# The 2006 AFRL/HEC Speaker Recognition Systems

**Raymond E. Slyh**
**Human Effectiveness Directorate**
**Air Force Research Laboratory**

---

# Team Members

- **Eric Hansen, Human Effectiveness Directorate, Air Force Research Laboratory (AFRL/HEC)**

- **Brian Ore, General Dynamics Advanced Information Systems**

- **Raymond Slyh, Human Effectiveness Directorate, Air Force Research Laboratory (AFRL/HEC)**
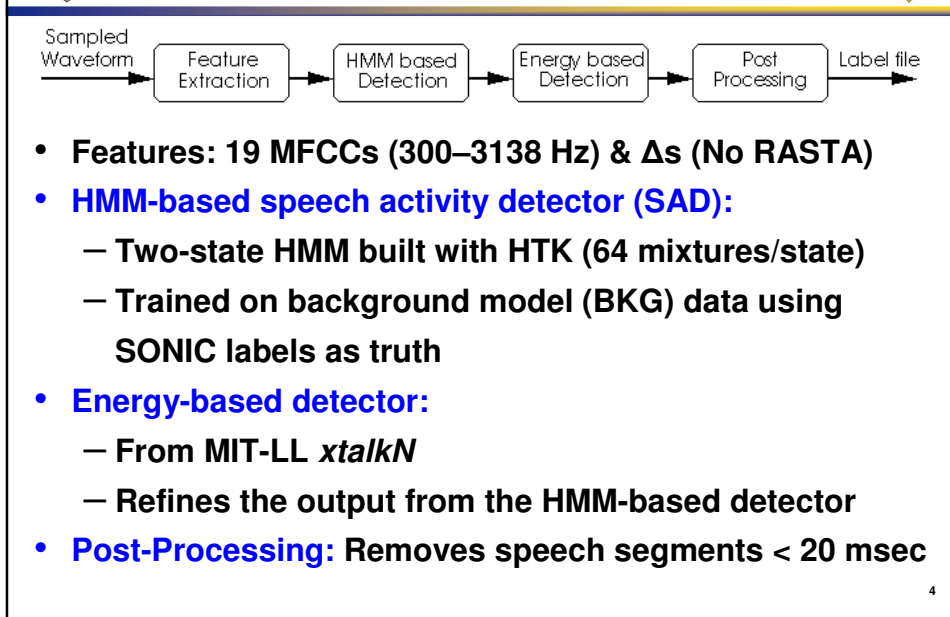
# Components of Submitted Systems

| Conditions | | TESTING | | |
|---|---|---|---|---|
| | | 10sec4w | 1conv4w | 1conv2w |
| **T R A I N I N G** | 10sec4w | FMBWF0/GMM MFCC/GMM MFCC/SVM | | |
| | 1conv4w | FMBWF0/GMM MFCC/GMM MFCC/SVM | FMBWF0/GMM MFCC/GMM MFCC/SVM MFCC/PS-GMM WLM | FMBWF0/GMM MFCC/GMM MFCC/SVM |
| | 3conv4w | FMBWF0/GMM MFCC/GMM MFCC/SVM | FMBWF0/GMM MFCC/GMM MFCC/SVM MFCC/PS-GMM WLM | FMBWF0/GMM MFCC/GMM MFCC/SVM |
| | 8conv4w | FMBWF0/GMM MFCC/GMM MFCC/SVM | FMBWF0/GMM MFCC/GMM MFCC/SVM MFCC/PS-GMM WLM | FMBWF0/GMM MFCC/GMM MFCC/SVM |
| | 3conv2w | | FMBWF0/GMM MFCC/GMM MFCC/SVM | FMBWF0/GMM MFCC/GMM MFCC/SVM |

**KEY**

**FMBWF0:** F1–F3, BW1– BW3, log(F0)

**MFCC:** Mel-Frequency Cepstral Coeffs & Δs

**GMM:** Gaussian Mixture Models

**SVM:** Support Vector Machines

**PS-GMM:** Phoneme- Specific GMMs

**WLM:** Language Modeling on Transcripts Output by SONIC

3

---

# MFCC/HMM SAD

Sampled Waveform → Feature Extraction → HMM based Detection → Energy based Detection → Post Processing → Label file

- **Features: 19 MFCCs (300–3138 Hz) & Δs (No RASTA)**
- **HMM-based speech activity detector (SAD):**
  - **Two-state HMM built with HTK (64 mixtures/state)**
  - **Trained on background model (BKG) data using SONIC labels as truth**
- **Energy-based detector:**
  - **From MIT-LL *xtalkN***
  - **Refines the output from the HMM-based detector**
- **Post-Processing: Removes speech segments < 20 msec**

4

# GMM-Based Systems

- **Gaussian mixture models from MIT Lincoln Laboratory (MIT-LL) system with:**
    - **2048 mixtures per model**
    - **Diagonal covariance matrices**
- **T-norm applied to output scores**
- **(Initial) speaker & T-norm models built using MAP adaptation from BKG with:**
    - **Relevance factor of 16**
    - **Only mixture means adapted**

5

# GMM-Based Systems: Models

- **BKG:**
    - **16 hours of data balanced for gender & channel from:**
        - **NIST 2001–2003 Evals (digital cell, electret, & carb.)**
        - **OGI National Cellular Corpus (for analog cellular)**
    - **Gender/channel models used for feature mapping**
- **T-norm:**
    - **Other than 10sec4w training:**
        - **Gender-dependent: 120 models per gender**
        - **Single conversation sides from NIST 2001–2003 Evals**
    - **10sec4w training:**
        - **240 gender-independent models**
        - **First 30 sec of data from original set of models**

6

# MFCC/GMM System: Features

- **19 MFCCs every 10 msec with:**
  - **Bandwidth of 300–3138 Hz**
  - **No $0^{th}$ coefficient**
- **Applied RASTA filtering & calculated Δs of features**
- **Kept a frame if labeled as speech by MFCC/HMM SAD**
- **Applied feature mapping and mean & variance norm.**
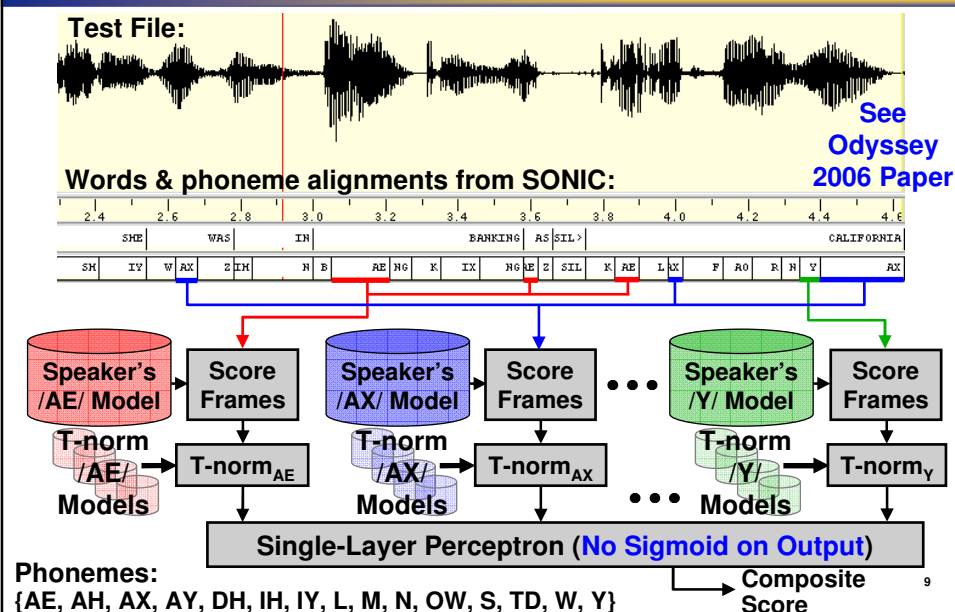
# FMBWF0/GMM System: Features

- **Every 10 msec:**
  - **Formant center frequencies (F1–F3) & bandwidths (BW1–BW3) using Snack toolkit from KTH**
  - **F0 & probability of voicing using *get_f0* from ESPS**
- **Kept a frame if:**

  **(speech) AND (voiced) AND (F0 < 250 Hz) AND {(F1, F2, F3) != (500 Hz, 1500 Hz, 2500 Hz)}**

- **Converted F1–F3 & BW1–BW3 to radians & took log(F0)**
- **Applied feature mapping (with channel picked by MFCCs)**

## MFCC/PS-GMM System

**Test File:**

**Words & phoneme alignments from SONIC:**

**See Odyssey 2006 Paper**

| 2.4 | 2.6 | 2.8 | 3.0 | 3.2 | 3.4 | 3.6 | 3.8 | 4.0 | 4.2 | 4.4 | 4.6 |

SHE | WAS | IN | BANKING | AS | SIL > | CALIFORNIA

SH | IY | W | AX | Z | IH | N | B | AE | NG | K | IX | NG | AE | Z | SIL | K | AE | L | AX | F | AO | R | N | Y | AX

**Speaker's /AE/ Model** → **Score Frames**

**T-norm /AE/ Models** → **T-norm$_{AE}$**

**Speaker's /AX/ Model** → **Score Frames**

**T-norm /AX/ Models** → **T-norm$_{AX}$**

• • •

**Speaker's /Y/ Model** → **Score Frames**

**T-norm /Y/ Models** → **T-norm$_{Y}$**

• • •

**Single-Layer Perceptron (No Sigmoid on Output)**

→ **Composite Score**

**Phonemes:**
**{AE, AH, AX, AY, DH, IH, IY, L, M, N, OW, S, TD, W, Y}**

9

---

## MFCC/SVM System

- **Features as in MFCC/GMM system**
- **Support vector machine classifier:**
  - **Generalized linear discriminant sequence kernel**
  - **From MIT-LL speech tools**
- **T-norm applied to scores (with T-norm models built using same data as for GMM systems)**

10

# WLM System

- **Used (English) transcripts generated by SONIC**
- **Pseudo sentence breaks were added**
- **Bigram language models with back-off**
- **CMU-Cambridge Language Modeling Toolkit with top 20,000 words, Witten-Bell discounting, & zero cut-offs**
- **Score a test file vs. claimant model as:**

$$\frac{1}{K}\sum_{k=1}^{K}\log(\mathrm{Pr}_{\mathrm{Claimant}}(k)) - \log(\mathrm{Pr}_{\mathrm{Background}}(k))$$

  **where K is the number of matching bigrams**

- **100 gender-independent two-conversation T-norm models from SWB II**

11

---

# Splitting NIST 2004 Control Files

**An Original NIST 2004 Control File**

**Split Into 10 Pieces**

**Testing file for split i: Let $S_{T,i}$ be the set of all speakers of the test files and target models**

**Sort Based on Test File Speaker Identities**

**Make "Disjoint" Train File**

**Make "Disjoint" Train File**

**Training file for split i: Let $S_{R,i}$ be the set of all speakers of the test files and target models**

**Disjoint: $S_{T,i} \cap S_{R,i} = \emptyset$**

12

# Four-Wire Fusion & Thresholds

- **For each split:**
  - **Built a single-layer perceptron (SLP) on training file**
  - **Applied SLP to system scores for the test file**
- **Concatenated score files for the ten splits**
- **Determined threshold for minDCF (this was the threshold used for the 2006 Eval)**
- **Built new SLP over the entire control file for the condition (this was the SLP used for the 2006 Eval)**
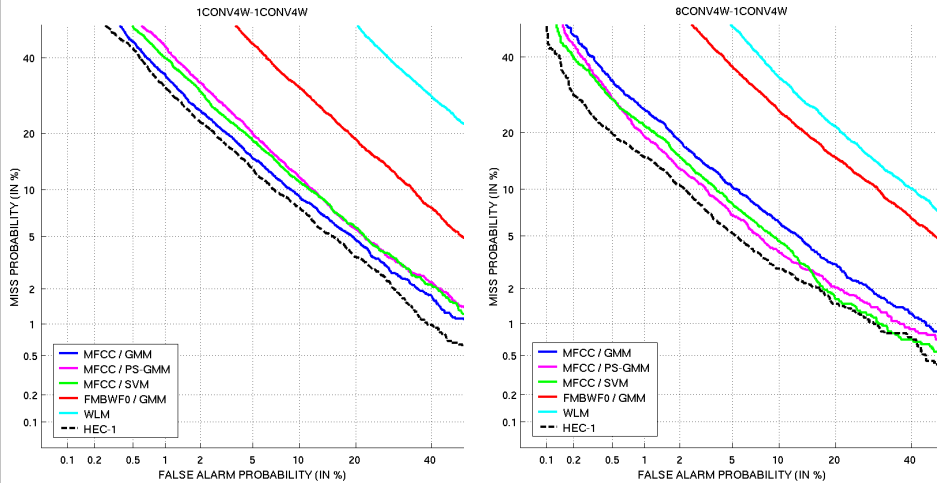- **SLPs built using LNKnet from MIT-LL**

13

---

# 10sec4w Testing



**Score combination provided considerable benefit for 10sec4w training but less benefit for larger amounts of training data**

14

# 1conv4w Testing



**MFCC/PS-GMM system outperformed MFCC/GMM system for 8conv4w training even though it used only 15 out of 50 *English* phonemes**

---

# Unsupervised Adaptation

- **HEC-2 System: MFCC/GMM system with & without unsupervised adaptation (UA) of mixture means, $\overline{\mu}_m$ :**

$$\overline{\mu}_m^{\text{NEW}} = \alpha_m E_m(X) + (1 - \alpha_m)\overline{\mu}_m \qquad E_m(X): \text{ Mean of vectors prob. assigned to mixture } m$$

- **Initial speaker models built using MAP adaptation from BKG with:**

$$\alpha_m = \frac{n_m(X)}{n_m(X) + r}$$

$n_m(X):$ **Probabilistic "count" of vectors in mixture $m$**

$r:$ **Relevance factor = 16**

- **Updated speaker model built using MAP adaptation from current speaker model with:**

$$\alpha_m = \begin{cases} 0.1, & \beta < 0.1 \\ 0.5, & 0.5 < \beta \\ \beta, & \text{otherwise} \end{cases} \qquad \beta = \frac{T_T}{T_T + T_M}$$

$T_T:$ **# speech frames in test file**
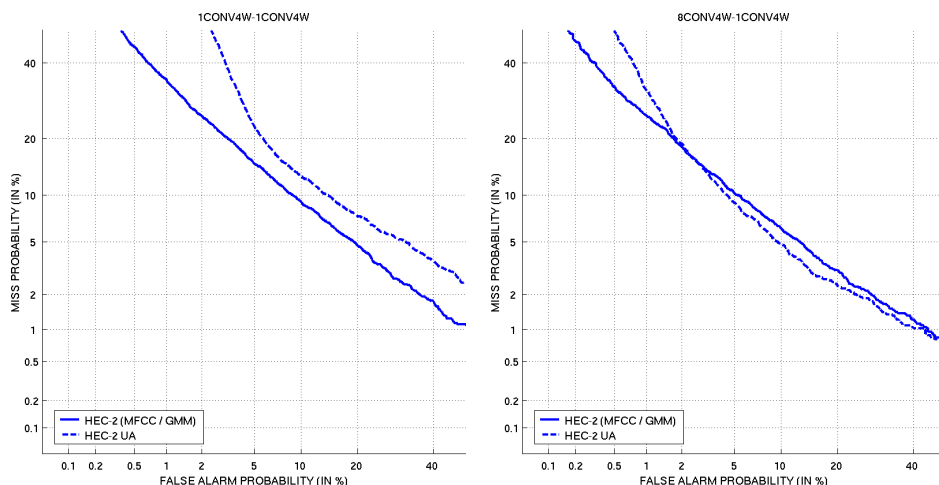
$T_M:$ **# speech frames used for current model**

- **See Odyssey 2006 paper for more details**
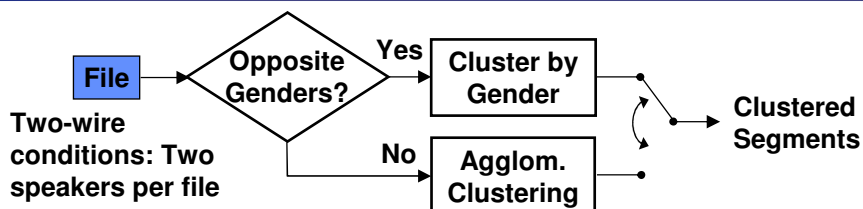
# Unsupervised Adaptation



- **Model updating threshold: minDCF threshold from NIST 2004 Eval data**
- **UA degraded performance: Need a different updating threshold?**

17

---

# Two-Wire Segmentation/Clustering



- **1conv2w testing:**
  - **If gender-based clustering used for a file: Test correct-gender cluster against target model**
  - **If agglomerative clustering used for a file: Cluster into three sets, test each set against the target model, & pick the highest score**
- **3conv2w training:**
  1) **Segment & cluster each of the three files individually**
  2) **Cluster across the three files**
  3) **Build model**

18

# Opposite-Gender Files & Clustering

- **Opposite-Gender File Determination:**
  - **MFCC/HMM SAD determines speech/non-speech segments**
  - **Score files against male, female, & BKG GMMs**
  - **If target speaker is male, label a file opposite-gender if:**

    **Score$_{BKG}$ – Score$_{Male}$ > Gender-dependent threshold**
  - **Similar procedure if target is female**
- **Gender-Based Clustering:**
  - **MFCCs, 300–3138 Hz, RASTA, Δs, but no feature mapping**
  - **Score each segment individually against male & female GMMs**
  - **Take top 90% of the segments of proper gender for target model**
- **See Odyssey 2006 paper for more details**

19

# Agglomerative Clustering Within File

- **MFCC/HMM SAD determined speech/non-speech segments**
- **64-mixture GMM trained with all speech vectors from the file using MFCCs band limited to 200–2860 Hz and Δs, but without RASTA filtering, feature mapping, or mean & variance normalization**
- **Weights then adapted for each speech segment**
- **In each clustering stage:**
  - **Let $X$ and $Y$ be two segments, and let $Z = X \cup Y$**
  - **$\forall X, Y$ calculate:** $$\Lambda(X,Y) = \frac{L(Z|\theta_Z)}{L(X|\theta_X)L(Y|\theta_Y)}$$  $L(X|\theta_X)$: **Likelihood of data for segment X given model for X**
  - **Merge the $X$ and $Y$ segments with the highest $\Lambda(X,Y)$**
- **Repeat the process until three sets of segments are left (presumably, one for each speaker and a "garbage" set)**

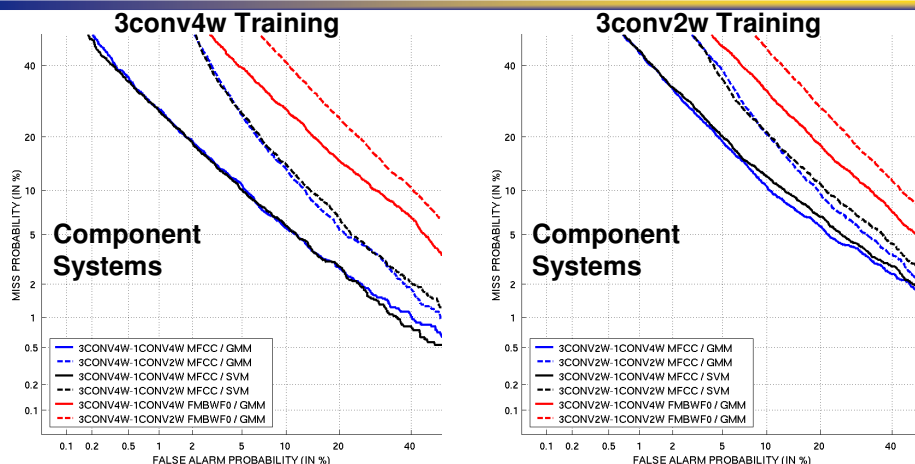20

## 3conv2w: Clustering Across Files

- **Features: 19 MFCCs with a bandwidth of 300–3138 Hz, RASTA, Δs, feature mapping, and mean & variance normalization**
- **If any files were segmented by gender:**
  - **Correct-gender segments used to build an initial speaker model**
  - **Segments from other files tested against the initial speaker model**
- **If no files were segmented by gender:**
  - **Models were built for each of the three segment sets in each file by using MAP adaptation of mixture means from BKG**
  - **Segments were scored against the models (from other files) & highest scoring segment/model pair was clustered**
  - **Segments from third file tested against the clustered model**

## Segmentation Results



- **Comparisons within a plot show effect of two-wire testing, while comparisons across the plots show the effect of two-wire training**
- **Substantial performance difference between 3conv4w & 3conv2w training and between 1conv4w & 1conv2w testing**

# Acknowledgements

- **MIT Lincoln Laboratory:**
  - **MFCC/GMM, MFCC/SVM, and feature mapping code**
  - **LNKnet**
- **Bryan Pellom, Univ. of Colorado at Boulder: SONIC speech recognizer & acoustic models**
- **Cambridge Univ.:**
  - **Statistical Language Modeling Toolkit (with CMU)**
  - **HTK**
- **KTH: Snack toolkit**

23