

# **Mining MR Image Data by Discriminative Methods for the Diagnosis of Dementia**

**Ceyhun Burak Akgül**

Former Marie Curie Postdoctoral Fellow @ Philips Research

[www.cba-research.com](http://www.cba-research.com)

[cb.akgul@gmail.com](mailto:cb.akgul@gmail.com)

# Motivation – 1/2

Diagnose dementia (e.g., Alzheimer's disease) from MR Images

Standard medical practice:

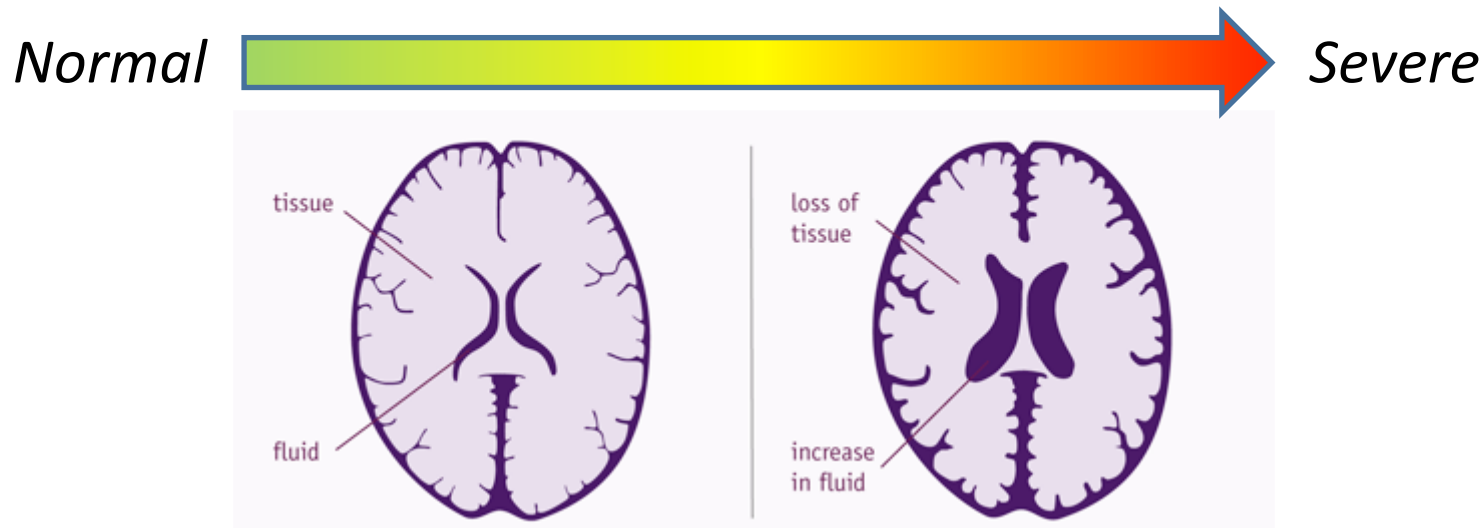
- ❑ patient history, collateral history from relatives
- ❑ clinical observations: neurological/neuropsychological features

**BUT:** does not often lead to an early diagnosis

**An emerging trend: Exploit imaging data  
HOW?**

# Motivation – 2/2

## Brain Atrophy?



- requires longitudinal data: MR scans at different time stamps
- requires complex mathematical modeling and algorithms
- should quantify minute changes (that human eye can't see)

***Or something else...***

# Data Mining Framework

## Representation

learn an image representation from data: analyze images

- at each location
- at several scales
- with several patterns

## Selection

discover image features using labeled data

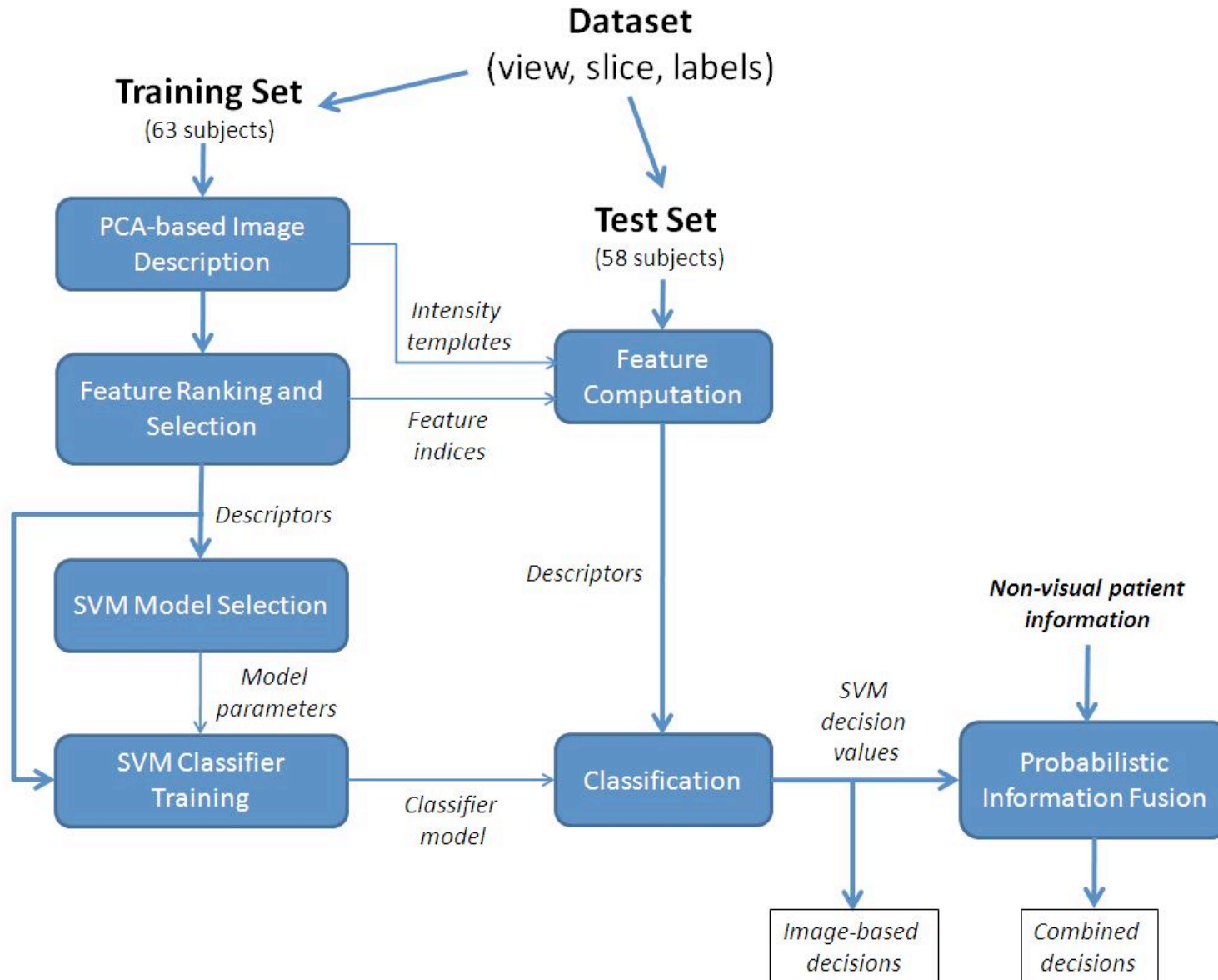
## Classification

characterize patient groups discriminatively

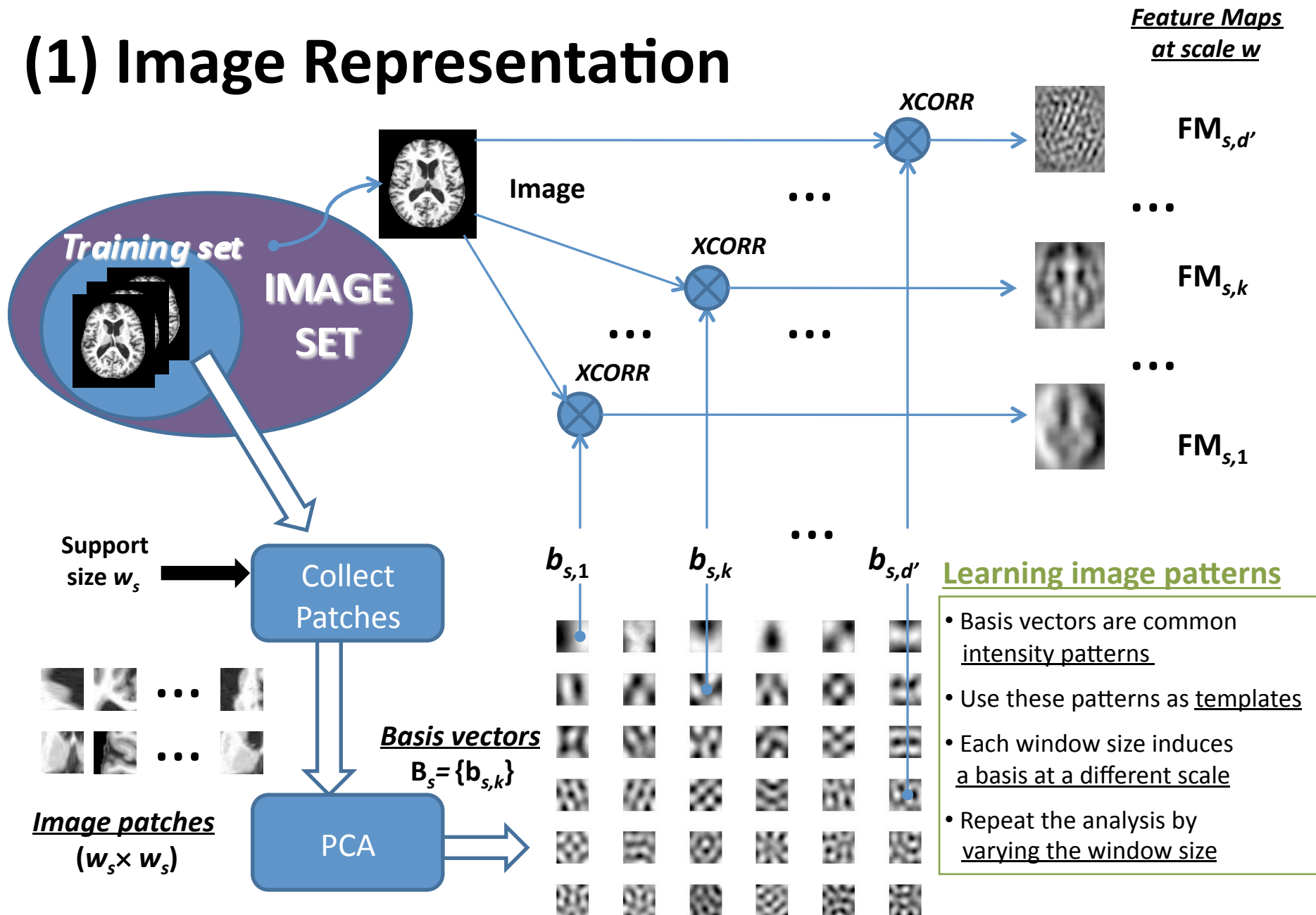
## Information Fusion

combine multiple (visual or non-visual) information sources

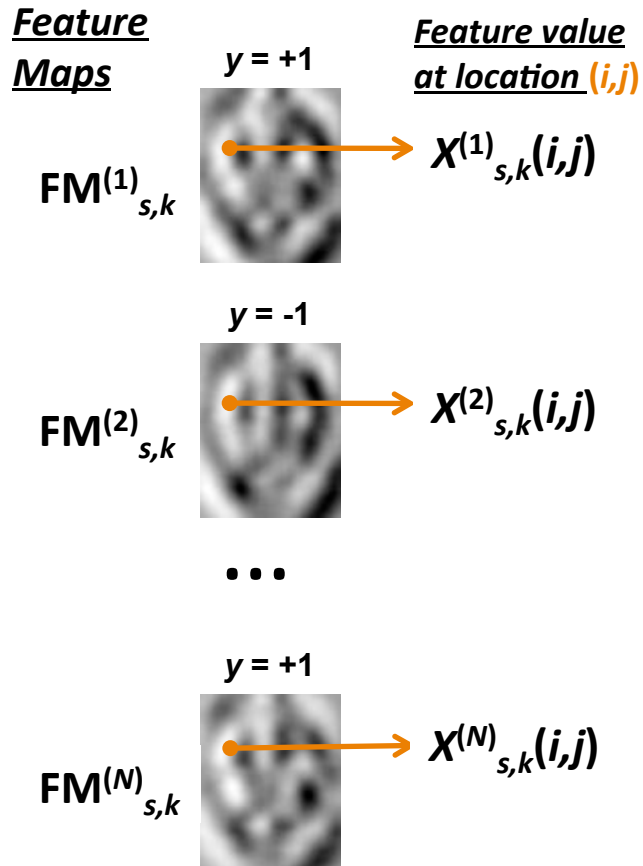
# Data Mining Framework: Overview



# (1) Image Representation



## (2) Feature Selection by Ranking – 1/3



Scale IDs:  $s=1,\dots,S$   
 Basis (Pattern) IDs:  $k=1,\dots,K_s$   
 Subject IDs:  $n=1,\dots,N$

- Each image is described by  $S \times K_s$  feature maps
- **At each pixel location, there are  $S \times K_s$  feature values**
- Each feature  $\leftrightarrow$  a distinct (scale, template)-pair
- At each location:
  - rank the features based on their “**usefulness**”
  - pick **the most “useful”** feature for description

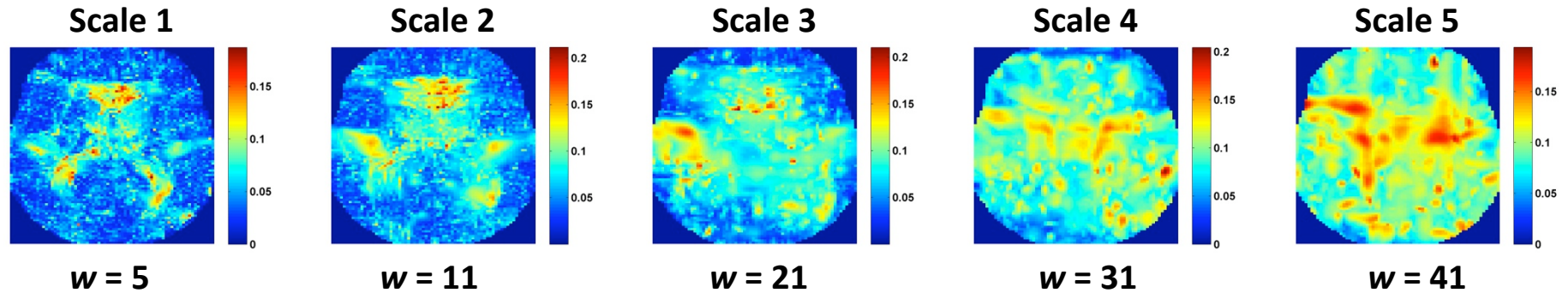
“Usefulness”  
 $\leftrightarrow$   
 Mutual information  
 between feature and diagnostic label

$$MI(x, y) = \sum_{y \in \{-1, +1\}} \int p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

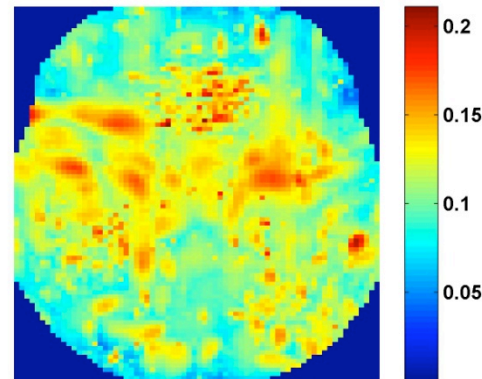
$x \in [0,1]$ : normalized feature value at location  $(i, j)$   
 $y \in \{-1, +1\}$ : diagnostic label of the image

## (2) Feature Selection by Ranking – 2/3

Maximum Mutual Information Maps<sup>(\*)</sup> at Different Scales



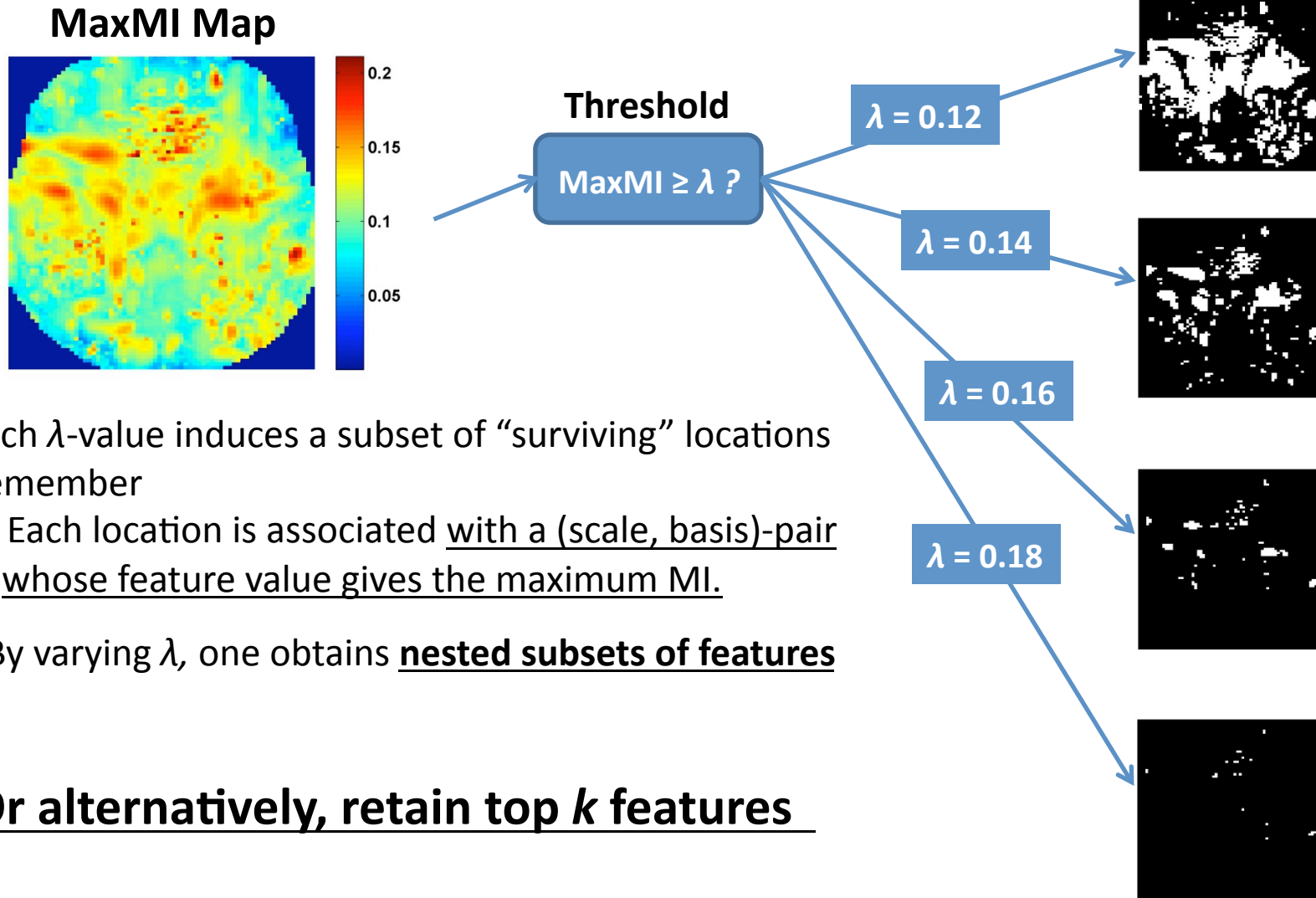
Maximum Mutual Information Map<sup>(\*)</sup>  
Combined over Scales



(\*) Map size:  $87 \times 70$



## (2) Feature Selection by Ranking – 3/3

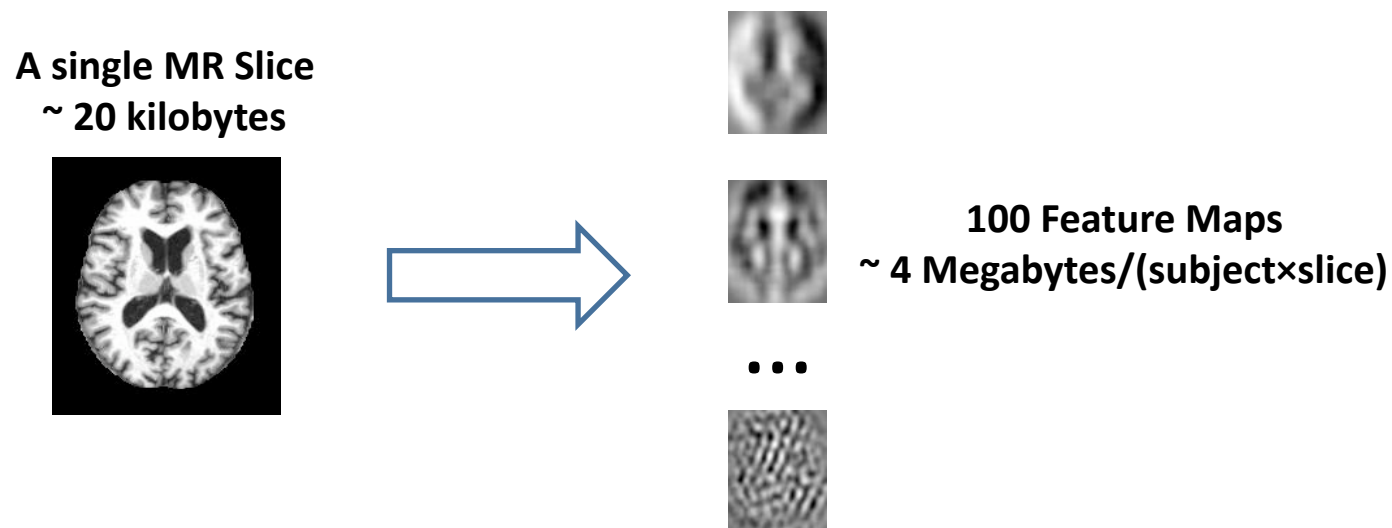


- Each  $\lambda$ -value induces a subset of “surviving” locations
- Remember
  - Each location is associated with a (scale, basis)-pair whose feature value gives the maximum MI.

→ By varying  $\lambda$ , one obtains nested subsets of features

- Or alternatively, retain top  $k$  features

# Amount of Data Processed: Some Facts



- 100 slices per subject ~ **400 Megabytes/subject**
- 121 subjects ~ **50 Gigabytes** TOTAL AMOUNT OF DATA PROCESSED
- 100 informative features/(subject×slice) selected as the descriptor  
< 1 kilobyte/(subject×slice)

## (3) Classification: SVM Basics

A non-linear SVM classifier  $F$  is indexed by two parameters  $(C, \gamma)$ :

- The parameter  $C$  trades off training error vs. classifier complexity
- The kernel parameter  $\gamma$  determines the class of functions  $F$  and affects class separation  
(in some sense, it also determines the classifier complexity)

**One has to specify the “best”  $(C, \gamma)$ -pair before testing the classifier.**

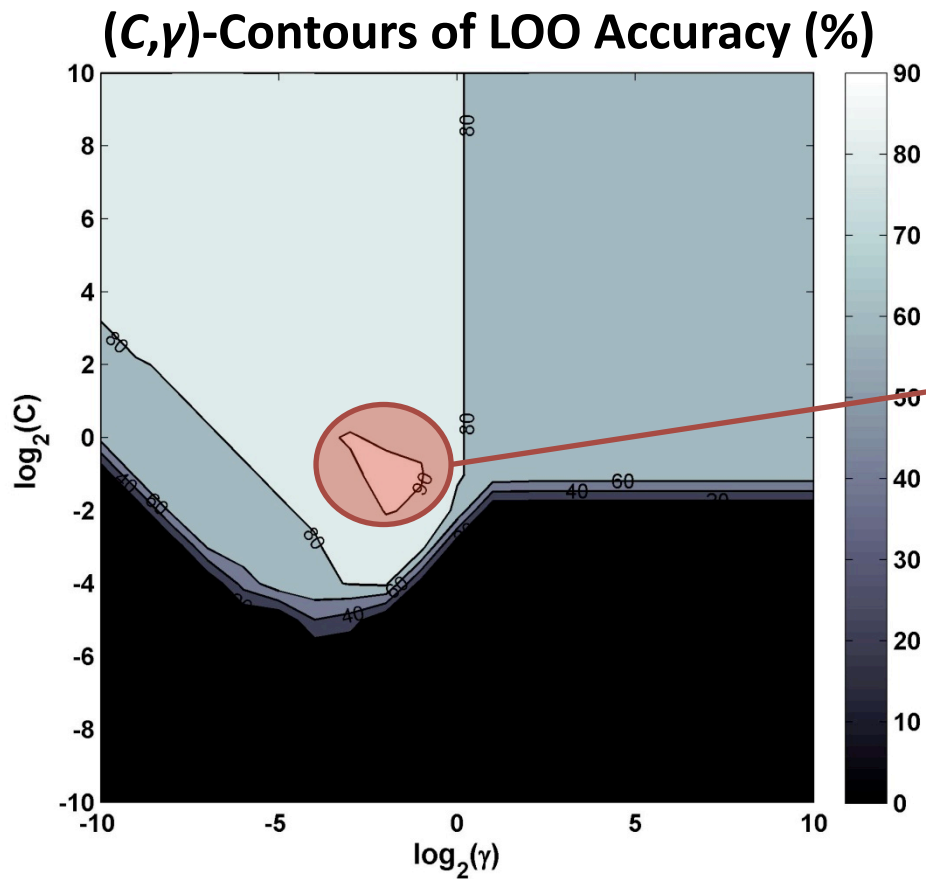
**A good empirical option**

$$(C, \gamma)^* = \operatorname{argmin} \operatorname{Err}_{CV}(F(C, \gamma))$$

*Err<sub>CV</sub>: Cross validation error*

# (3) Classification: Model Selection

- Leave-One-Out (LOO) cross-validation
- Initial search for the  $(C, \gamma)$ -parameters on a coarse grid



**$\{(C, \gamma) \text{ such that } \text{ACC} > \text{THRESH}\}$**

- Search on a finer grid
- Further heuristics – Look at:
  - Sensitivity
  - Fraction of SVs (model parsimony)
  - Specificity

## (4) Probabilistic Information Fusion

**Bayesian Theory:** The decision on the class label should be made on the conditional probability of the class label given all other relevant information.

Cognitive test scores, e.g., MMSE  
Age  
Gender  
Genetics ...

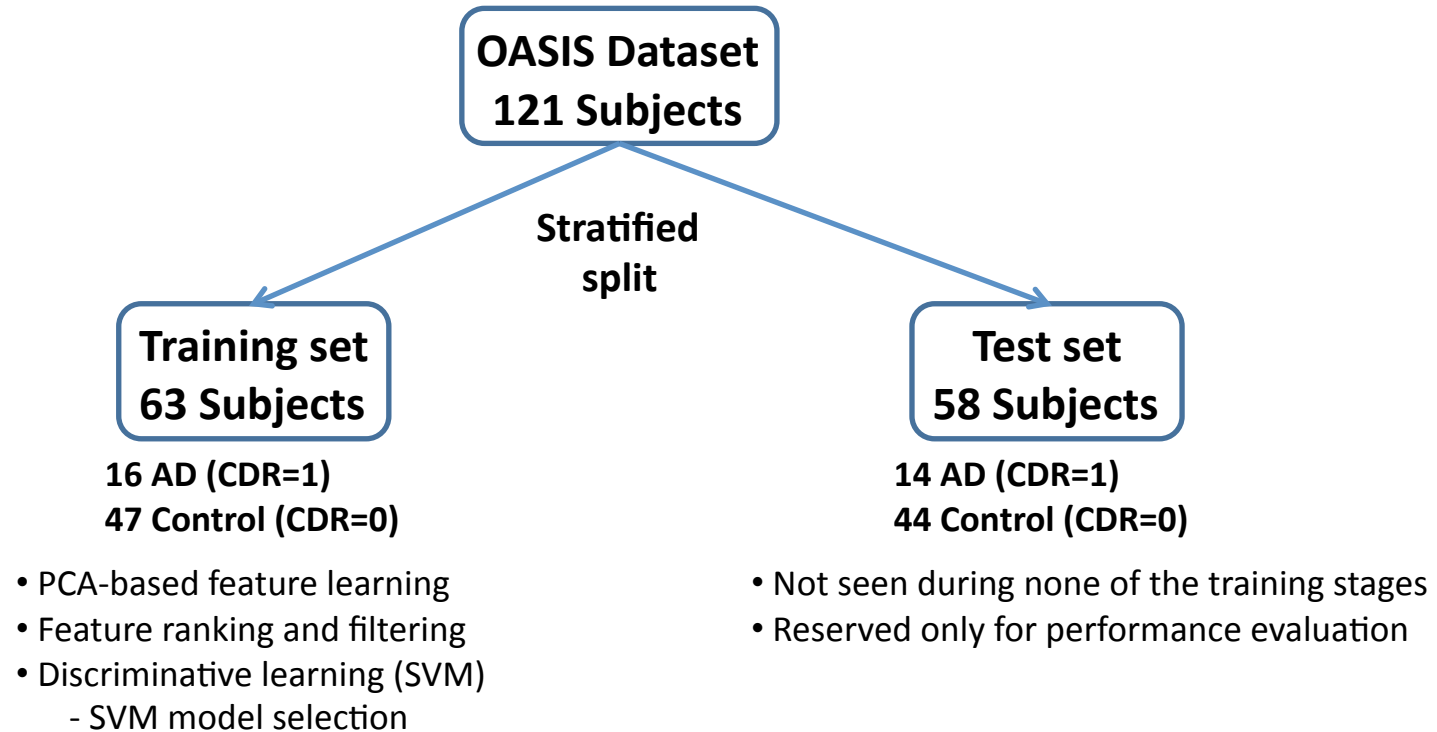
$$P(\text{label} \mid \text{info}) = P(\text{label} \mid \text{visual}, \text{non-visual})$$

$$\begin{aligned} P(\text{label} \mid \text{info}) &\propto P(\text{label}) \times P(\text{visual}, \text{non-visual} \mid \text{label}) \\ &= P(\text{label}) \times P(\text{visual} \mid \text{label}) \times P(\text{non-visual} \mid \text{label}) \\ &\propto \underbrace{P(\text{label} \mid \text{visual})}_{\text{derived from SVM outputs}} \times \underbrace{P(\text{non-visual} \mid \text{label})}_{\text{class-conditional distributions}} \end{aligned}$$

derived from SVM outputs

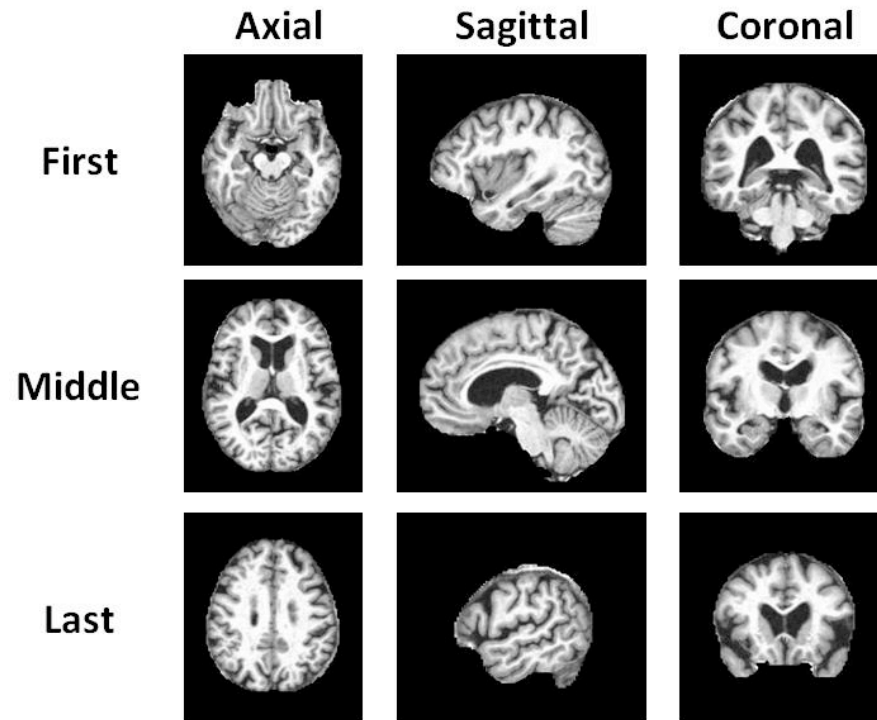
class-conditional distributions  
estimated from training data

# Experiments: Dataset



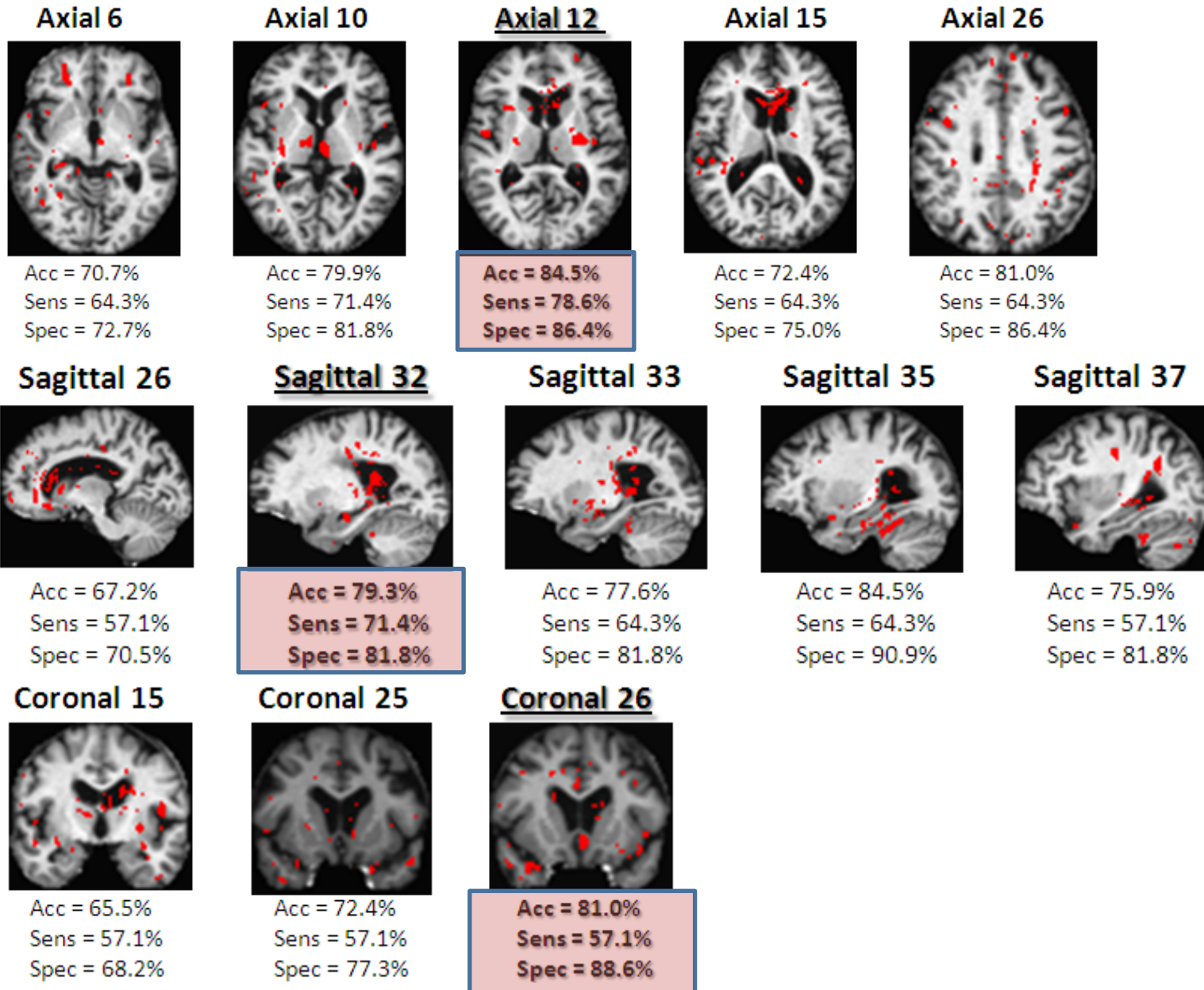
- CDR: Clinical Dementia Rating: normal → CDR = 0 moderate dementia → CDR = 1
- Stratified split keeps the class proportions the same in both sets (Control/AD ≈ 3)

# Experiments: MR Data



- **26 Axial + 46 Sagittal + 28 Coronal = 100 MR slices processed separately**
- **Each slice described by 100 informative image features**

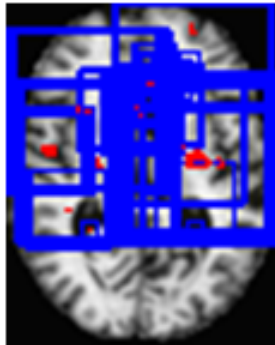
# Experiments: Discriminative Slices – 1/2





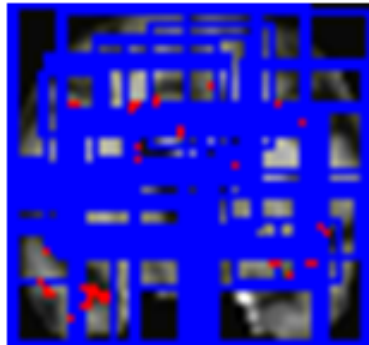
# Experiments: Discriminative Slices – 2/2

Axial 12



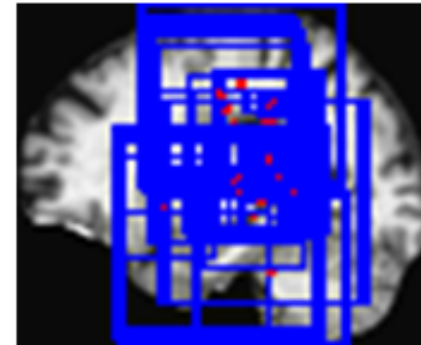
Acc = 84.5%  
Sens = 78.6%  
Spec = 86.4%

Coronal 26



Acc = 81.0%  
Sens = 57.1%  
Spec = 88.6%

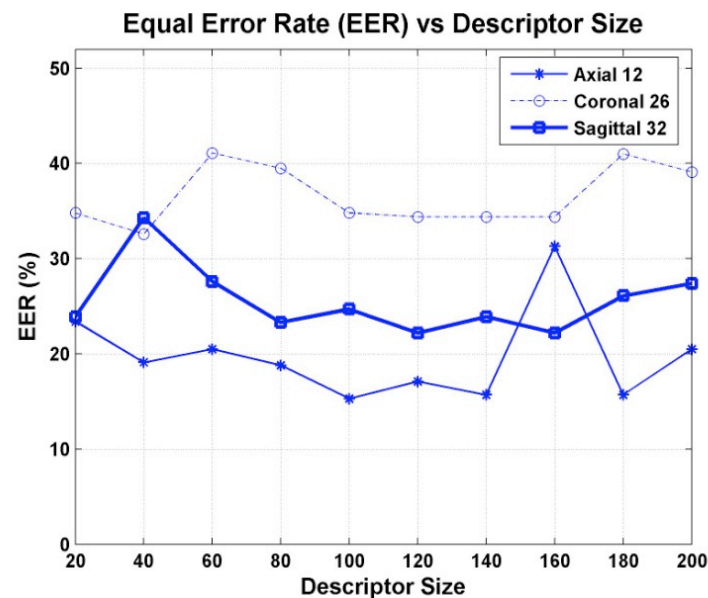
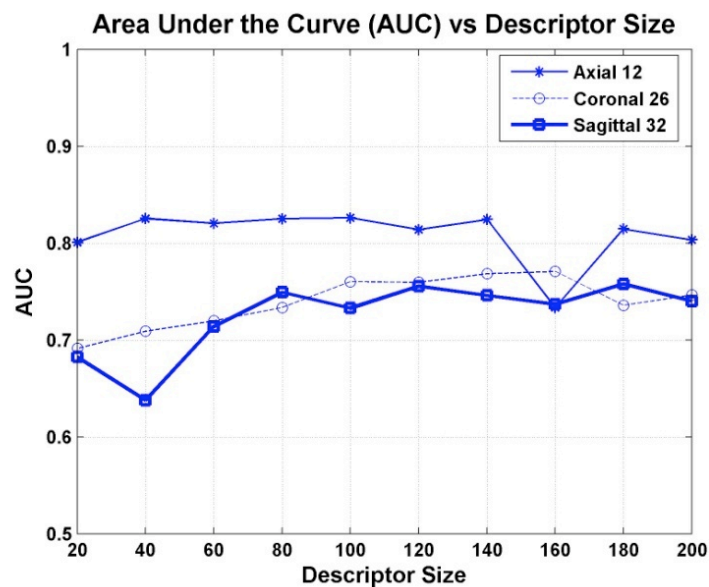
Sagittal 32



Acc = 79.3%  
Sens = 71.4%  
Spec = 81.8%

Axial 12 > Sagittal 32 > Coronal 26

# Experiments: ROC vs. Descriptor Size



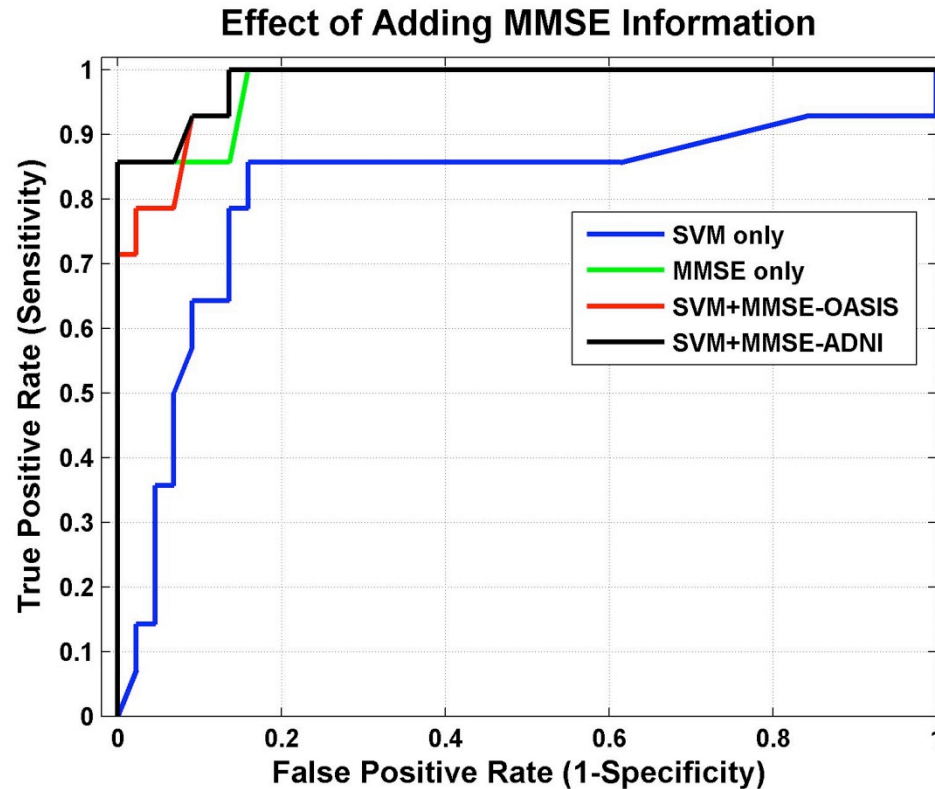
	Descriptor Size	Accuracy (%)	Sensitivity (%)	Specificity (%)
<b>Axial 12</b>	<b>100</b>	<b>84.5 (49/58)</b>	<b>85.7 (12/14)</b>	<b>84.1 (37/44)</b>
<b>Sagittal 32</b>	160	65.5 (38/58)	64.3 (9/14)	65.9 (29/44)
<b>Coronal 25</b>	160	77.6 (45/58)	78.6 (11/14)	77.3 (34/44)

**ROC:** Receiver Operating Characteristic: TPR vs. FPR

**AUC:** Area under the ROC curve

**EER:** Equal error rate (sensitivity = specificity)

# Experiments: Information Fusion – 1/2



**SVM-only:** Image-based decisions gleaned from SVM outputs

**MMSE-only:** MMSE-based decisions: *if  $MMSE < Thresh$ , then decide ill*

**SVM+MMSE-OASIS:** statistics estimated from OASIS training set (63 subjects)

**SVM+MMSE-ADNI:** statistics estimated from ADNI dataset (322 subjects)

# Experiments: Information Fusion – 2/2

## ROC Summary

	AUC	EER (%)	Accuracy (%)
SVM only	0.8260	15.3	84.5
MMSE only	0.9798	13.3	86.7
SVM+MMSE-OASIS	0.9798	8.7	91.3
SVM+MMSE-ADNI	0.9871	8.4	91.6

**SVM+MMSE-ADNI > SVM+MMSE-OASIS > MMSE-only > SVM-only**

- Information fusion is very useful indeed
- Reliable statistics!!! ADNI (322 subjects) > OASIS (63 subjects)

**229 controls**

**93 positives**

**47 controls**

**16 positives**

# Summary

- Data-driven image representation
  - Unsupervised learning of local image patterns via PCA
  - Localized, at several scales, with several patterns
- Feature ranking and filtering
  - Supervised: based on MI between scalar features and class labels
- Discriminative learning
  - SVM model selection via cross-validation and further heuristics
- Information fusion
  - Leverage image-only decisions by non-visual information
  - Generic: works with any kind of meta-data as long as statistics available

## Proof of concept:

A promising data-driven framework for the diagnosis of dementia with high predictive performance

# What's Next?

## Practical

Go validate these results clinically

Do these slices, locations, scales, patterns make sense?

Acquire larger sets of labeled data

Allocate higher computational resources

## Methodological

Other sparser image representations: ICA-based? NNMF-based?

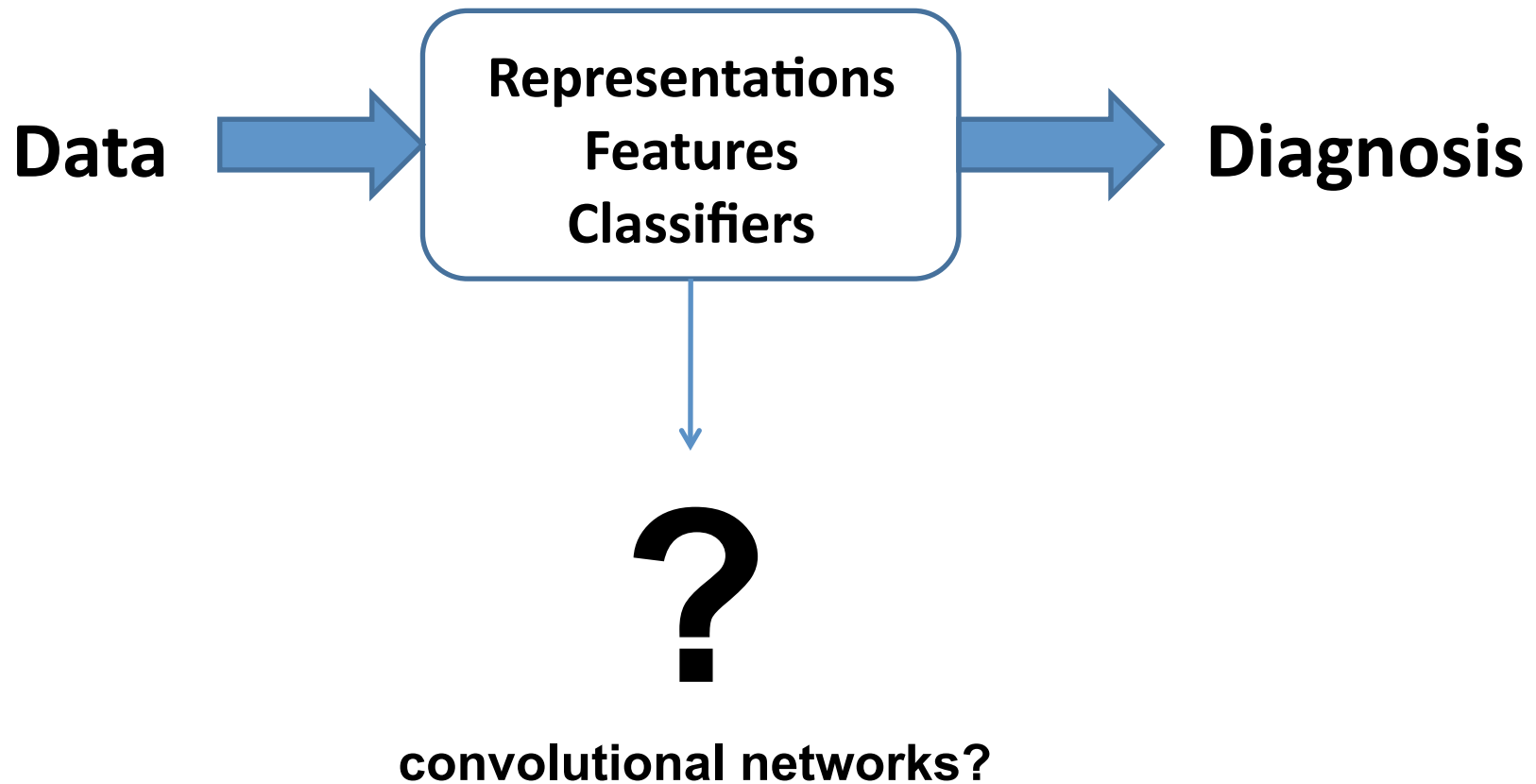
Multivariate feature selection

Model selection: Don't use one, average multiple models

Other classification schemes: AdaBoost

*Theoretical ...*

# What's Next? – *Theoretical*



***To conclude...***

There's nothing more practical than a good theory.

Lewin, 1952