

The 2006 France Telecom Research and Development Center (Beijing) Speaker Recognition System



June 26, 2006

Xianyu Zhao, Yuan Dong, Hao Yang, Jian Zhao

France Telecom
Research & Development

D1 - 23/04/2007

Outline



- S Overview
- S Front End
- S Speaker Modeling
- S Score Normalization
- S Speaker Segmentation and Clustering
- S Conclusion

France Telecom
Research & Development

La communication de ce document est soumise à autorisation de la R&D de France Télécom
D2 - 23/04/2007

GMM-UBM based Text Independent Speaker Verification System

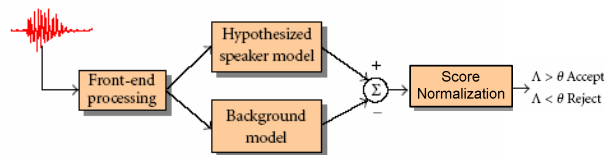


S Likelihood-ratio-based speaker verification system using GMM-UBM

Q Front End Feature extraction

Q Speaker (and alternative hypothesis) Modeling

Q Score normalization



France Telecom
Research & Development

La communication de ce document est soumise à autorisation de la R&D de France Télécom
D3 - 23/04/2007

Front-End Processing: Aim and Structure

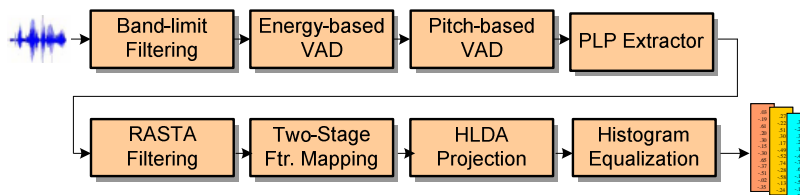


S The aim of robust front-end processing:

Q To reduce the influences of transmission channel, microphone type and environmental noises

Q To facilitate the following speaker modeling with GMM

S The FTRD SRE'06 Front-End



France Telecom
Research & Development

La communication de ce document est soumise à autorisation de la R&D de France Télécom
D4 - 23/04/2007

Front-End Processing: Some Details



S Two-Stage Feature Mapping

QStage-1 (Channel Mapping): Land, GSM, CDMA, Cellular, TDMA, Cordless

QStage-2 (Microphone Mapping): Speaker phone, Head phone, Handheld, Ear-bud

S HLDA Projection

QOriginal Feature: PLP_0 + Delta + Delta-Delta + Delta-Delta-Delta (52-dim)

QProjected Feature: 51-dim

S Histogram Equalization

QReference Cumulative Distribution Function (CDF) was estimated using all the development data

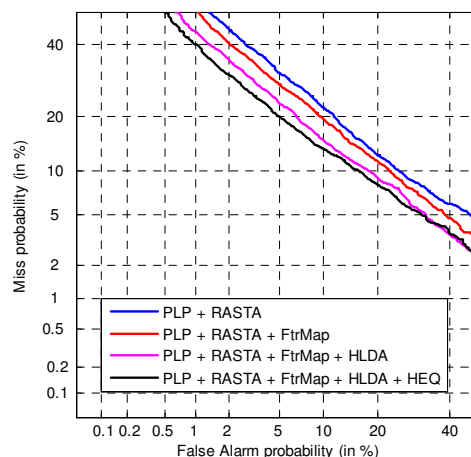
France Telecom
Research & Development

La communication de ce document est soumise à autorisation de la R&D de France Télécom
D5 - 23/04/2007

Front-End Processing: Performance



S NIST SRE04 1side-1side



France Telecom
Research & Development

La communication de ce document est soumise à autorisation de la R&D de France Télécom
D6 - 23/04/2007



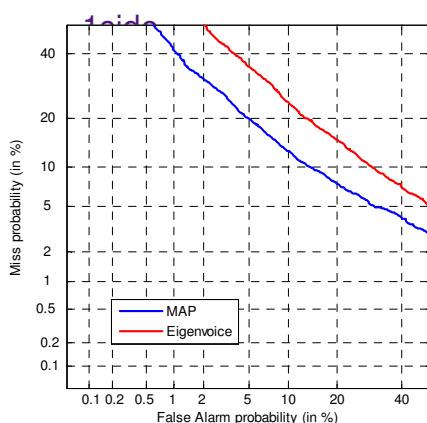
Speaker Modeling

- S The UBM model is a gender independent GMM with 2048 mixture components trained using about 40 hours of data from the Switchboard and SRE'2004 evaluation database.
- S Speaker models are obtained by adapting from the UBM with their individual training data through
 - Q Bayesian learning approach (Maximum a Posteriori adaptation): good asymptotic behavior
 - Q Speaker clustering based approach (Eigenvoice adaptation): rapid adaptation

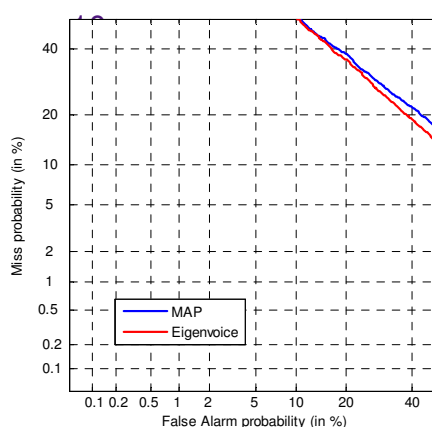


Speaker Modeling: Performance

S NIST SRE04 1side-



S NIST SRE04 10sec-





Score Normalization

- S Each verification score is normalized by subtracting the mean and then dividing by the standard deviation of imposter score distribution.
- S Estimation of imposter score distribution:
 - QTNorm: Each utterance is compared against to a set of imposter models to estimate the imposter score distribution.
 - QATNorm: Similar to TNorm, but the set of imposter models is dependent on the hypothesized speaker of a test segment.



Score Normalization: some details

- S Multi-language imposter pool

	English	Mandarin	Arabic	Russian	Spanish	Total
Male	268	31	32	14	9	354
Female	349	24	27	31	37	468
Total	617	55	59	45	46	822

- S TNorm: for male segment speaker, all 354 male imposters are used; for female segment speaker, all 468 female imposters are used.



Score Normalization: some details

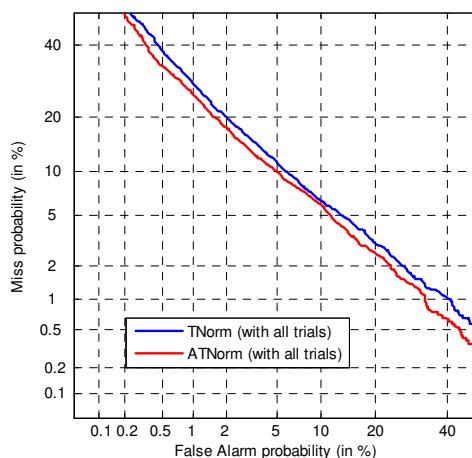
- S ATNorm: for each male target speaker, 55 nearest imposters are selected from the male imposter pool; for each female target speaker, 55 nearest female imposters are selected from the pool.
- S The following cross-model log likelihood ratio distance is used to define the neighborhood of two models

$$d(\lambda_i, \lambda_j) = -\frac{1}{N_i} \log \left(\frac{p(x_i | \lambda_j)}{p(x_i | \lambda_{UBM})} \right) - \frac{1}{N_j} \log \left(\frac{p(x_j | \lambda_i)}{p(x_j | \lambda_{UBM})} \right)$$



Score Normalization: Performance

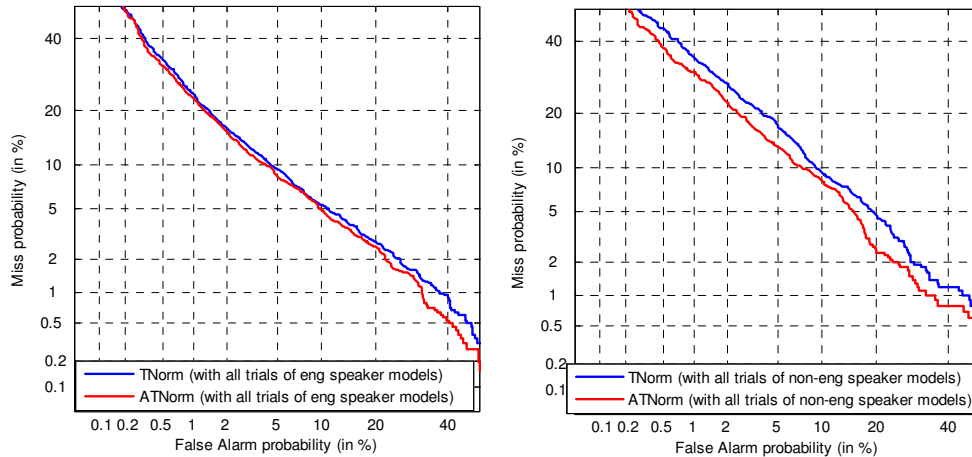
- S NIST SRE06 1conv4w-1conv4w test



Score Normalization: Performance



S NIST SRE06 1conv4w-1conv4w test



France Telecom
Research & Development

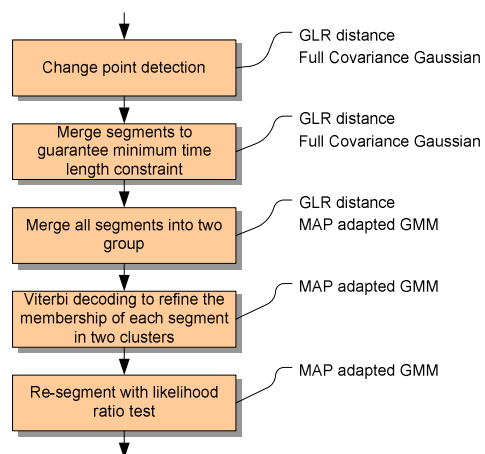
La communication de ce document est soumise à autorisation de la R&D de France Télécom
D13 - 23/04/2007

Speaker Segmentation



- S MFCC + energy, no channel compensation
- S GLR distance is used for change point detection and agglomerative clustering
- S Two-Stage agglomerative clustering: as more data are in cluster, more complicated models are used

France Telecom
Research & Development

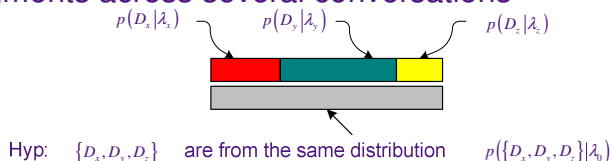


La communication de ce document est soumise à autorisation de la R&D de France Télécom
D14 - 23/04/2007



Speaker Clustering

- S MFCC + Delta + RASTA: reduce channel mismatch across different conversations
- S A form of GLR distance is employed to cluster multiple segments across several conversations



$$GLR(D_x, D_y, D_z) = \log \left(\frac{p(D_x | \lambda_x) \cdot p(D_y | \lambda_y) \cdot p(D_z | \lambda_z)}{p(\{D_x, D_y, D_z\} | \lambda_0)} \right)$$

- S Gender of target speaker is taken into account to make decision

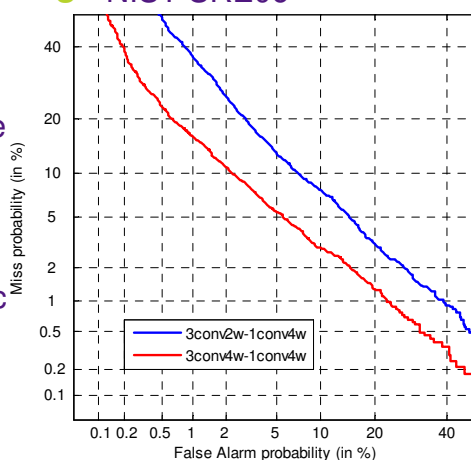
France Telecom
Research & Development

La communication de ce document est soumise à autorisation de la R&D de France Télécom
D15 - 23/04/2007



Summed Channel SRE Performance

- S There is still large gap between 2-wire and 4-wire training; our speaker segmentation/custering process need improvements



France Telecom
Research & Development

La communication de ce document est soumise à autorisation de la R&D de France Télécom
D16 - 23/04/2007

Conclusion



- S It is the first time that FTRD (Beijing) participates NIST SRE and we present a simple GMM-UBM likelihood ratio based system.
- S Some improvements have been made on front-end, model adaptation and score normalization processes.
- S To catch up with the state-of-art, there are still a lot of things to do, e.g. prosodic/idiolectal features, session variability modeling, cooperation with ASR, etc.

Thanks



- S To NIST for providing the wonderful evaluation platform!
- S To all pioneers in speaker recognition/verification research area from who/whose publications we have learned a lot!