

# REGULARIZED SPARSE KERNEL SLOW FEATURE ANALYSIS



Bernstein Focus:  
Neurotechnology  
Berlin

Wendelin Böhmer<sup>1,2</sup>, Steffen Grünewälder<sup>1,3</sup>, Hannes Nickisch<sup>1,4</sup> and Klaus Obermayer<sup>1,2</sup>

(1) Neural Information Processing Group, Technische Universität Berlin, Germany

(2) Bernstein Focus: Neurotechnologie, Bernstein Center Berlin, Germany

(3) Centre for Computational Statistics and Machine Learning, University College London, United Kingdom

(4) Philips Research Laboratories, Hamburg, Germany



## ABSTRACT

This paper develops a kernelized slow feature analysis (SFA) algorithm. SFA is an unsupervised learning method to extract features which encode latent variables from time series. Generative relationships are usually complex, and current algorithms are either not powerful enough or tend to over-fit. We make use of the kernel trick in combination with sparsification to provide a powerful function class for large data sets. Sparsity is achieved by a novel matching pursuit approach that can be applied to other tasks as well. For small but complex data sets, however, the kernel SFA approach leads to over-fitting and numerical instabilities. To enforce a stable solution, we introduce regularization to the SFA objective. Feature extraction is demonstrated on a vowel classification task.

## OBJECTIVE

$$\min \sum_i E_t[y_i^2(x_t)]$$

$$s.t. \quad \forall i: E_t[y_i(x_t)] = 0$$

$$\forall i, j: E_t[y_i(x_t)y_j(x_t)] = \delta_{ij}$$

$y_i(x)$  can be from an arbitrary function class

Given an **infinite time series** and **sufficient function class**, SFA features span a **Fourier basis** in the **space of the latent variables** [I].

## SLOW FEATURE ANALYSIS

## REPRODUCING KERNEL HILBERTSPACES

$$\min \frac{1}{n-1} \text{tr}(A^T K D D^T K^T A)$$

$$s.t. \quad \frac{1}{n} A^T K \mathbf{1} = 0$$

$$\frac{1}{n} A^T K K^T A = I$$

$$y_i(x) = \sum_t A_{it} \kappa(x, x_t) - c_i$$

Only **small training-sets** feasible due to complexity. For any **finite training-set**, kernel SFA exhibits **over-fitting** and **numerical instability** when applied on a test-set [III].

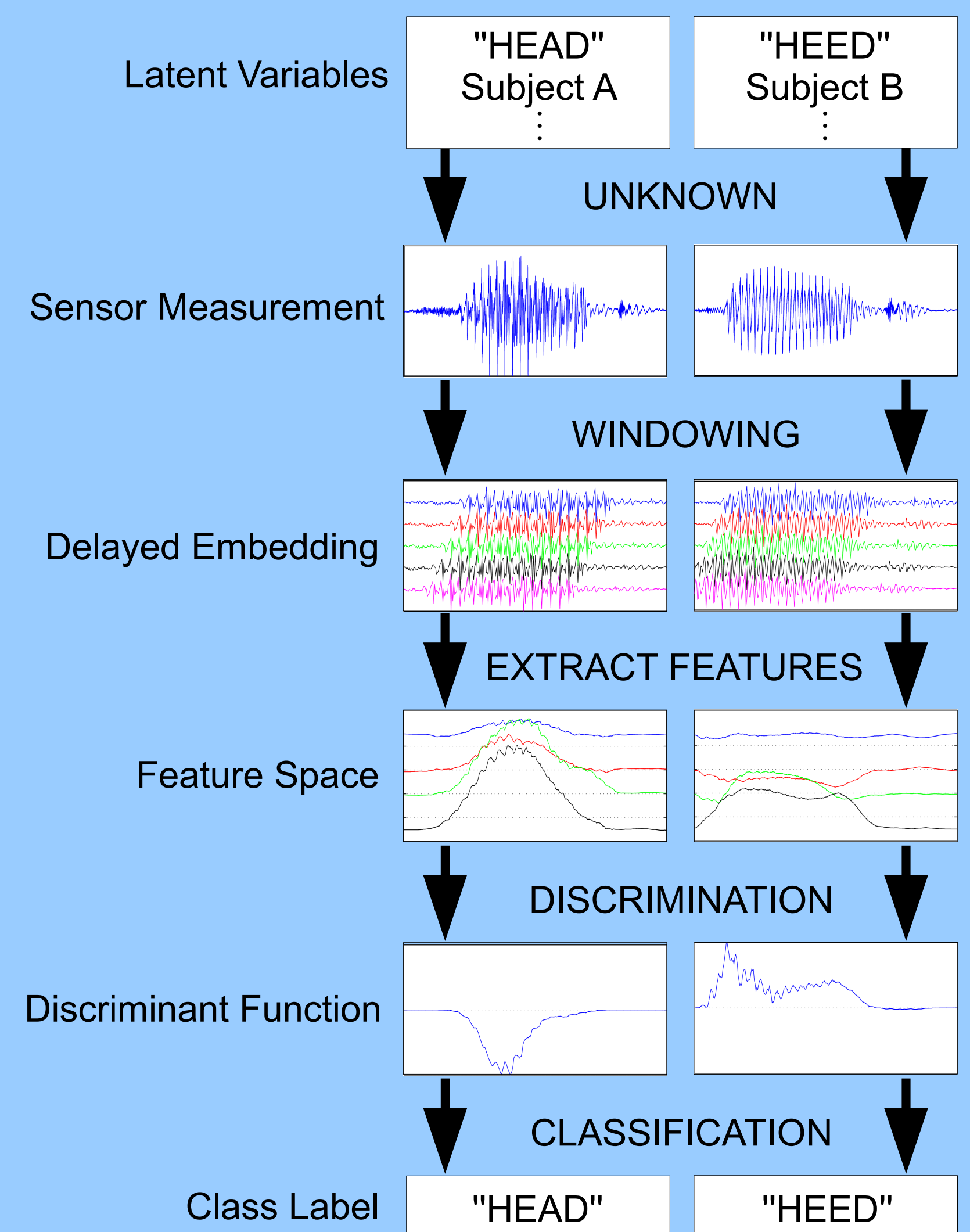
## FUNCTION CLASS

• Non-linear **PCA** does not encode latent variables

• Non-linear **SFA** does

## UNSUPERVISED NON-LINEAR FEATURE EXTRACTION

- **Classification** or **regression** w.r.t. **latent variables**
  - Latent variables **non-linearly embedded** in data
  - Discriminant/regression function **non-linear** in the space of latent variables
- Utilize knowledge from **unlabelled data**
  - **Unsupervised** construction of features from data
  - Functional basis on **manifold** of latent variables



## RSK-SFA

Slow  
Feature  
Analysis

## PENALIZE COMPLEX FUNCTIONS BY PENALIZING HILBERT NORM

$$\min \sum_i [E_t[y_i^2(x_t)] + \lambda \|y_i\|_H^2] \equiv$$

$$\min \frac{1}{n-1} \text{tr}(A^T K D D^T K^T A) + \lambda \text{tr}(A^T \bar{K} A)$$

$$\bar{K}_{ij} = \kappa(z_i, z_j)$$

Optimal constant  $\lambda$  can be **very small**. Adapts **better** to the objective. No computational overhead. No speed-up.

## PREVENT COMPLEX FUNCTIONS BY RESTRICTION TO SUBSPACE

$$y_i(x) \in \text{span}(\{\kappa(x, z_j)\}_{j=1}^m)$$

$$\{z_j\}_{j=1}^m \subset \{x_t\}_{t=1}^n, \quad m \ll n$$

$$K_{jt} = \kappa(z_j, x_t)$$

Removes the computational **bottle-neck**. Regularizes the solution **indirectly**. Depends strongly on **sparse subset selection**.

## SPARSENESS

## SPARSE SUBSET SELECTION

- $\{z_j\}_{j=1}^m$  should span the maximal possible subspace

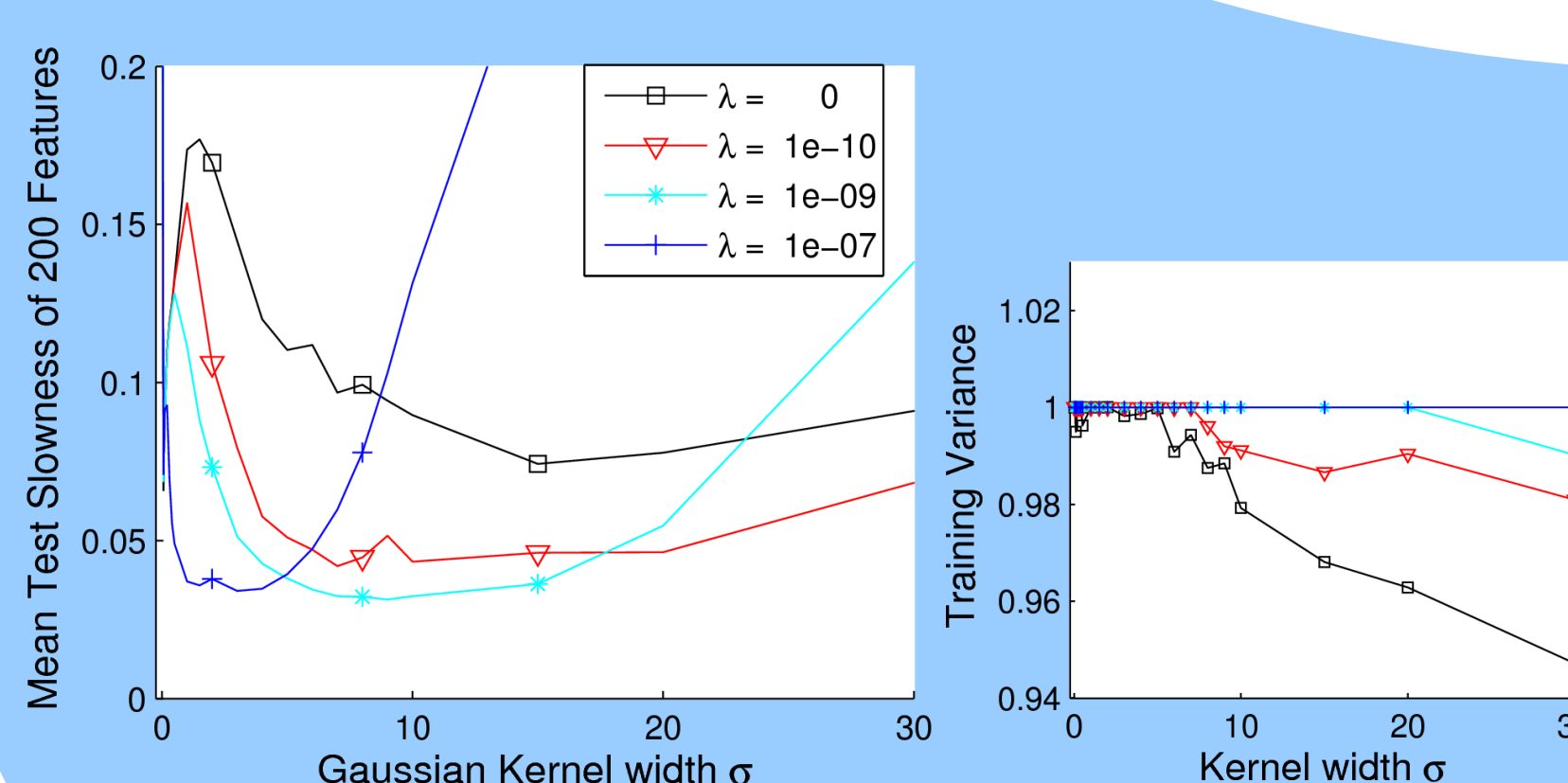
$$\epsilon_t^i := \min_{\alpha} \left\| \kappa(\cdot, x_t) - \sum_{j=1}^m \alpha_j \kappa(\cdot, x_{t_j}) \right\|^2$$

- Equivalent to the sparse kernel PCA problem
  - *Matching Pursuit for Sparse Kernel PCA* [V]  $O(n^2 m)$
  - *Online Maximization of the Affine Hull* [IV]  $O(n m^2)$
- *Matching Pursuit of Online MAH*  $O(n m^2)$

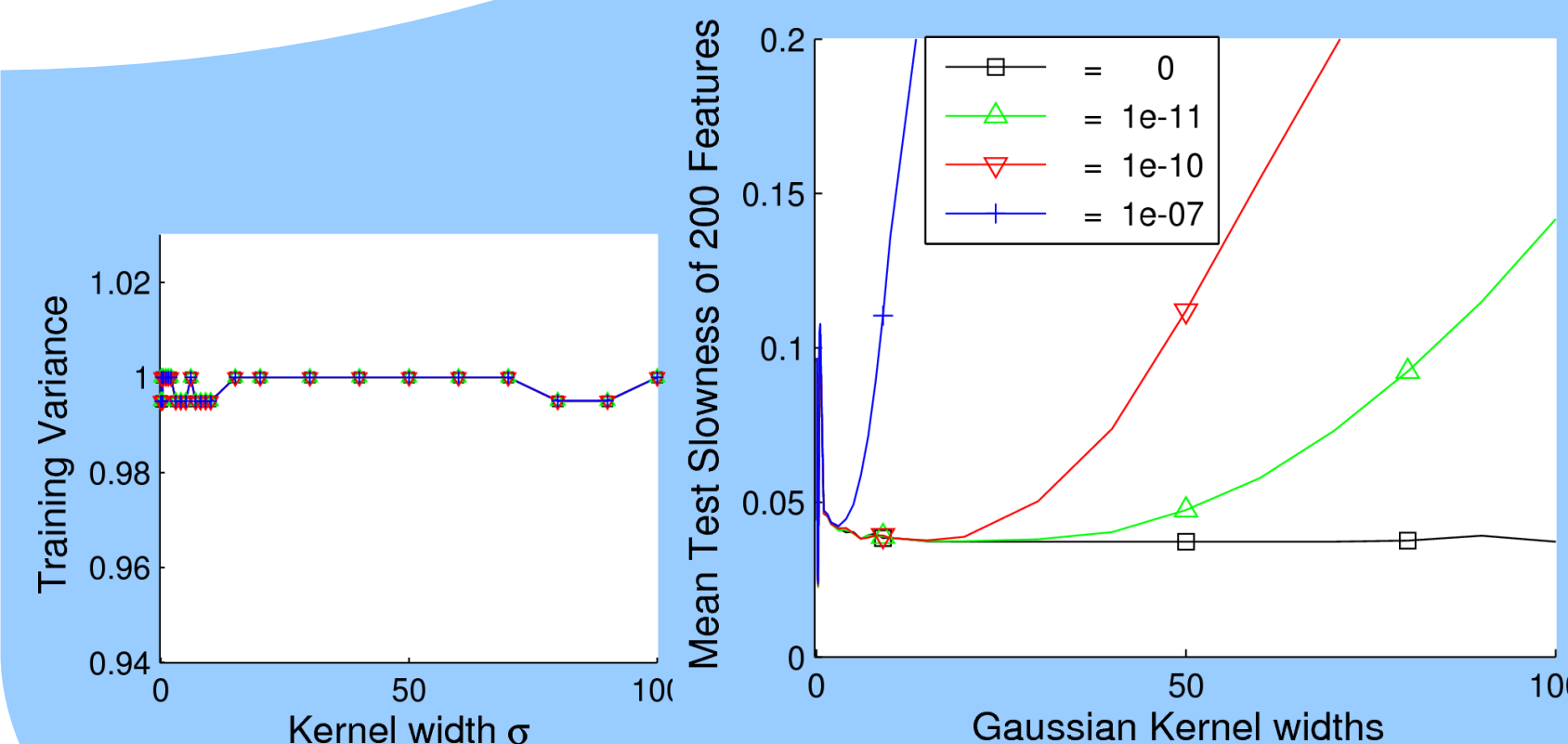
$$i_{j+1} := \text{argmin}_i \left\| [\epsilon_1^{i \cup j}, \dots, \epsilon_n^{i \cup j}] \right\|_{\infty} \approx \text{argmax}_i \epsilon_i^i$$

$$\epsilon_t^{i \cup j} = \epsilon_t^i - \frac{1}{\epsilon_j} (K_{tj} - K_{ti} (K_{ii})^{-1} K_{ij})^2$$

## RESULTS

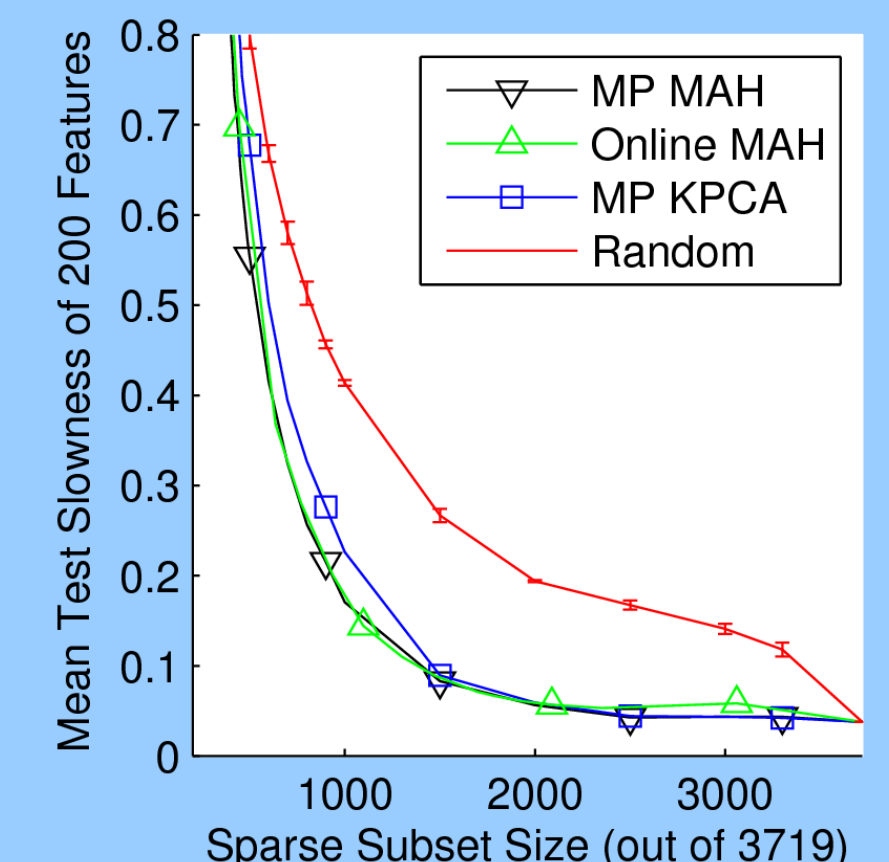


## RESULTS

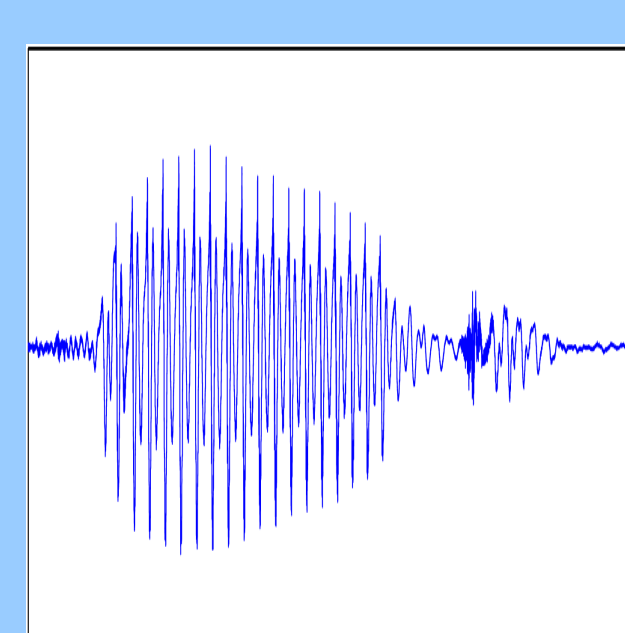


## RESULTS

- All algorithms outperform random selection
- MP MAH exhibits **same performance**
- MP MAH returns **ordered subset**
- **Regularization by sparseness** becomes much easier

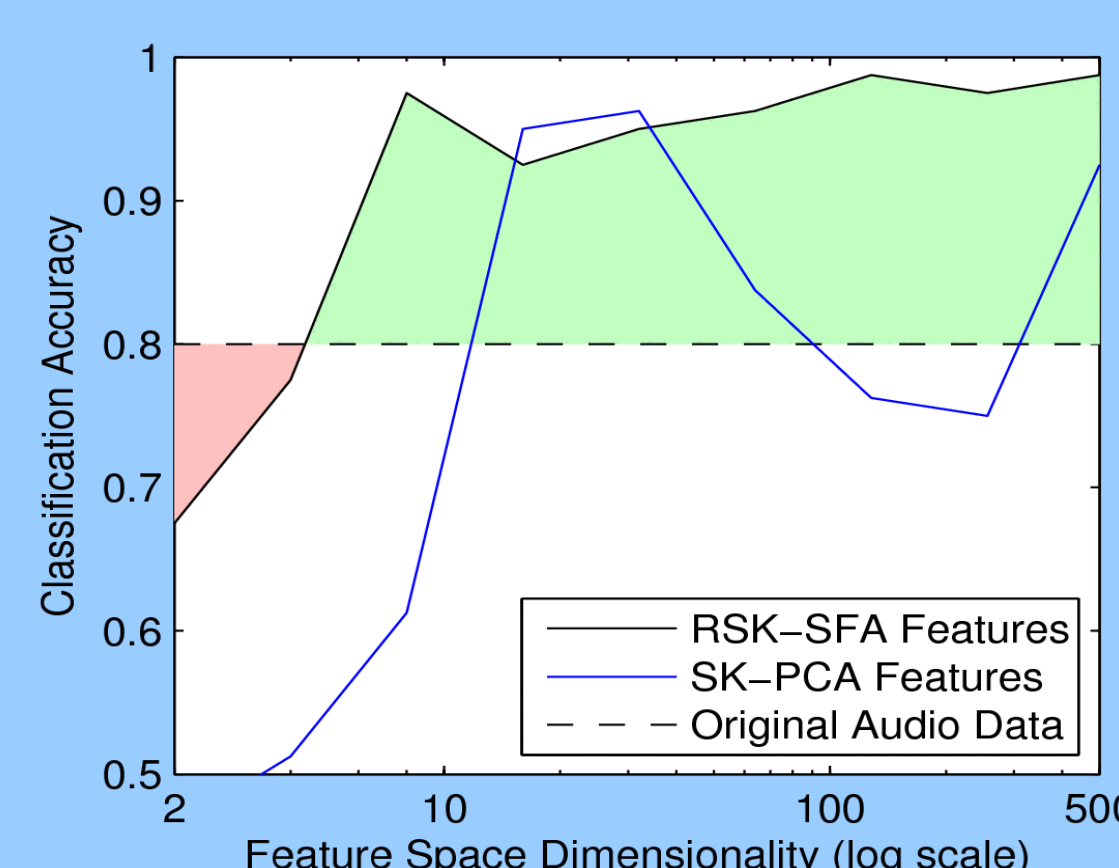
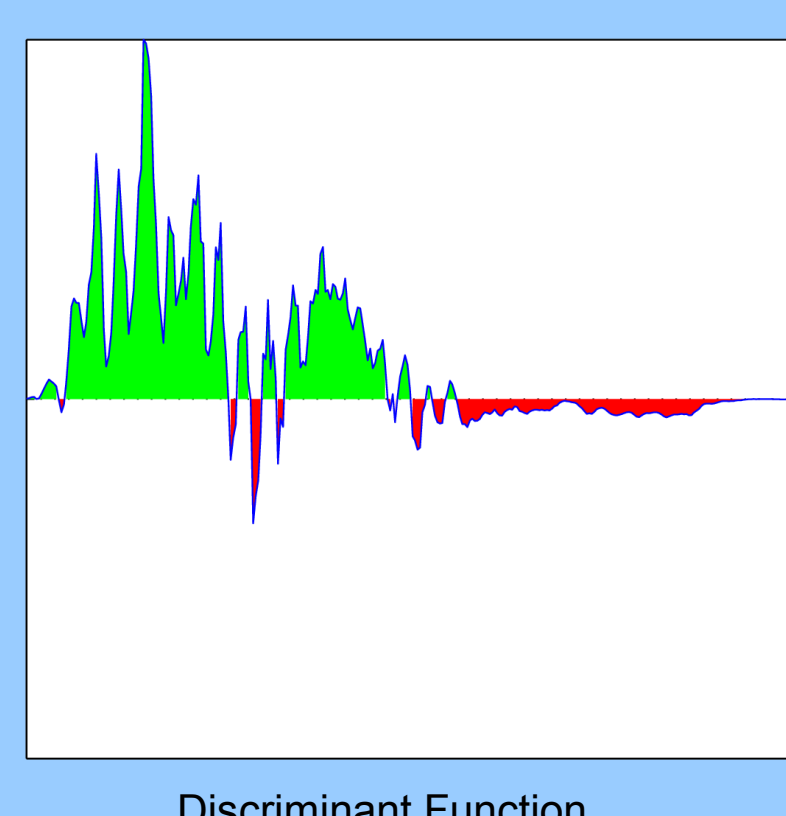


## FEATURE VALIDATION: VOWEL CLASSIFICATION



Delayed Embedding  
RSK-SFA Features

QDA  
(quadratic discriminant analysis)



## REFERENCES

- I. Wiskott: *Slow feature analysis: a theoretical analysis of optimal free responses*. Neural Computation 15(4):2147-2177, (2002)
- II. Bray and Martinez: *Kernel based extraction of slow features: complex cells learn disparity and translation invariance from natural images*. NIPS 15:253-260, (2002)
- III. Fukumizu, Bach and Gretton: *Statistical consistency of kernel canonical correlation analysis*. Journal of Machine Learning Research 8:361-383, (2007)
- IV. Csato and Oppor: *Sparse on-line gaussian processes*. Neural Computation 14(3):641-668, (2002)
- V. Hussain and Shawe-Taylor: *Theory of matching pursuit*. NIPS 21:721-728, (2008)

contact: Wendelin Böhmer <wendelin@cs.tu-berlin.de>

This work has been supported by the IGP on Human-Centric Communication at TU Berlin and the German Federal Ministry of Education and Research (grant 01GQ0850)