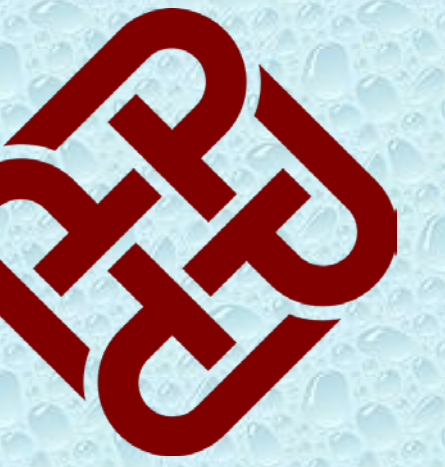


Likelihood-Ratio Empirical Kernels for I-Vector Based PLDA-SVM Scoring



Man-Wai MAK and Wei RAO

Dept. of Electronic and Information Engineering, The Hong Kong Polytechnic University



Summary

- Likelihood ratio (LR) scoring in PLDA speaker verification systems only uses the information of background speakers **implicitly**. This paper exploits the notion of empirical kernel maps to incorporate background speaker information into the LR scoring process **explicitly**.
- We train a scoring SVM for each target speaker based on empirical kernels.
- We demonstrate that a number of target-speaker i-vectors can be generated by an utterance partitioning and resampling technique, resulting in much better scoring SVMs.
- Results suggest that incorporating background speaker information into PLDA scoring through speaker-dependent SVMs can boost the performance of i-vector based PLDA systems significantly.

Methods

• Likelihood Ratio Score for Gaussian PLDA

Given a set of D -dim length-normalized i-vectors:

$$X = \{x_{ij}; i = 1, \dots, N, j = 1, \dots, H_i\},$$

we aim to estimate the latent variables $Z = \{z_i; i = 1, \dots, N\}$ and parameters $\omega = \{\mu, W, \Sigma\}$ of a factor analyzer:

$$x_{ij} = \mu + Wz_i + \varepsilon_{ij}$$

Given a test i-vector x_t and target-speaker's i-vector x_s , the likelihood ratio (LR) score is:

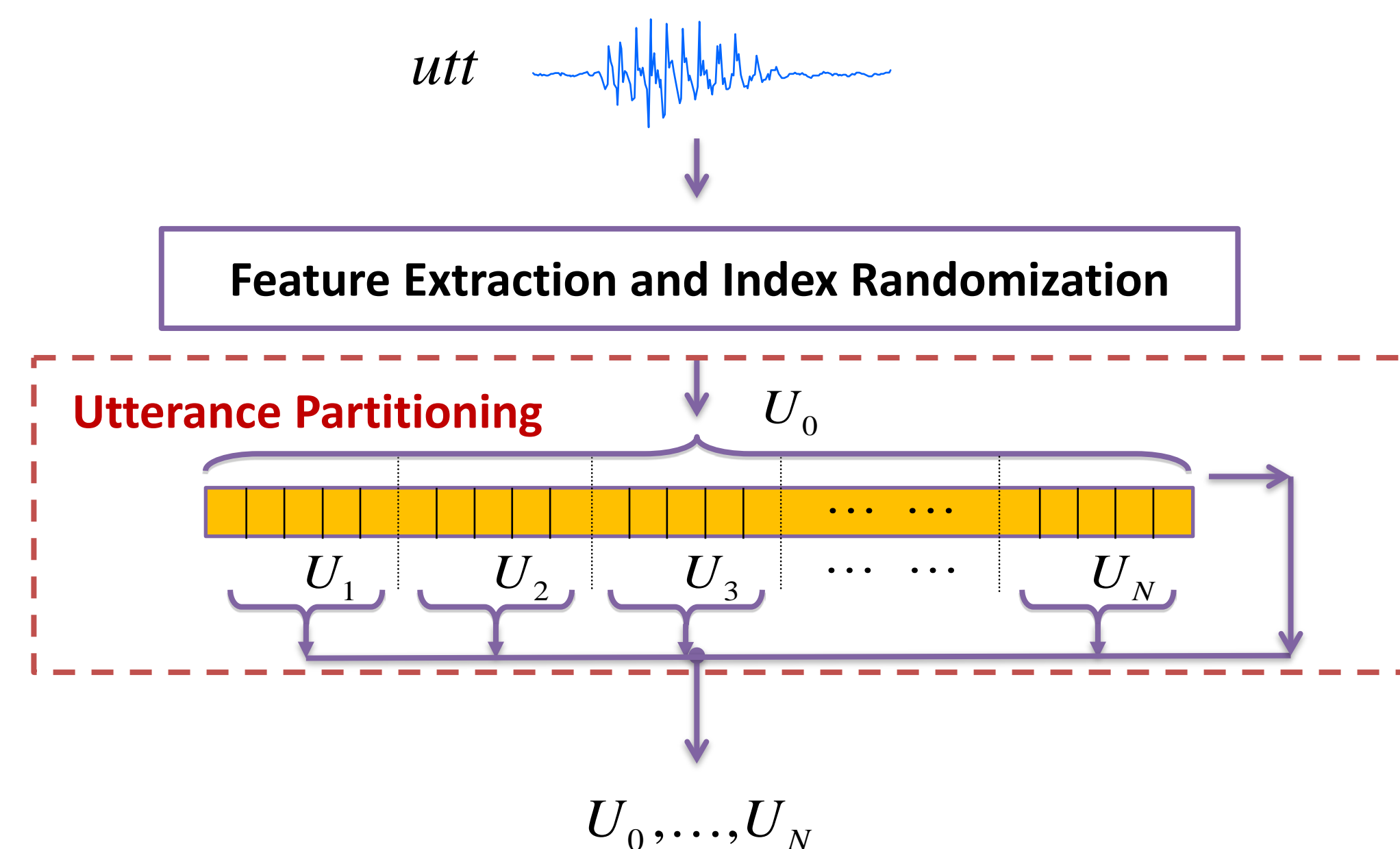
$$S_{LR}(x_s, x_t) = \text{const} + x_s^T Q x_s + x_t^T Q x_t + 2x_s^T P x_t$$

where

$$P = \Lambda \Gamma (\Lambda - \Gamma \Lambda^{-1} \Gamma)^{-1}; \Lambda = W W^T + \Sigma; Q = \Lambda^{-1} - (\Lambda - \Gamma \Lambda^{-1} \Gamma)^{-1}; \Gamma = W W^T$$

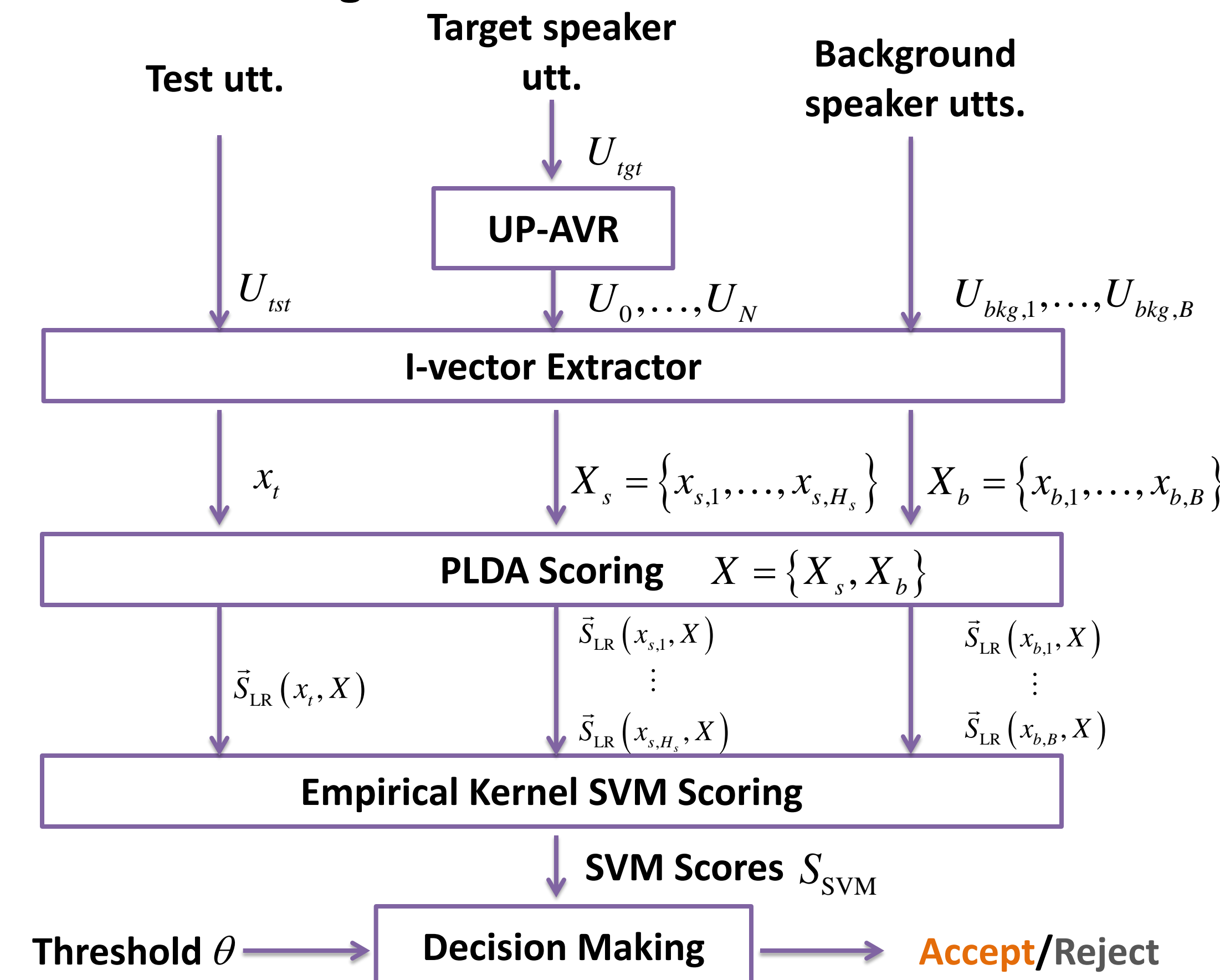
• UP-AVR

Background information



Methods

• PLDA-SVM Scoring



$$S_{SVM}(x_t, X_s, X_b) = \sum_{j \in SV_s} \alpha_{s,j} K(x_t, x_{s,j}) - \sum_{j \in SV_b} \alpha_{b,j} K(x_t, x_{b,j}) + d_s$$

Empirical LR Kernel I:

$$K(x_t, x_{s,j}) = \mathbb{K}(\vec{S}_{LR}(x_t, X_s), \vec{S}_{LR}(x_{s,j}, X_s))$$

where

$$X_s = \{x_{s,1}, \dots, x_{s,H_s}\}; \vec{S}_{LR}(x_t, X_s) = \begin{bmatrix} S_{LR}(x_t, x_{s,1}) \\ S_{LR}(x_t, x_{s,2}) \\ \vdots \\ S_{LR}(x_t, x_{s,H_s}) \end{bmatrix}$$

Empirical LR Kernel II:

$$K(x_t, x_{s,j}) = \mathbb{K}(\vec{S}_{LR}(x_t, X), \vec{S}_{LR}(x_{s,j}, X))$$

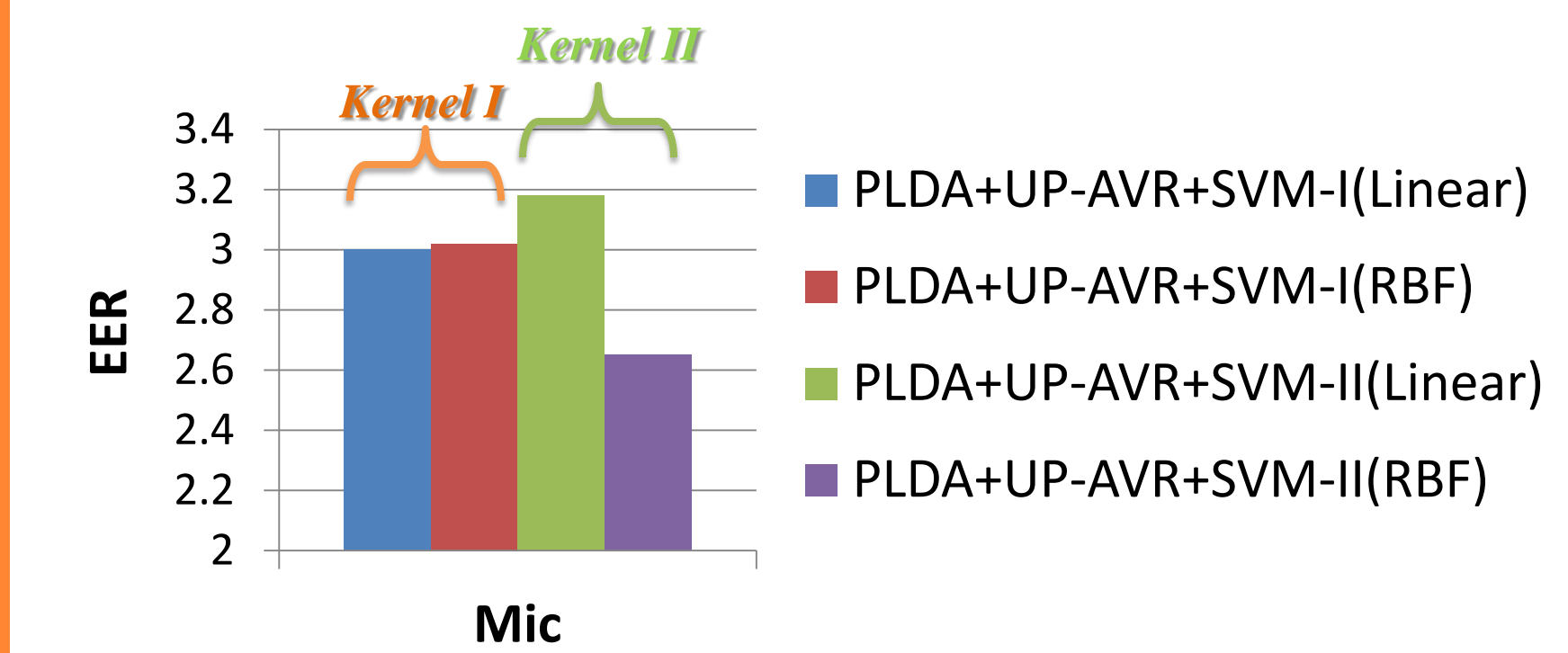
where

$$X = \{X_s, X_{b'}\}; \vec{S}_{LR}(x_t, X) = \begin{bmatrix} S_{LR}(x_t, x_{s,1}) \\ \vdots \\ S_{LR}(x_t, x_{s,H_s}) \\ S_{LR}(x_t, x_{b,1}) \\ \vdots \\ S_{LR}(x_t, x_{b,B'}) \end{bmatrix}$$

$$X_{b'} = \{x_{b,1}, \dots, x_{b,B'}\};$$

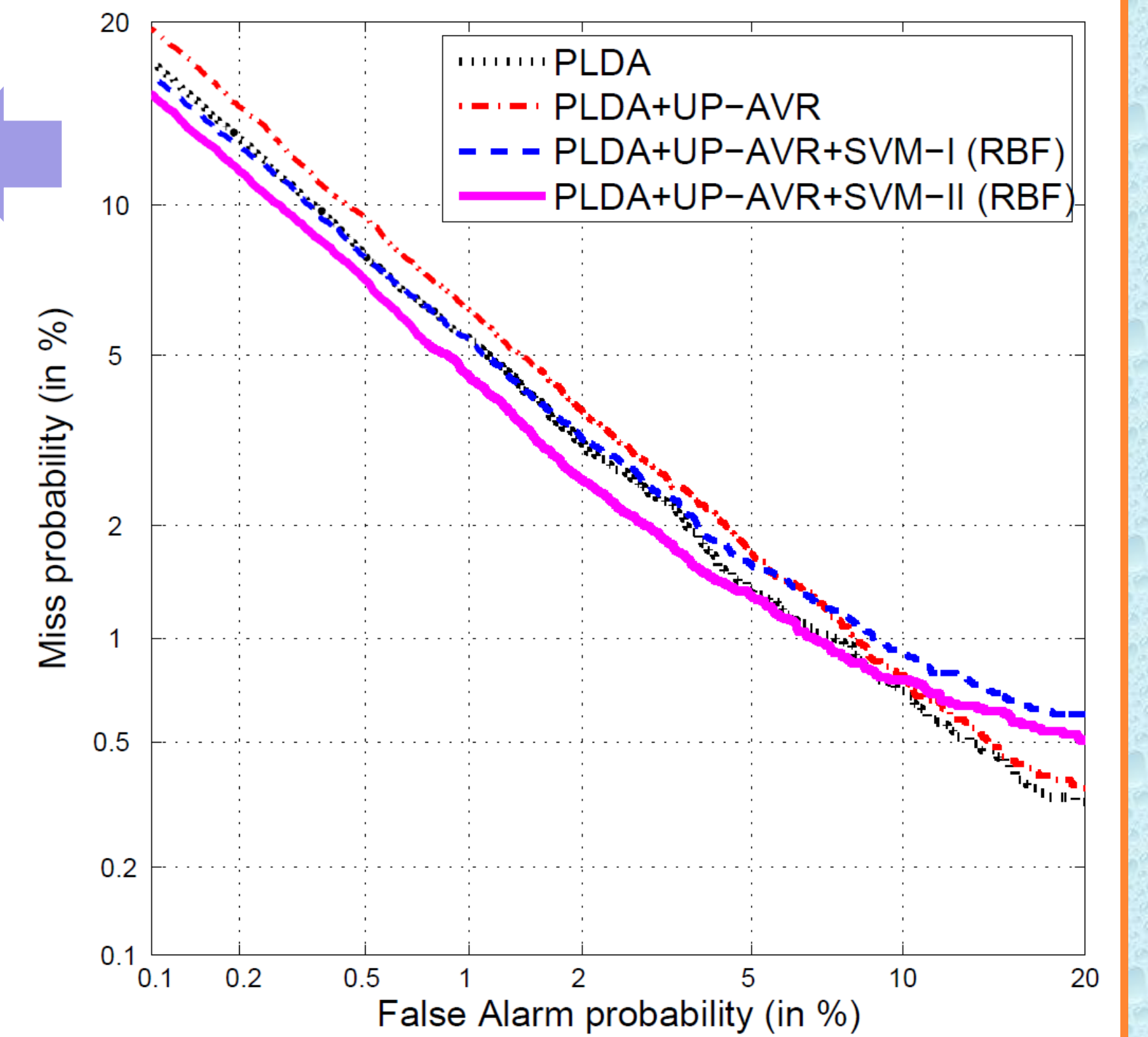
$$B' \leq B;$$

Results

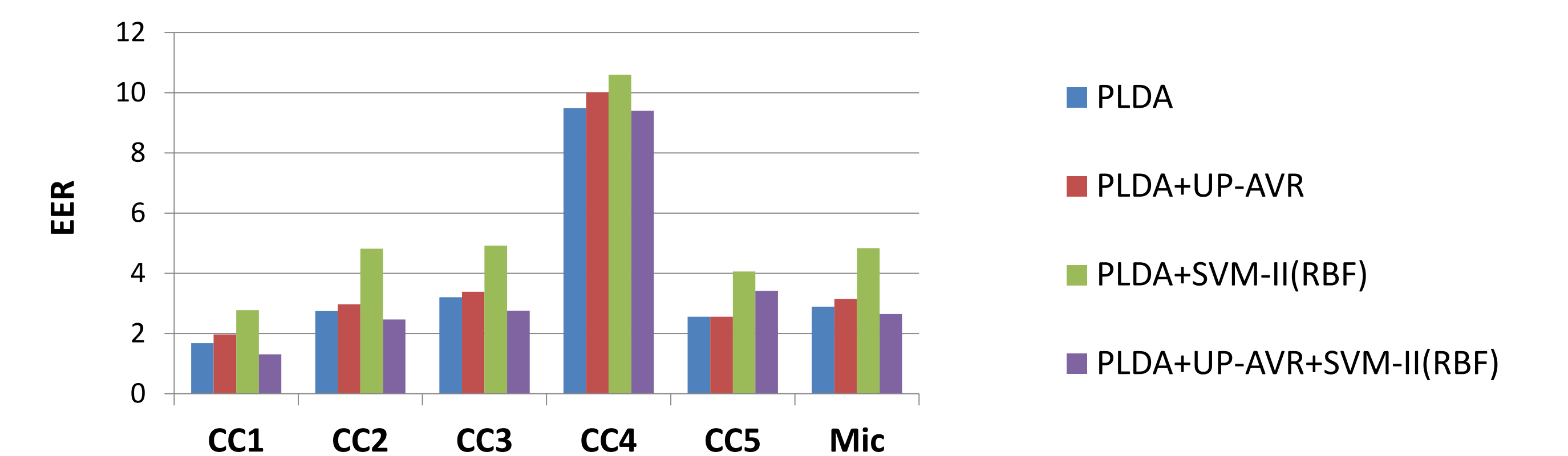


➤ Linear kernel performs better in the empirical LR kernel I.
➤ Non-linear kernel (RBF) performs better in the empirical LR kernel II.

PLDA+UP-AVR+SVM-II(RBF) significantly outperforms PLDA LR scoring.



➤ UP-AVR not only helps to alleviate the data-imbalance problem in SVM training, but also enriches the information content of the scoring vectors by increasing the number of LR scores derived from the target speaker.
➤ UP-AVR is not beneficial to LR scoring.



References

- [1] P. Kenny, "Bayesian speaker verification with heavy-tailed priors," in *Proc. of Odyssey: Speaker and Language Recognition Workshop, Brno, Czech Republic*, June 2010.
- [2] D. Garcia-Romero and C.Y. Espy-Wilson, "Analysis of i-vector length normalization in speaker recognition systems," in *Interspeech'2011*, pp. 249–252, 2011.
- [3] W. Rao and M. W. Mak, "Boosting the performance of i-vector based speaker verification via utterance partitioning," *IEEE Trans. on Audio, Speech and Language Processing*, vol. 21, no. 5, pp. 1012–1022, 2013.