# Using LSF Features for Speaker Verification in Noise

Pujita Raman and Dr. A. A. (Louis) Beex

DSP Research Laboratory

Wireless @ Virginia Tech
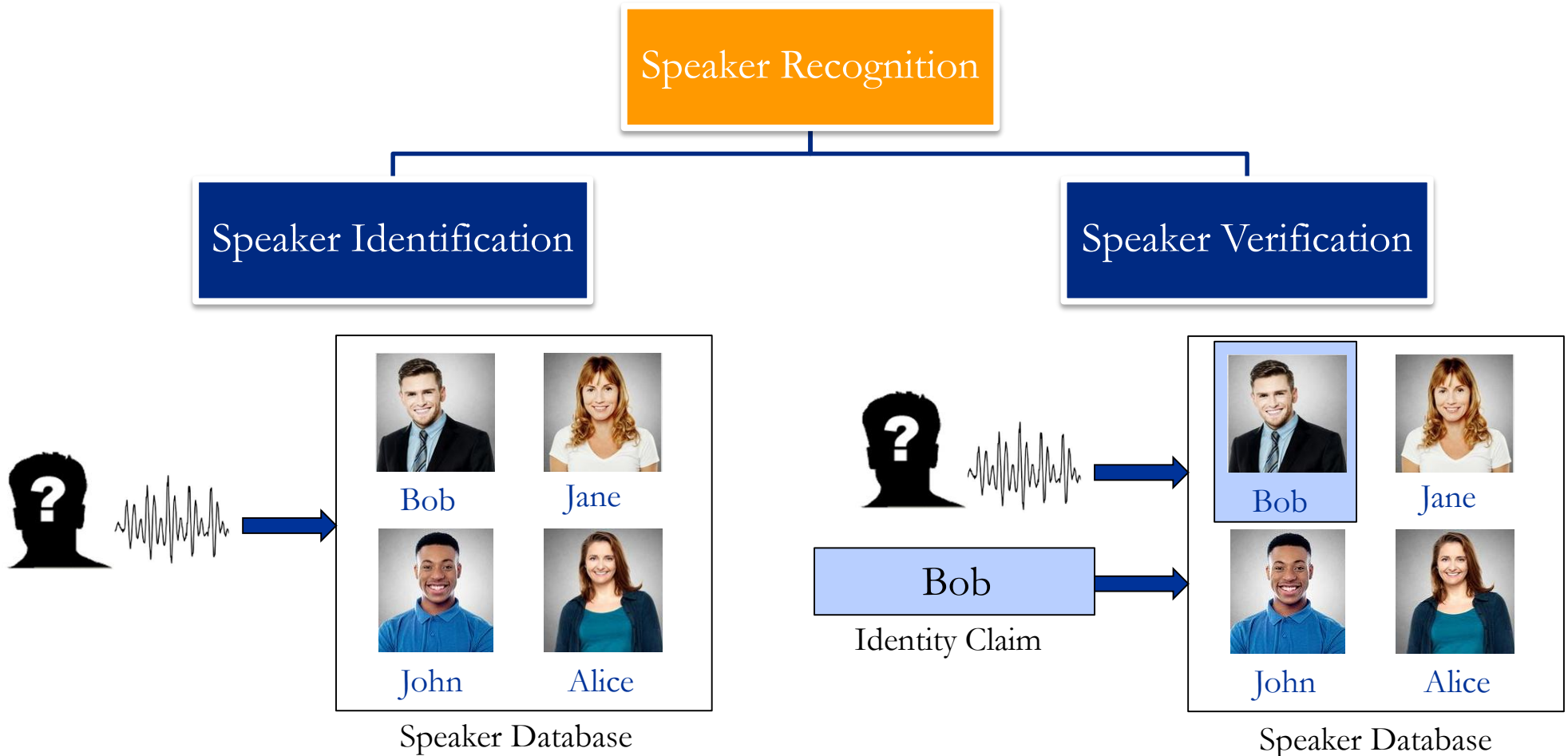
Wireless @ Virginia Tech

# Outline

- **Introduction**

- Motivation

- Feature Extraction

- Speaker Verification Framework

- Speaker Verification Results

- Conclusions

Wireless @ Virginia Tech

# Speaker Recognition

*The task of determining a speaker's identity using information extracted from his/her voice.*



**Speaker Recognition**
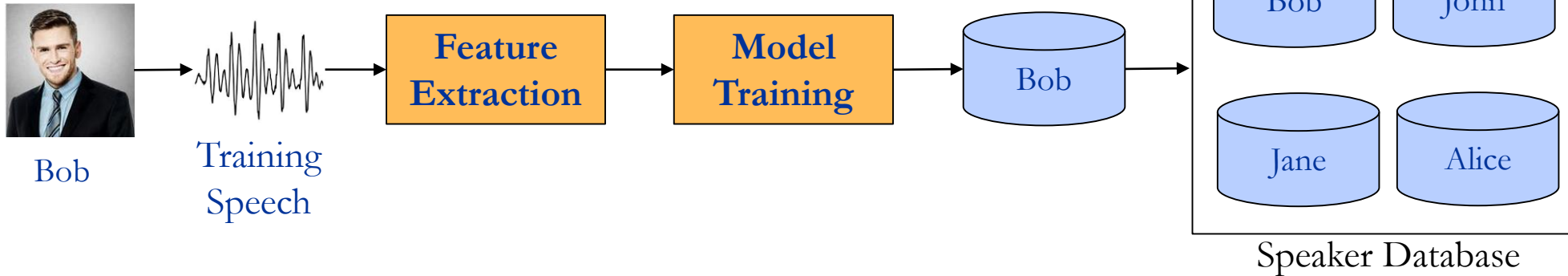
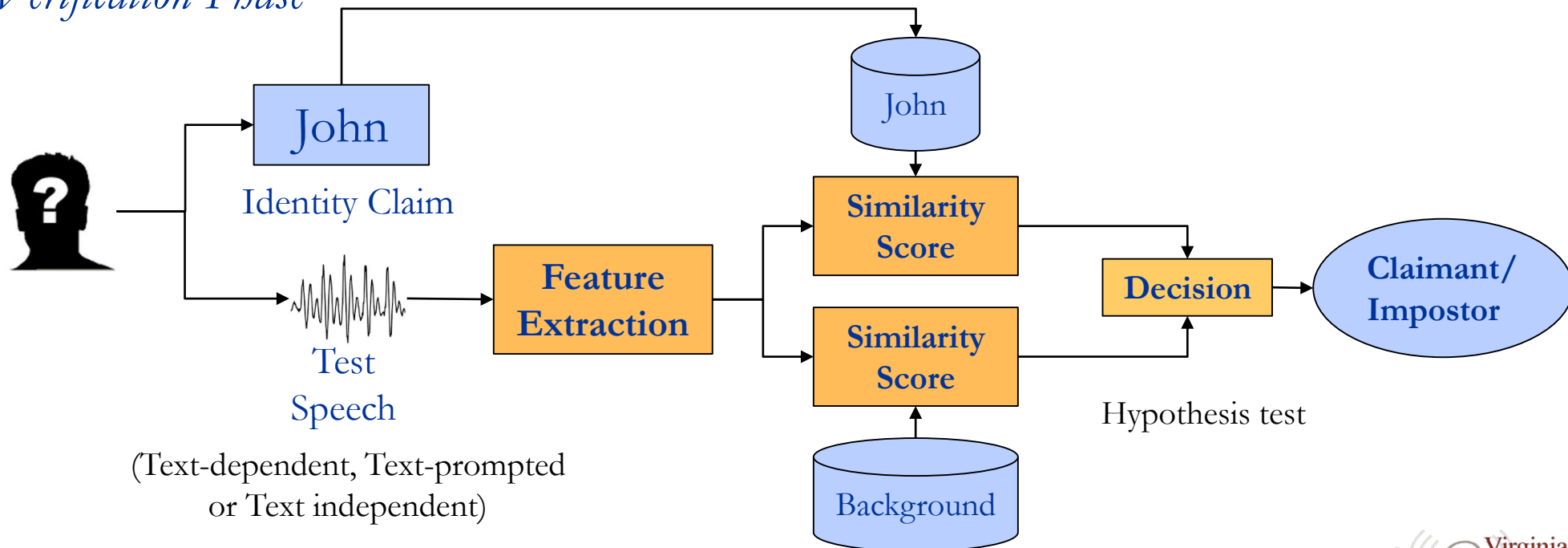**Speaker Identification**

**Speaker Verification**

Bob    Jane

John    Alice

Speaker Database

Bob

Identity Claim

Bob    Jane

John    Alice

Speaker Database

*Whose voice is this?*

*Is this Bob's voice?*

Wireless @ Virginia Tech

# Speaker Verification System

*Training/Enrollment Phase*



Bob → Training Speech → **Feature Extraction** → **Model Training** → Bob → Speaker Database (Bob, John, Jane, Alice)

Speaker Database

---

*Verification Phase*

John — Identity Claim

Test Speech
(Text-dependent, Text-prompted or Text independent)

→ **Feature Extraction** → **Similarity Score** / **Similarity Score** → **Decision** → **Claimant/ Impostor**

John

Background

Hypothesis test

Wireless @ Virginia Tech

# Outline

- Introduction
- **Motivation**
- Feature Extraction
- Speaker Verification Framework
- Speaker Verification Results
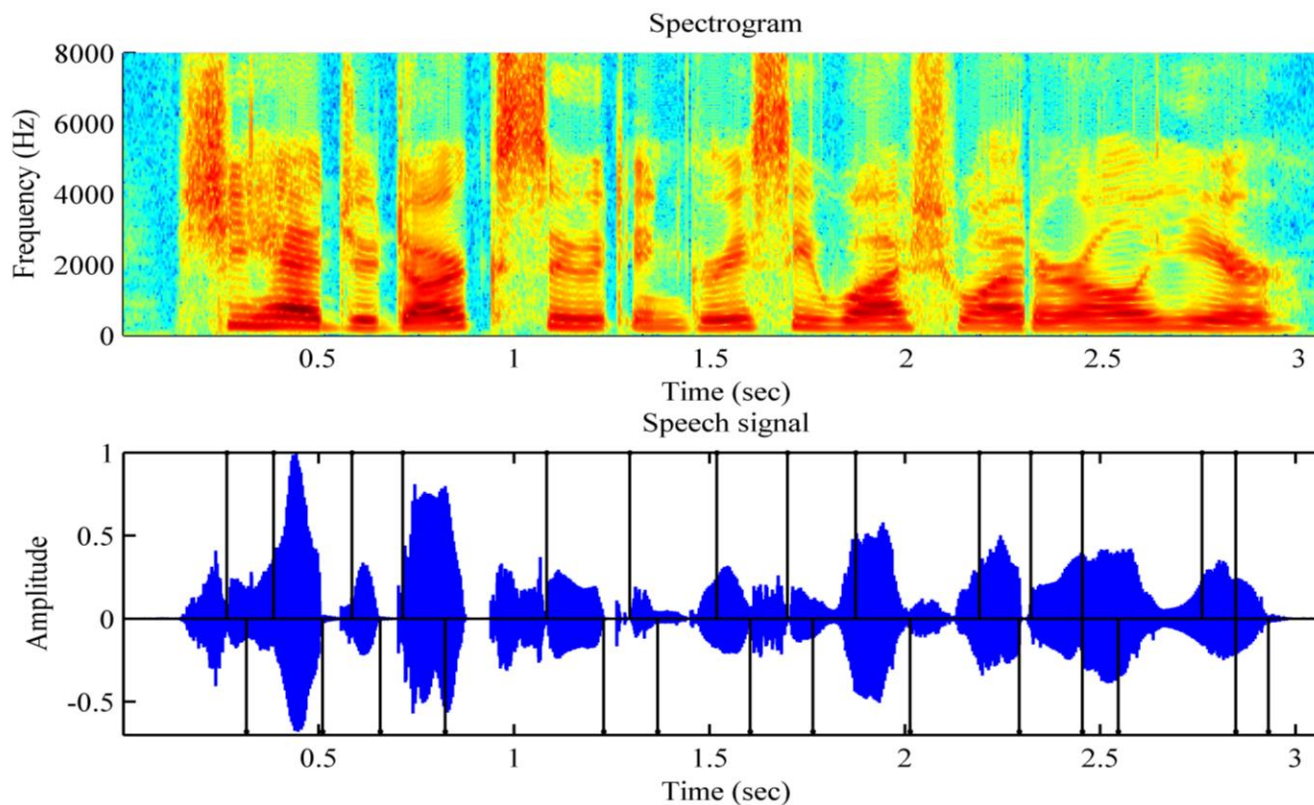- Conclusions

Wireless @ Virginia Tech

# Motivation

- State-of-the-art SV systems provide near perfect performance under clean conditions.

- Performance deteriorates in the presence of background noise.

- Noise-robustness improved by feature/model compensation and signal enhancement techniques.

- Drawbacks:
  - Require extensive training
  - Computationally expensive
  - Make assumptions about the noise characteristics.

*Can we improve performance by utilizing only <u>important zones</u> of speech, and discarding less important zones during verification?*

*Wireless @* Virginia Tech

# Relative Importance of Speech Zones
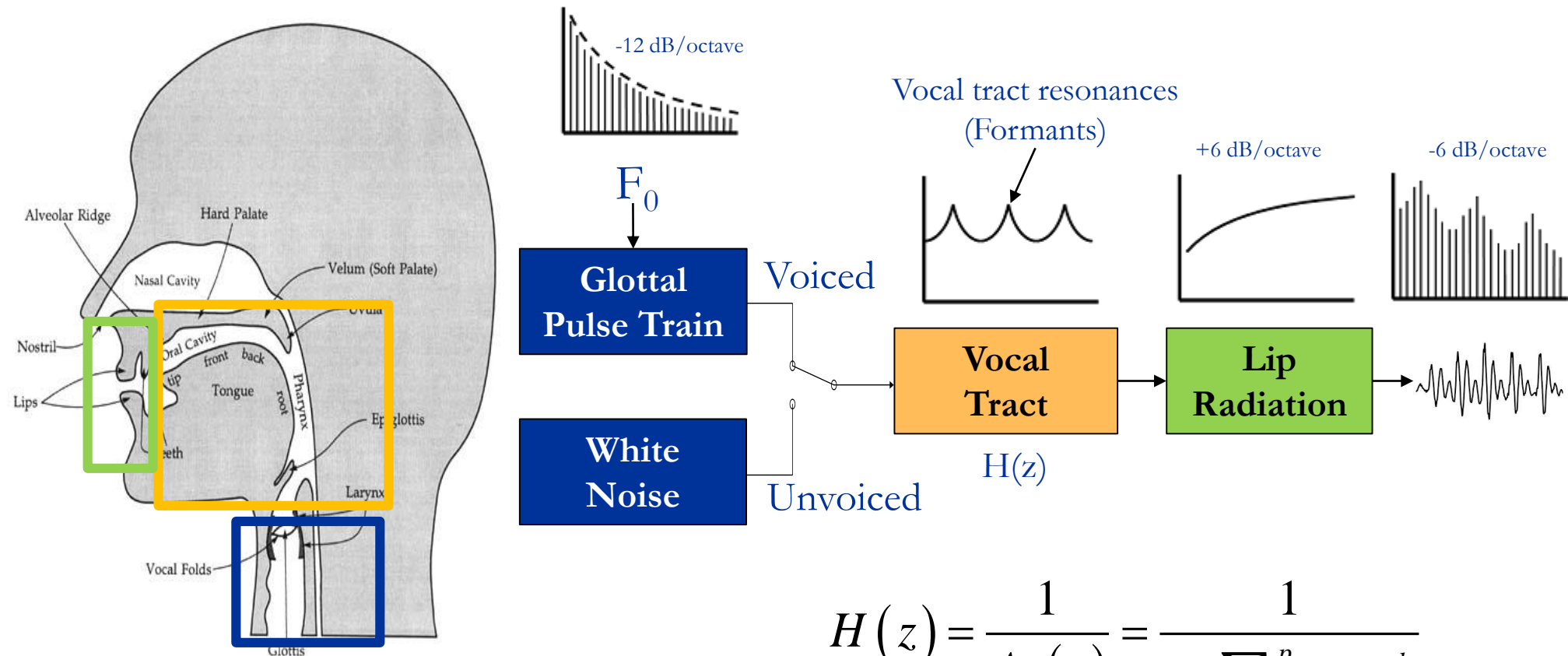
*Relative Importance = Amount of speaker-specific information*

- Co-articulation: the way a speaker moves from one sound to another is speaker specific.
- Dynamic transition regions are more speaker-specific than steady regions.
- We consider consonant-vowel (CV) and vowel-consonant (VC) transitions as vowels are easy to identify under noise.

# Outline

- Introduction
- Motivation
- **Feature Extraction**
- Speaker Verification Framework
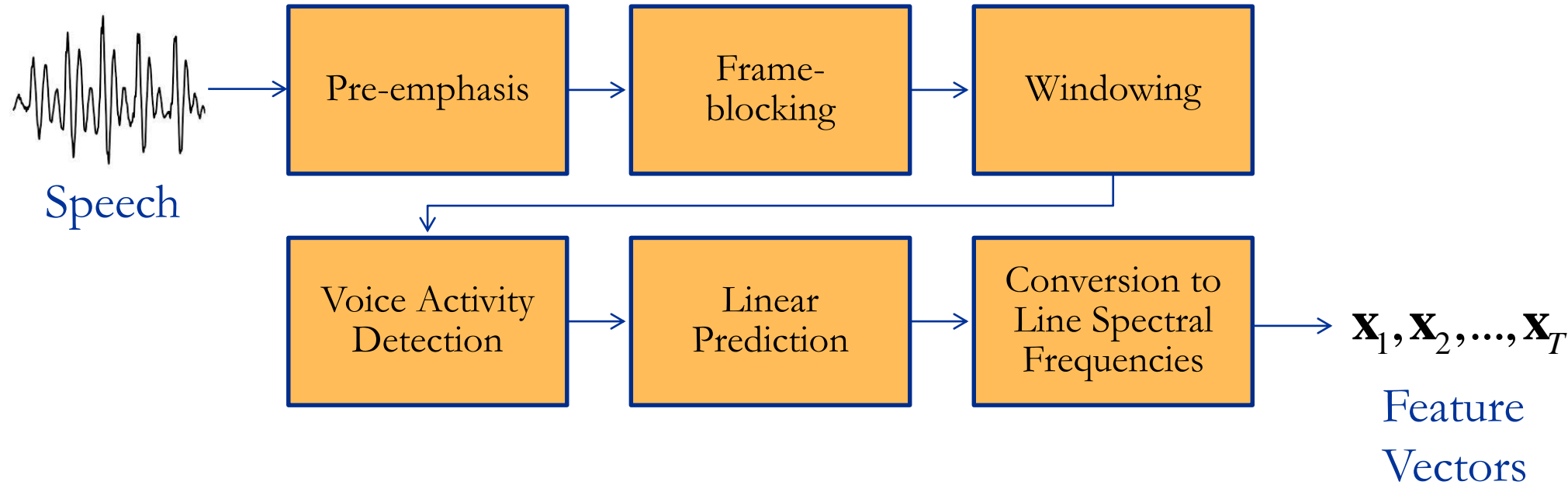- Speaker Verification Results
- Conclusions

Wireless @ Virginia Tech

# Source-Filter Model of Speech Production



$$H(z) = \frac{1}{A_p(z)} = \frac{1}{1 + \sum_{k=1}^{p} a_k z^{-k}}$$

Linear Prediction Coefficients

DSP Research Laboratory

Wireless @ Virginia Tech

# Feature Extraction



Speech

Pre-emphasis → Frame-blocking → Windowing → Voice Activity Detection → Linear Prediction → Conversion to Line Spectral Frequencies → $\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_T$

Feature Vectors

Wireless @ Virginia Tech

# Line Spectral Frequencies

Speech is a combination of two resonance conditions – vocal tract closed at the glottis and vocal tract open at the glottis.

$$A_p(z) = 1 + \sum_{k=1}^{p} a_k z^{-k} \quad \Longrightarrow \quad A_p(z) = \frac{P(z) + Q(z)}{2}$$

Closed glottis:

$$P(z) = A_p(z) + z^{-(p+1)} A_p(z^{-1})$$

(Symmetric)

$$\theta_{Pk} = e^{j\omega_{Pk}}, \quad 1 \le k \le p+1$$

Open glottis:

$$Q(z) = A_p(z) - z^{-(p+1)} A_p(z^{-1})$$

(Anti-symmetric)
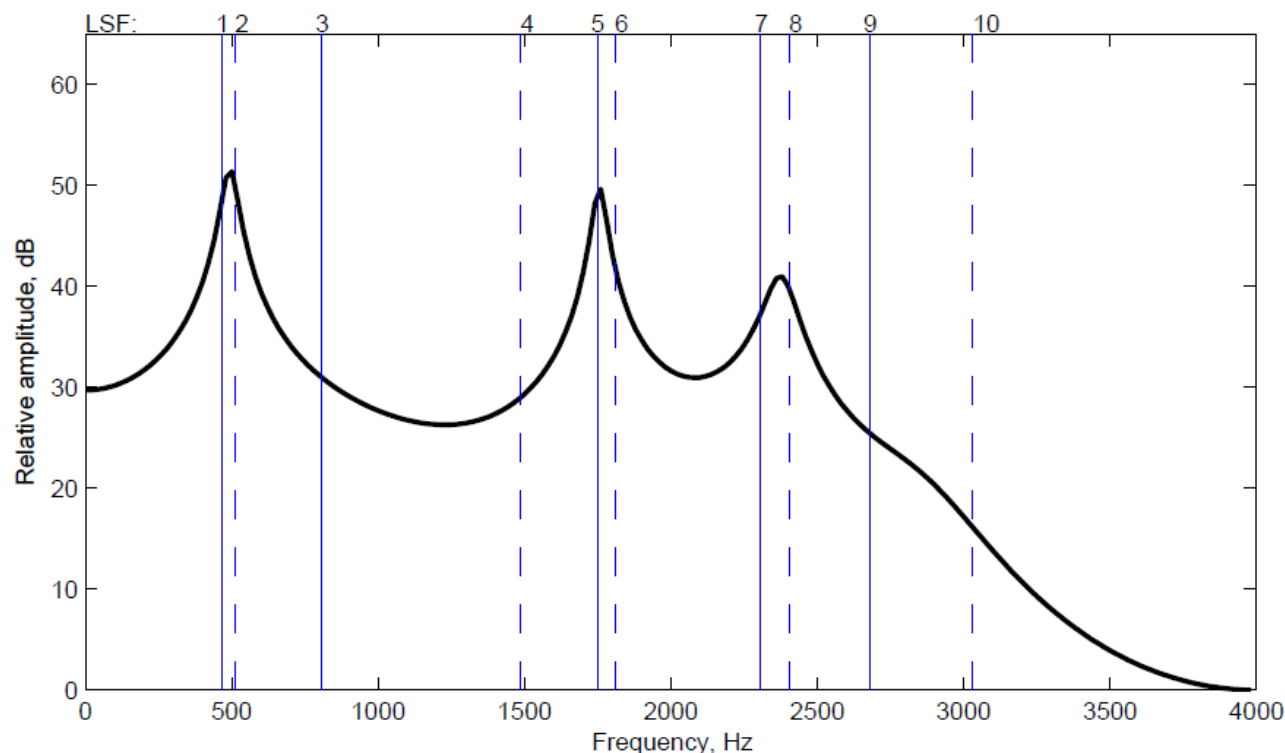
$$\theta_{Qk} = e^{j\omega_{Qk}}, \quad 1 \le k \le p+1$$

LSF Feature:
$$\mathbf{x} = \left[ \omega_{P1}\, \omega_{Q1}\, \omega_{P2}\, \omega_{Q2} \ldots \omega_{\frac{Pp}{2}}\, \omega_{\frac{Qp}{2}} \right]$$

- Efficient representation
- Good quantization properties
- Can be interpolated

Interlacing :
$$0 < \omega_{P1} < \omega_{Q1} < \omega_{P2} < \omega_{Q2} \ldots . < \omega_{\frac{Pp}{2}} < \omega_{\frac{Qp}{2}} < \pi$$

Wireless@ Virginia Tech

# Visualizing LSFs

- Line Spectral Frequencies (LSFs) are spectral features.
- Every formant is bracketed by an LSF pair
- If a pair of LSFs are far from each other, the magnitude response will be relatively flat around the two LSF.

DSP Research Laboratory

Wireless @ Virginia Tech

# Outline

- Introduction
- Motivation
- Feature Extraction
- **Speaker Verification Framework**
- Speaker Verification Results
- Conclusions

*Wireless* @ Virginia Tech

# Gaussian Mixture Model

- A Gaussian Mixture Model (GMM) $\lambda$ is a linear weighted sum of $M$ Gaussian components

$$p\left(\mathbf{x}|\lambda\right) = \sum_{m=1}^{M} p_m g_m\left(\mathbf{x}\mid\boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m\right) \qquad \sum_{m=1}^{M} p_m = 1 \qquad \begin{aligned}\boldsymbol{\Sigma}_m &\in \mathbb{R}^{D\times D}\\ \boldsymbol{\mu}_m &\in \mathbb{R}^{D}\end{aligned}$$

$$X = \left\{\mathbf{x}_t \in \mathbb{R}^D : 1 \le t \le T\right\} \implies \lambda = \left\{p_m, \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m\right\}, \quad m = 1, 2, \ldots, M$$

- Maximize GMM Likelihood: $p\left(X \mid \lambda\right) = \prod_{t=1}^{T} p\left(\mathbf{x}_t \mid \lambda\right)$

- **Expectation-Maximization Algorithm:** $\boxed{p\left(X \mid \lambda^{(i+1)}\right) \ge p\left(X \mid \lambda^{(i)}\right)}$

Wireless@ Virginia Tech

# Speaker Verification – Enrollment Phase

$H_0$ : Speech is from the hypothesized speaker – Speaker Model

$H_1$ : Speech is not from the hypothesized speaker – Background Model

- The Universal Background Model (UBM) is a 256 component GMM.
- Trained by pooling speech from 462 speakers in the TIMIT corpus.
- Speaker Models - GMMs obtained by Maximum a posteriori (MAP) adaptation of the UBM means
- Tighter coupling– better performance, faster scoring.

DSP Research Laboratory

Wireless @ Virginia Tech

# Speaker Verification – Testing Phase

*The speaker verification task is a simple <u>hypothesis test</u>*

Given a set of test features $X = \left\{ \mathbf{x}_t \in \mathbb{R}^D : 1 \le t \le T \right\}$

$H_0 : X$ is from speaker $S$ 　　　　　　　　　　　$H_1 : X$ is not from speaker $S$

$$\Phi_s = \sum_{t=1}^{T} \log p(\mathbf{x}_t \mid \lambda_s)$$ 　　　　$$\Phi_{ubm} = \sum_{t=1}^{T} \log p(\mathbf{x}_t \mid \lambda_{ubm})$$



Log-likelihood ratio:

$$\Lambda(X) = \Phi_s - \Phi_{ubm}$$

$$\Lambda(X) \begin{cases} \ge \theta & accept\ H_0 \\ < \theta & reject\ H_0 \end{cases}$$

Wireless @ Virginia Tech

# Speaker Verification – Performance Evaluation



Miss: Rejecting a target trial

$$E_{miss} = n_{miss} / n_t$$

False Alarm: Accepting an impostor trial

$$E_{fa} = n_{fa} / n_i$$

Equal Error Rate (%): Point at which probability of miss equals probability of false alarm.

DSP Research Laboratory

Wireless @ Virginia Tech

# Outline

- Introduction
- Motivation
- Feature Extraction
- Speaker Verification Framework
- **Speaker Verification Results**
- Conclusions

*Wireless* @ Virginia Tech

# Experimental Setup

| Parameter | Description |
|---|---|
| **Number of speakers (S)** | 168 |
| **Training set of each speaker** | All SX, SI sentences from TIMIT corpus (~3 seconds x 8) |
| **Test set of each speaker** | SA sentences from TIMIT (~3 seconds x 2) |
| **Feature Type** | LSF |
| **Feature Dimension/Order (p)** | 20 |
| **Frame Length (L)** | 20 msec |
| **Frame Shift (δ)** | 10 msec |
| **Number of GMM Components (M)** | 256 (UBM adapted GMM) |
| **GMM Covariance Type** | Nodal and Diagonal |
| **Noise Corpus** | SPIB Noise Dataset |

| # Speakers | # Target Trials | # Impostor Trials | # Total Trials |
|---|---|---|---|
| 168 | 168 x 2 = 336 | 168 x 167 x 2 = 56112 | 56112 + 336 = 56448 |

Wireless @ Virginia Tech

# SV System Performance

1) Static features - LSF
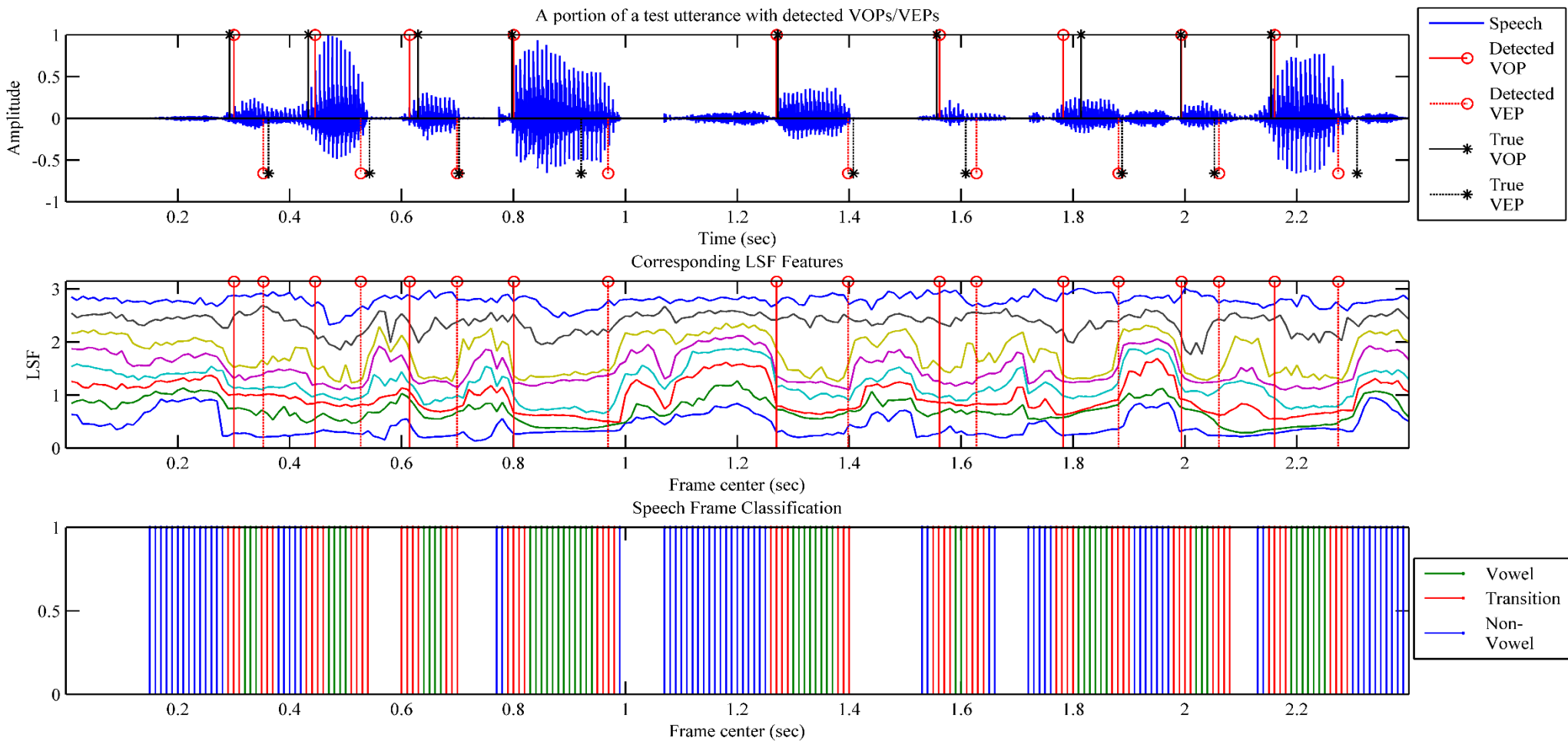2) Dynamic features - ΔLSF
3) Score Level Fusion -  LSF + ΔLSF

$$\Lambda_f(X) = \alpha \Lambda_{LSF}(X) + (1-\alpha)\Lambda_{\Delta LSF}(X)$$
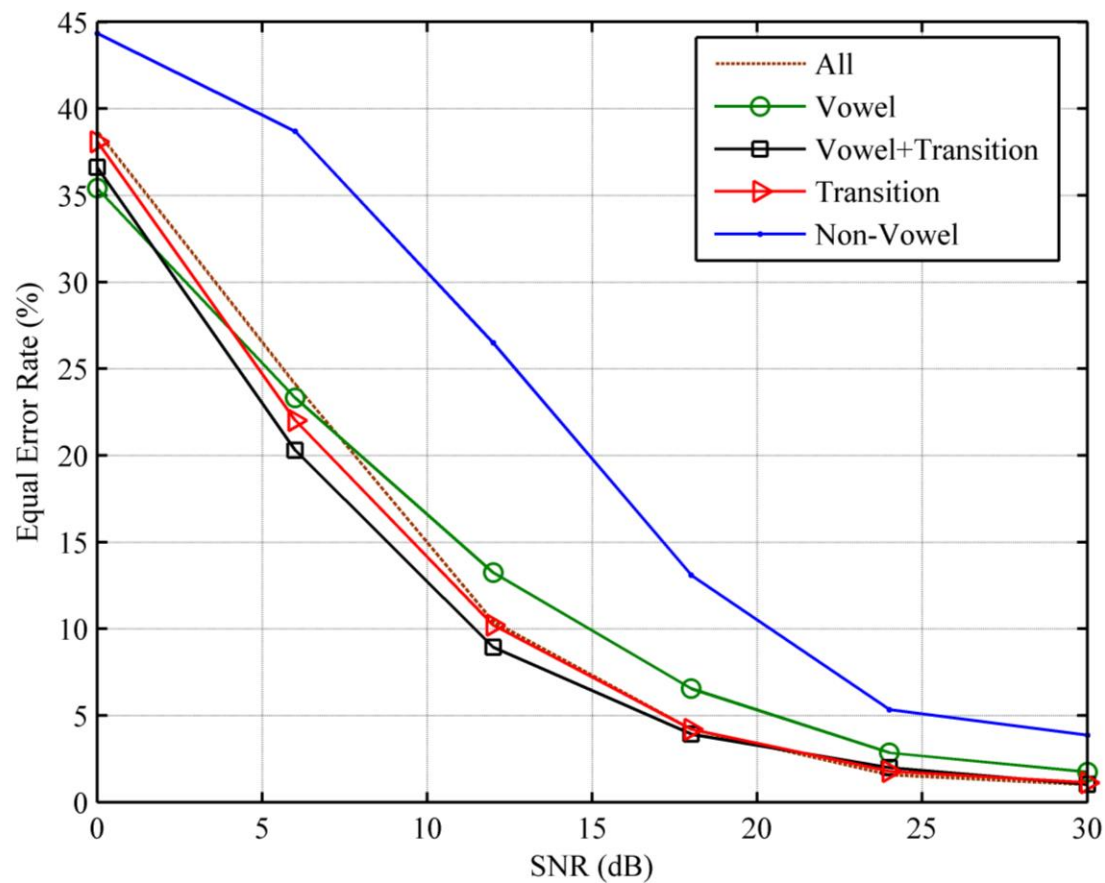


Performance of Score-level Fusion based SV system

- Baseline EER=0.86%
- Score-level fusion improves performance under noise

Wireless @ Virginia Tech

# Discriminative Power of Speech Zones



A portion of a test utterance with detected VOPs/VEPs

Corresponding LSF Features
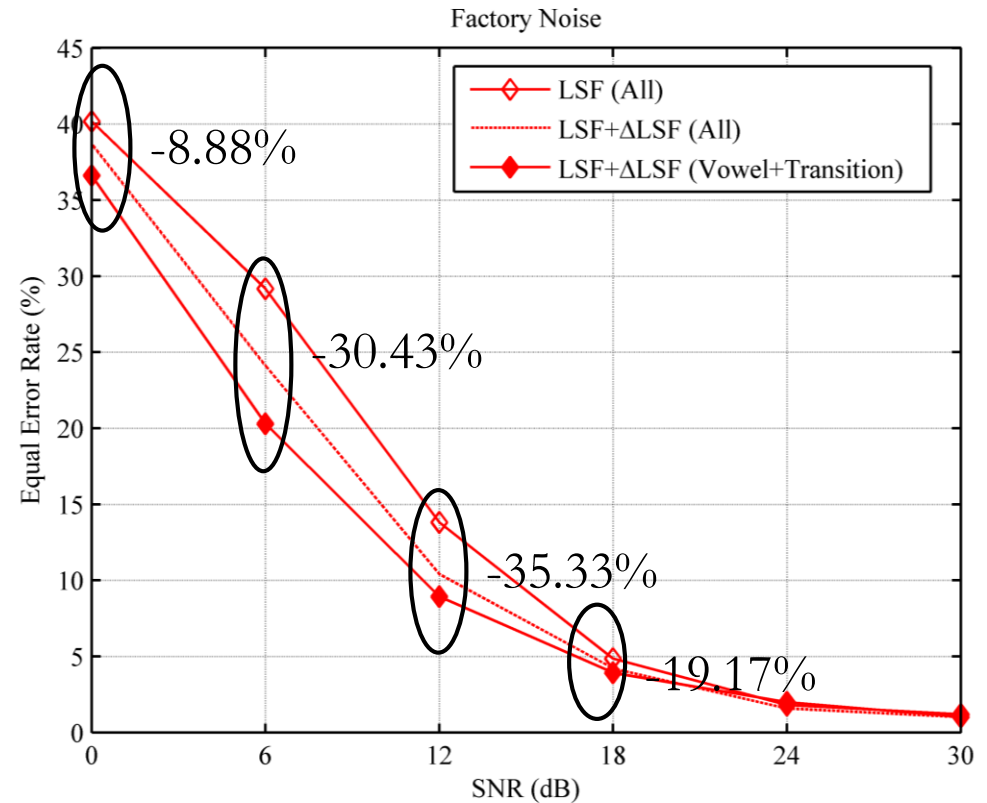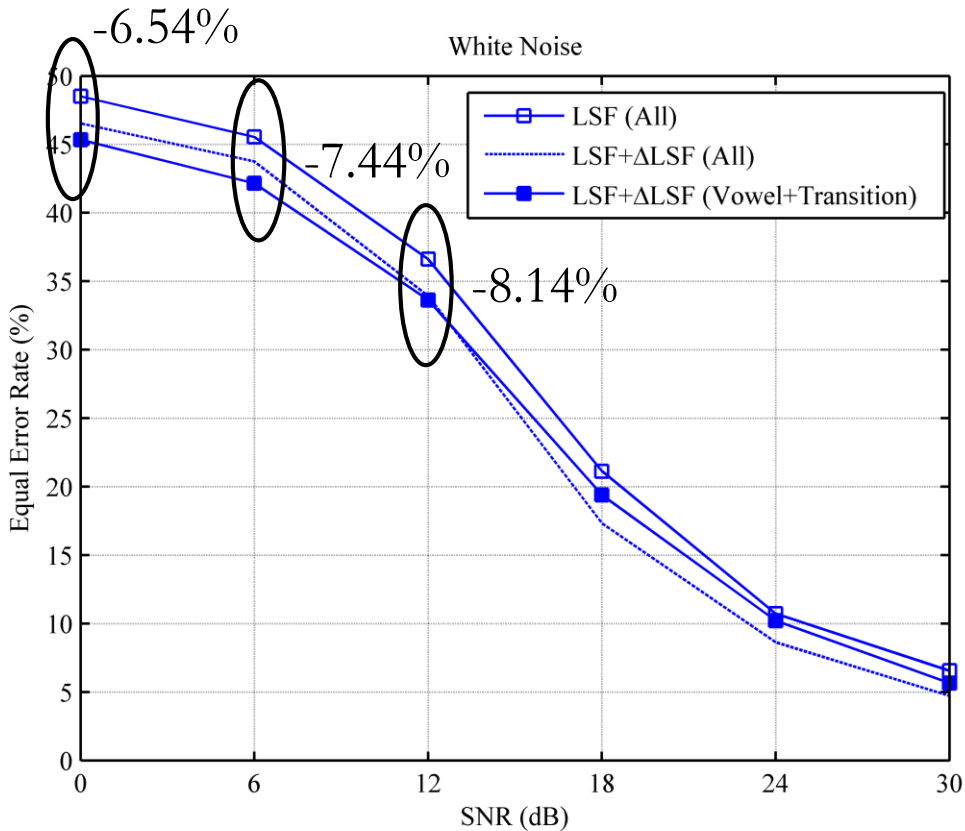
Speech Frame Classification

# Discriminative Power of Speech Zones

$X_{tr}$ is the set of features from transition frames $\Longrightarrow$ $\Phi_{s,X_{tr}} = \sum_{\mathbf{x} \in X_{tr}} \log p(\mathbf{x} \,|\, \lambda_s)$



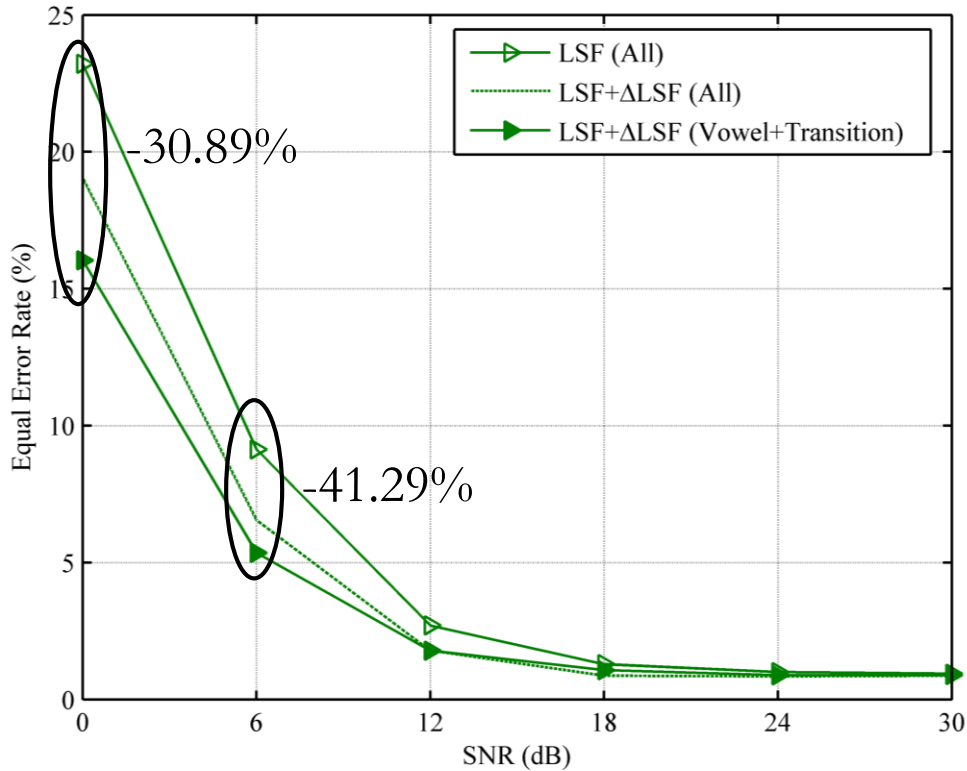- High SNR – transitions are most speaker discriminative
- Low SNR – vowels are most speaker discriminative
- Frame-level selection- not much benefit in high SNR
- Scoring on **vowel + transition** frames improves performance in low SNR.

Wireless @ Virginia Tech
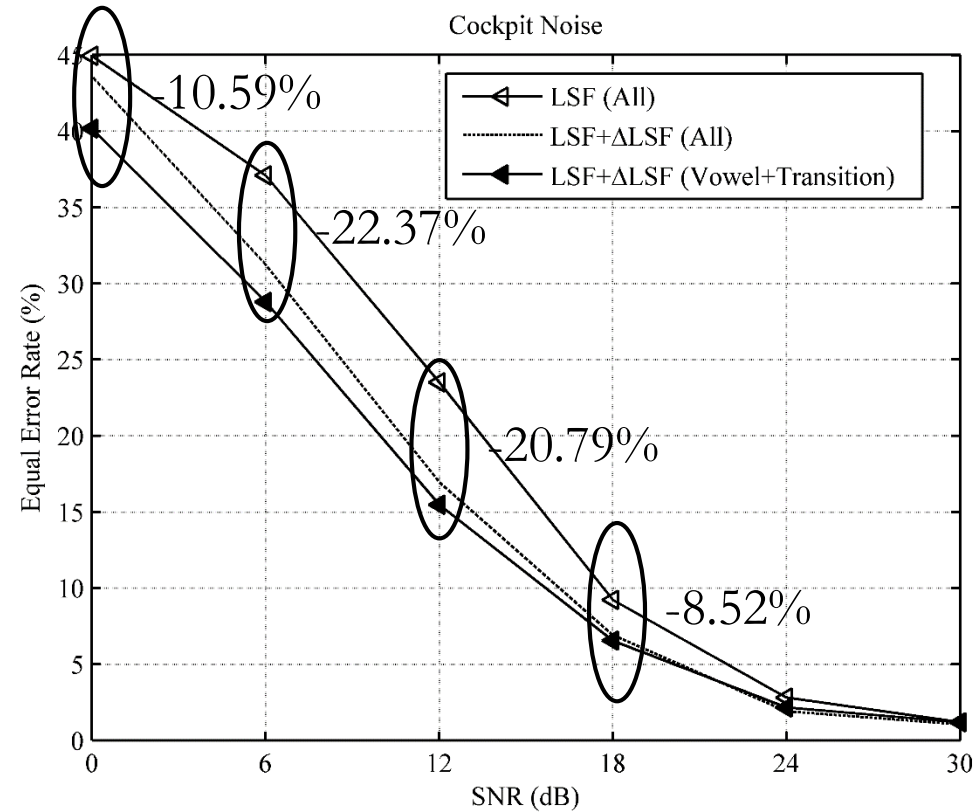
# Performance Improvement

# Performance Improvement

# Outline

- Introduction
- Motivation
- Feature Extraction
- Speaker Verification Framework
- Speaker Verification Results
- **Conclusions**

*Wireless* @ Virginia Tech

# Conclusions

- An automatic, text-independent speaker verification (SV) system was developed using Line Spectral Frequency (LSF) features.
- The performance of the SV system was evaluated under noise.
- Score-level fusion was used to combine complementary information from static and dynamic LSF features.
- Speaker-discriminative power of vowel, transition and non-vowel regions were investigated.
- Transition regions are the most speaker-discriminative under high SNR conditions
- High-energy vowel regions are most speaker-discriminative under low SNR conditions.
- Under noisy conditions, the performance of the score-level fusion based SV systems can be improved substantially by scoring exclusively on a combination of transition and vowel frames.
- Future work
  - Investigate the effect of training speaker models using transition zones.
  - Improve the algorithm used to localize transition zones.

Wireless @ )))Virginia Tech