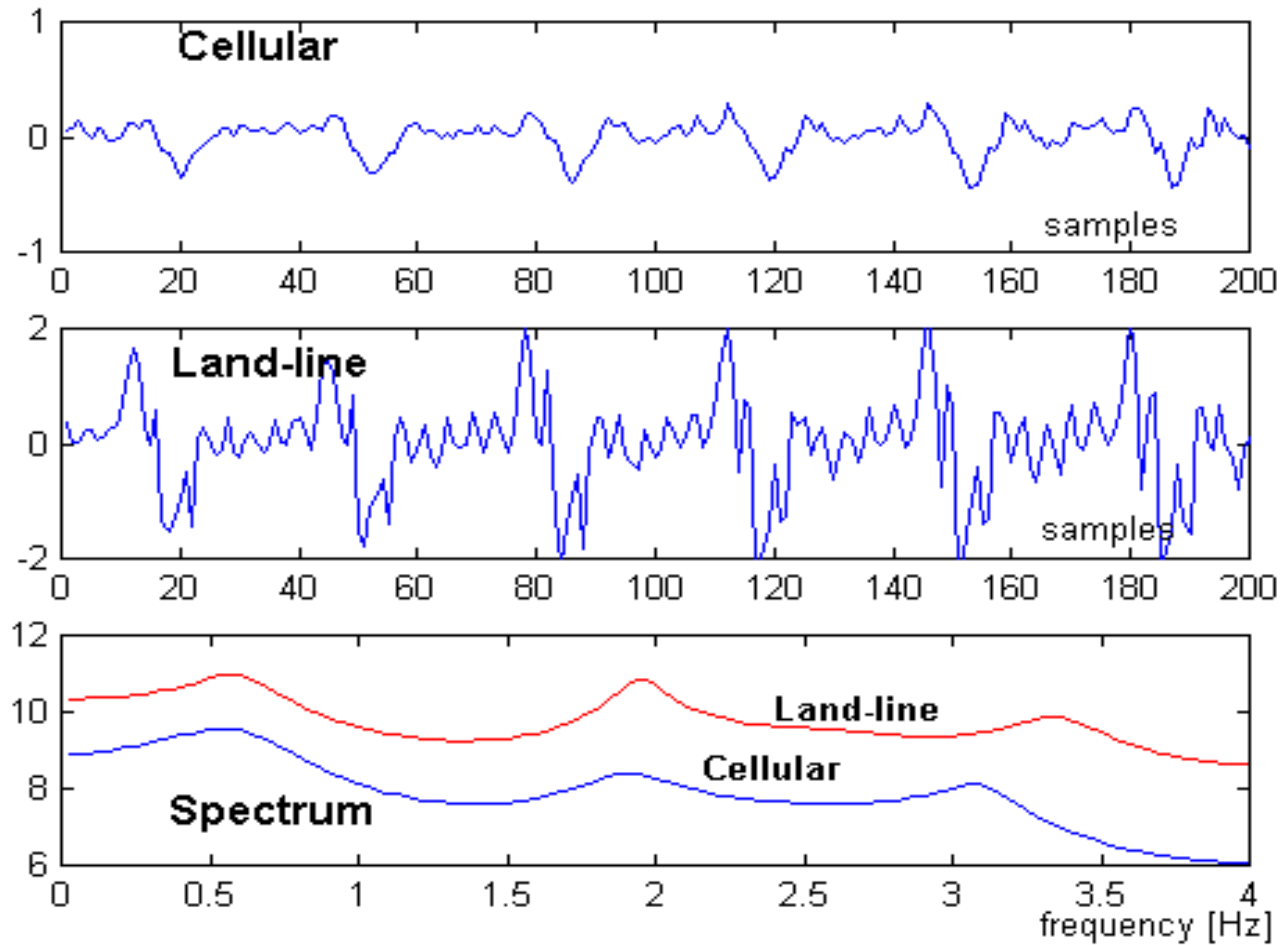
1 min. Land-line & Cellular
simultaneously



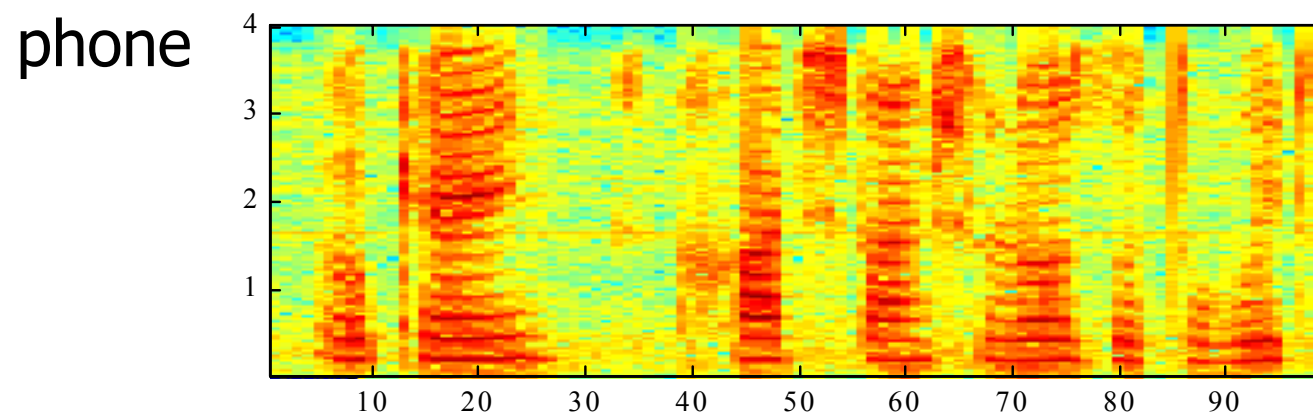
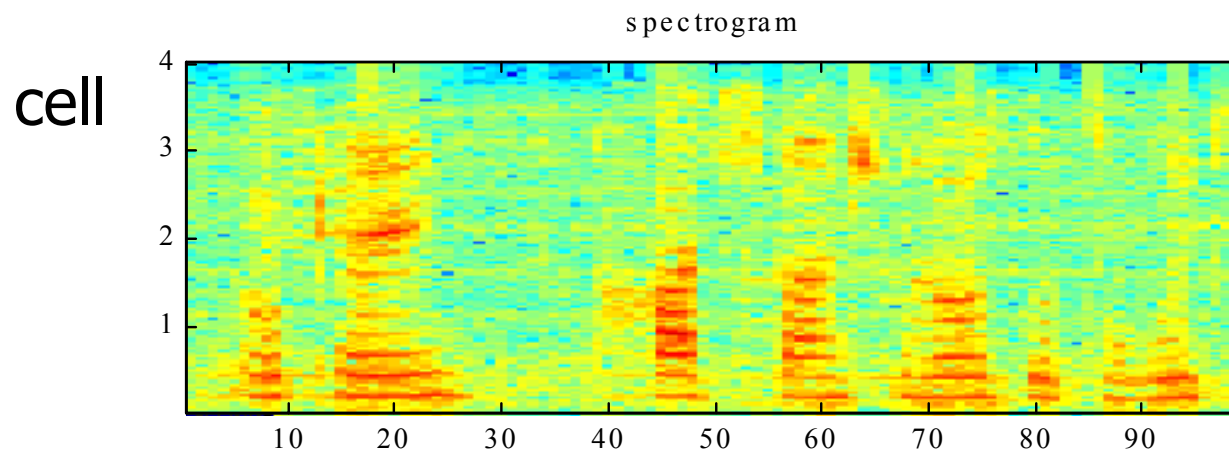
Dual-media Speech Signal



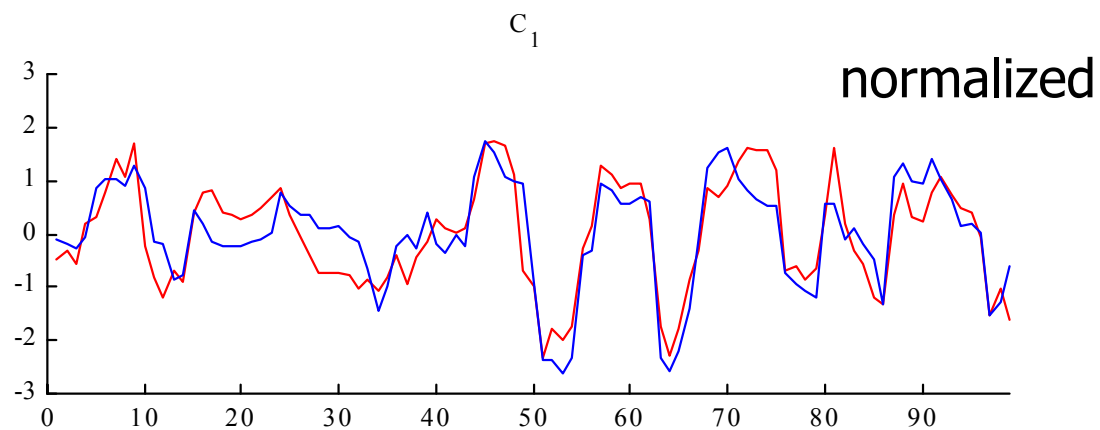
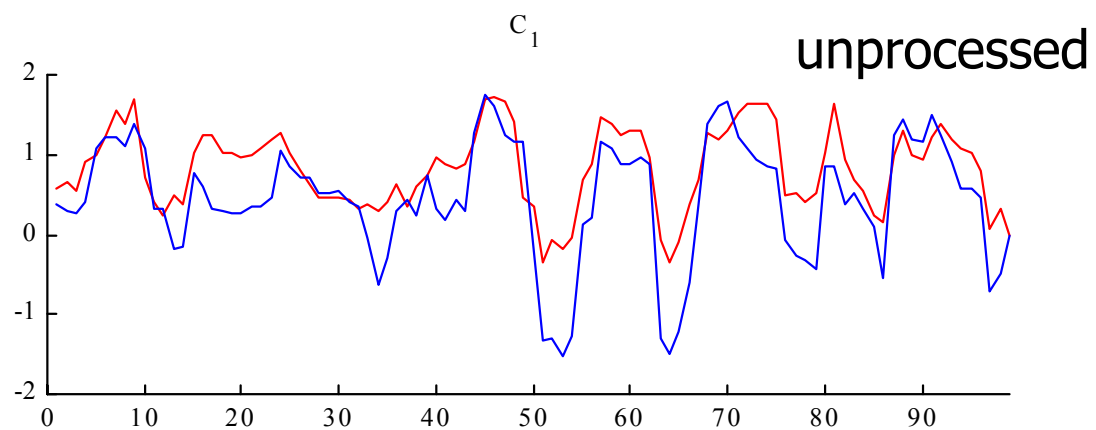
Dual-Media Speech Signal



Dual-media Spectrogram



Dual-media Cepstrum Track



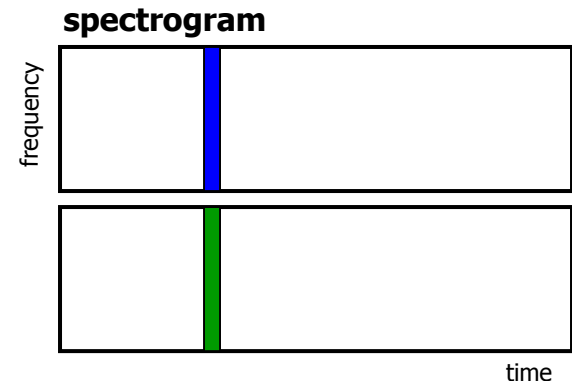
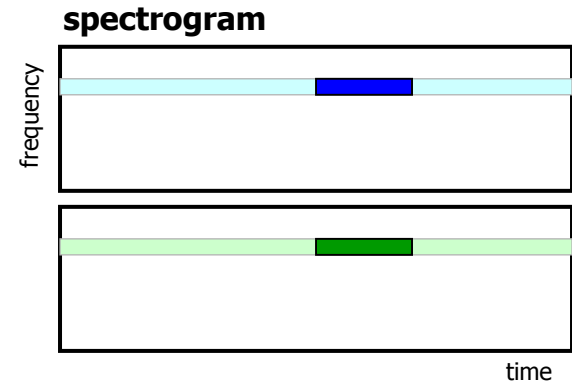
Dual-Channel Mappings

Goal:

- Suppress channel variability
- Preserve phonetic variability

Methods:

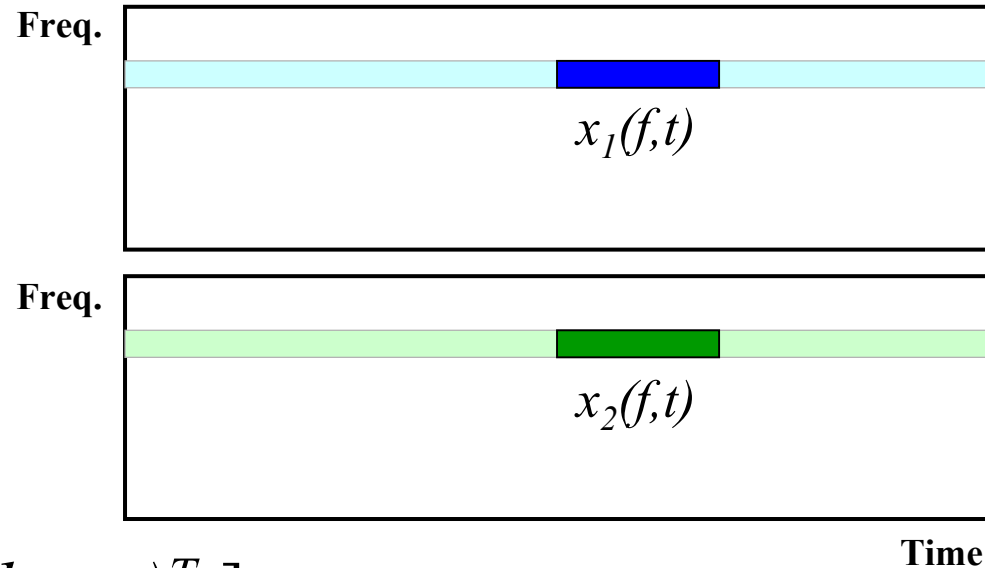
- Wiener-like FIR filters
- Spectral basis vectors



Hermansky et. al., ICASSP 95
Malayath et. al., DSP journal, 2000

Data-Driven Temporal Filters

$$Y = X^T h$$



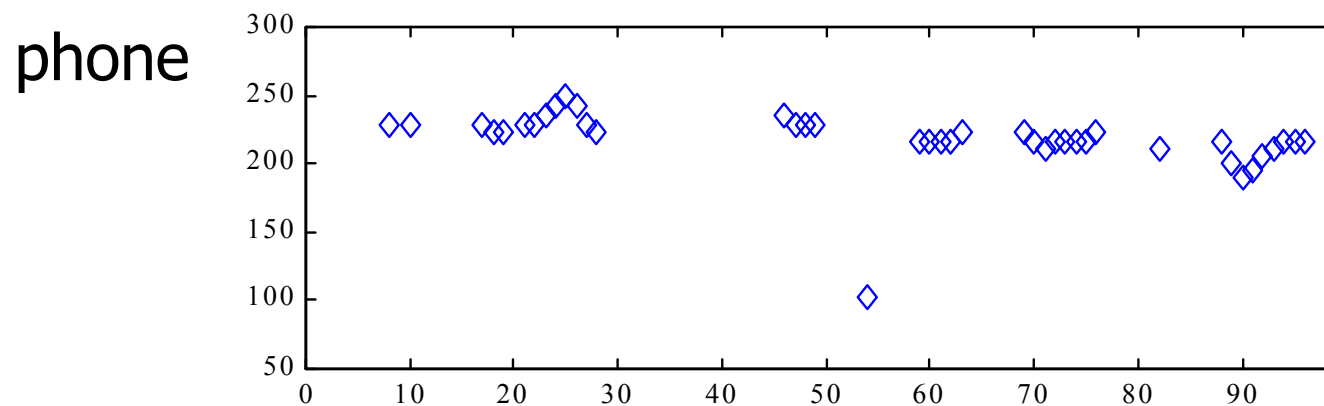
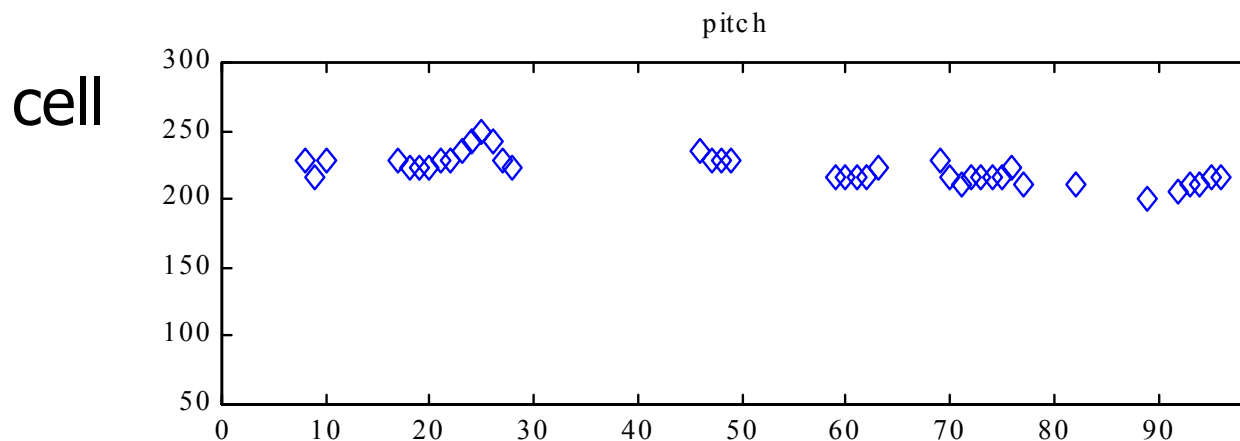
$$d = x_1 - x_2$$

$$\Sigma_{ch} = \text{E} [(d - \mu_d) (d - \mu_d)^T]$$

$$\Sigma_p = \text{E} [(\mu_p - \mu) (\mu_p - \mu)^T]$$

$$[\Sigma_{ch}^{-1} \Sigma_p] e = \lambda e$$

Dual-media Pitch Track



Cellular Effects

- Spectral and temporal distortions
- Noisy environment
- Speaking style
- Cutoffs
- Miniature microphone