



UNSUPERVISED LEARNING OF ACOUSTIC FEATURES VIA DEEP CANONICAL CORRELATION ANALYSIS



UNIVERSITY of WASHINGTON

Weiran Wang¹

¹TTI-Chicago

Raman Arora²

²Johns Hopkins University

Karen Livescu¹

³University of Washington

Jeff A. Bilmes³

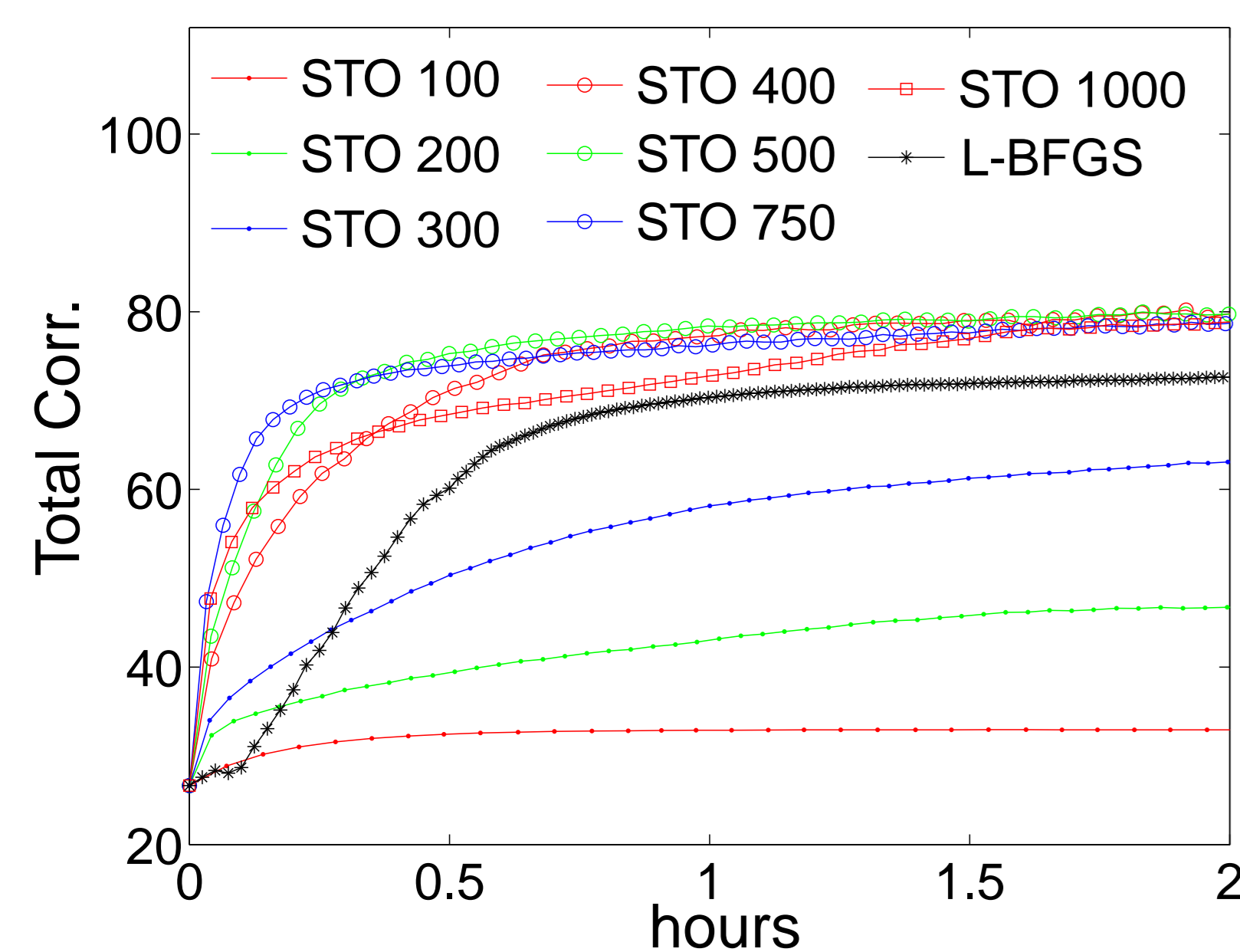
Presented by Sadaoki Furui

1 Overview

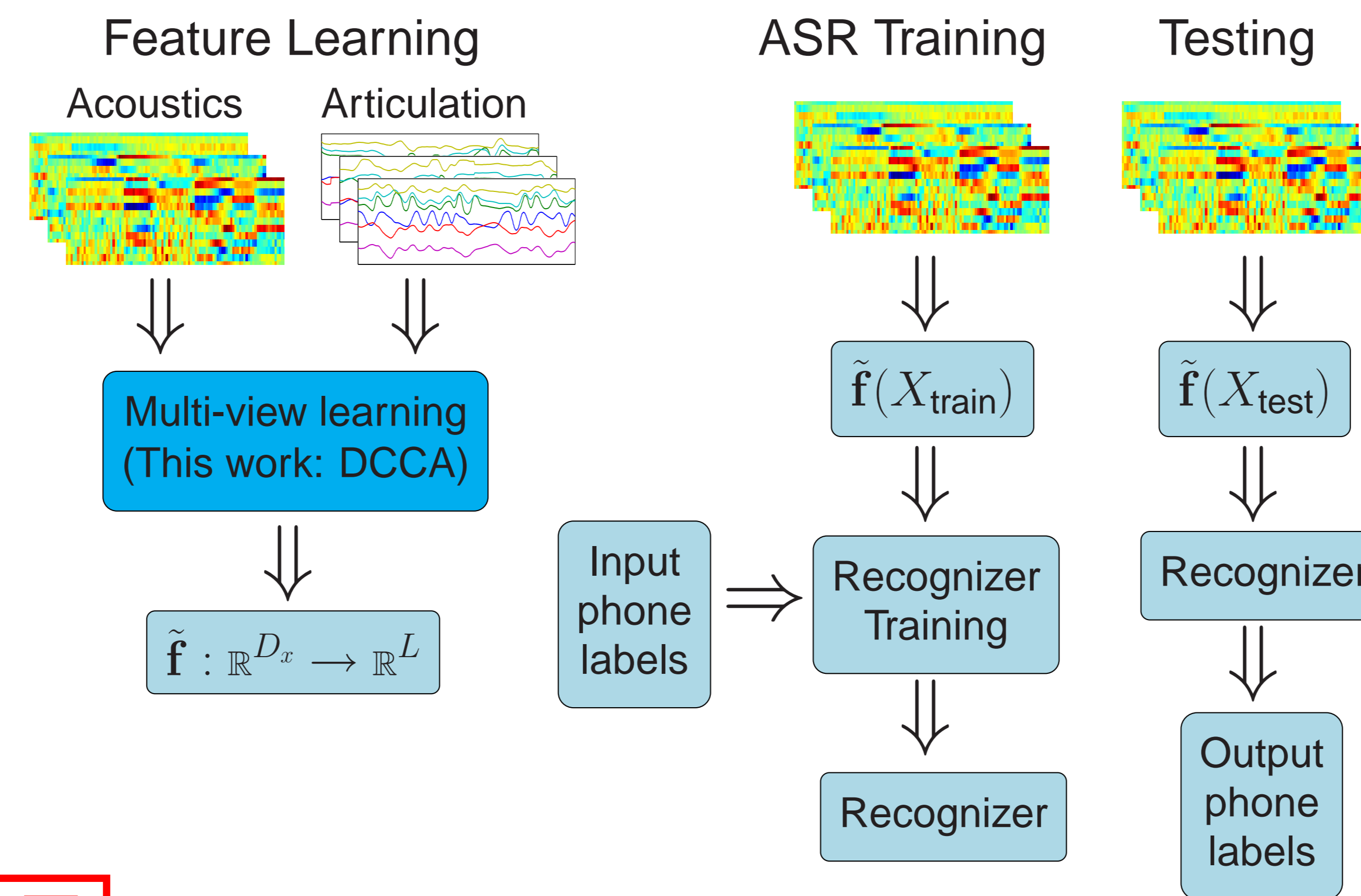
- Background**
 - Can we learn better acoustic features if we have access to multi-view data external to recognizer training/test data?
 - Here the views are acoustics and articulation.
 - Learned transformations are applied to **acoustics-only** data for recognizer training and testing.
 - Previous work: Improved recognition with transformations learned via canonical correlation analysis (CCA) and kernel CCA [1,2].
 - Intuition:
 - * 2nd view helps isolate signal from noise.
 - * Like articulatory inversion, but using latent articulatory space.
- This work**
 - We apply deep CCA (DCCA), where the feature mapping is a deep neural network (DNN) [3].
 - DCCA significantly improves phone recognition, without access to test speakers' articulatory data.
 - New stochastic optimization algorithm for large-scale DCCA.

4 Stochastic Optimization

- DCCA objective is a constrained loss that does not decompose over the training samples \Rightarrow not a good fit for stochastic gradient descent, but batch training is very slow.
- We use a minibatch stochastic approach with large minibatches for stable estimates of covariance matrices and gradient.



2 Multi-View Acoustic Feature Learning

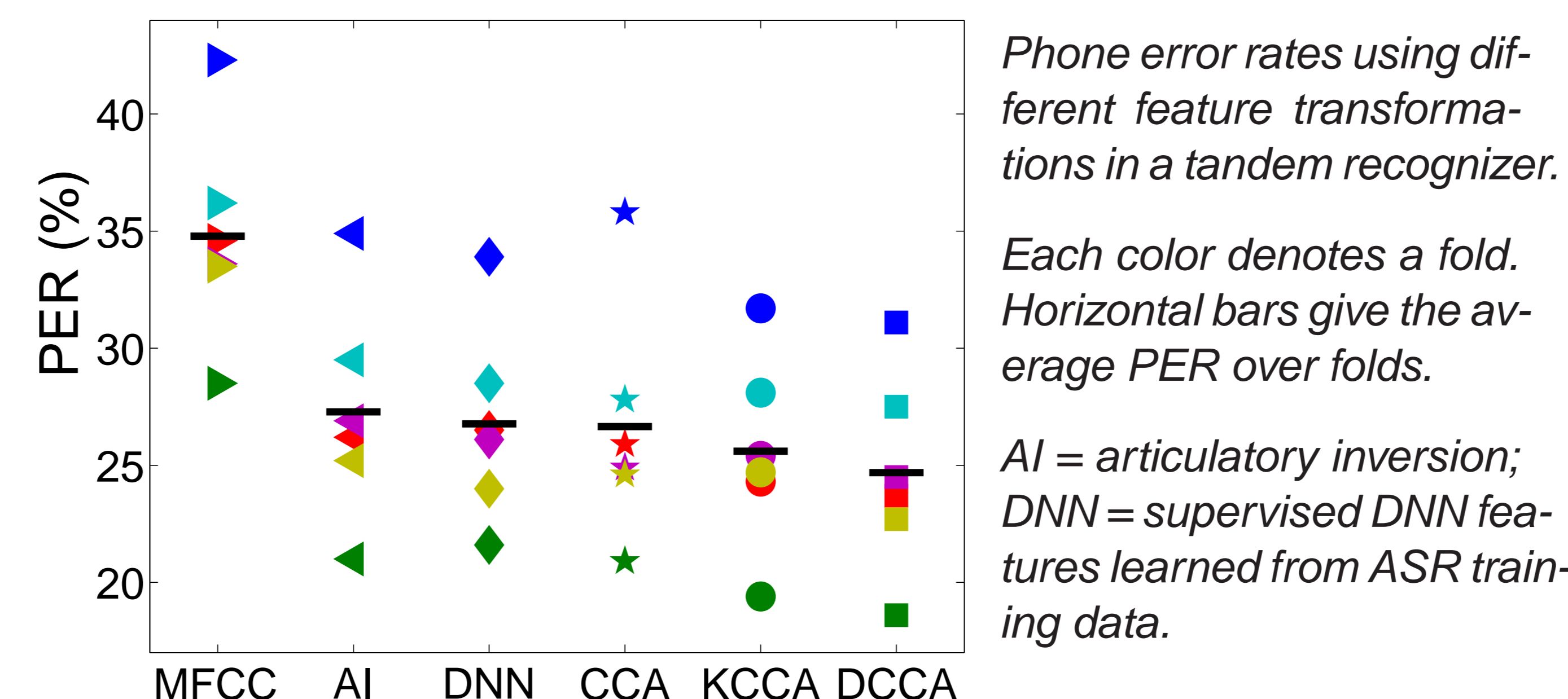


5 Experimental Results

Data: Wisconsin X-ray Microbeam Database, 47 speakers with ~ 50 utterances each, divided into 35/8/2/2 for feature learning/ASR train/tune/test. 6-fold cross-validation for ASR.

Acoustic input: 13 MFCCs + $\Delta + \Delta\Delta \times 7$ frames (273D).

Articulatory input: x, y displacements of 8 pellets $\times 7$ frames (112D) + per-speaker mean/variance normalization.

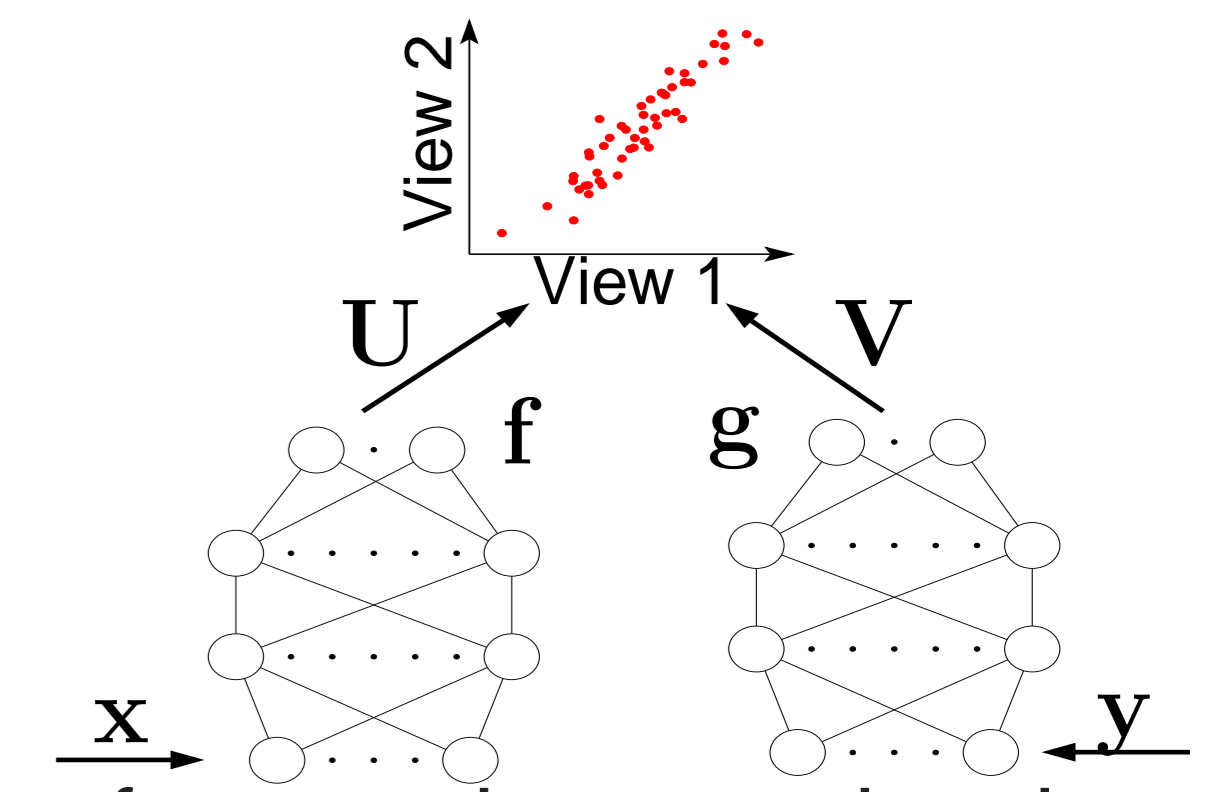


- DCCA is the best performer in all folds. All DCCA improvements over other feature types are significant at $p < 0.05$.
- Code is available from <http://ttic.edu/livescu>

3 Deep Canonical Correlation Analysis

- **Training data:** $\{(x_i, y_i)\}_{i=1}^N$ where $x_i \in \mathbb{R}^{D_x}$ and $y_i \in \mathbb{R}^{D_y}$ are input features for i^{th} frame. Here x = acoustics (View 1) and y = articulatory measurements (View 2).
- **Feature mappings:** $f : \mathbb{R}^{D_x} \rightarrow \mathbb{R}^{d_x}$ and $g : \mathbb{R}^{D_y} \rightarrow \mathbb{R}^{d_y}$, optionally parameterized by W_1, W_2 , for View 1 and View 2 respectively. Let $F = f(X) = [f(x_1), \dots, f(x_N)]$, $G = g(Y) = [g(y_1), \dots, g(y_N)]$.
- **CCA objective:** Find linear projections $U \in \mathbb{R}^{d_x \times L}$ and $V \in \mathbb{R}^{d_y \times L}$, and optionally parameters of f, g , with maximal canonical correlation:

$$\begin{aligned} & \max_{U, V, W_1, W_2} \text{tr}(U^T F G^T V) \\ \text{s.t. } & U^T \left(\frac{1}{N} F F^T + r_x I \right) U = I, \\ & V^T \left(\frac{1}{N} G G^T + r_y I \right) V = I. \end{aligned}$$



where r_x, r_y are regularization coefficients for covariance estimation.

CCA variants

- Linear CCA (CCA): $f(x) = x$ and $g(y) = y$.
- Kernel CCA (KCCA): f/g are feature maps induced by kernels k_x/k_y .
- Deep CCA (DCCA): f and g are outputs of DNNs.
- For fixed f, g , optimal (U, V) given via singular value decomposition.
- **Final features:** $\tilde{f}(x) = U^T f(x)$.

6 Conclusions

- DCCA significantly improves over previous multi-view methods, articulatory inversion, and supervised DNN features.
- Improvement over articulatory inversion suggests predicting details of articulation is not important or useful.
- Stochastic optimization allows DCCA to scale well to large data.
- Future work: Apply to hybrid ASR, new domains; incorporate supervision [4]; further analysis of stochastic training and network types

References

- [1] Bharadwaj, Arora, Livescu, and Hasegawa-Johnson. Multi-view acoustic feature learning using articulatory measurements. IWSML 2012.
- [2] Arora and Livescu. Multi-view CCA-based acoustic features for phonetic recognition across speakers and domains. ICASSP 2013.
- [3] Andrew, Arora, Livescu, and Bilmes. Deep canonical correlation analysis. ICML 2013.
- [4] Arora and Livescu. Multi-view learning with supervision for transformed bottleneck features. ICASSP 2014.