

# Novel Fusion Methods for Pattern Recognition

Muhammad Awais,  
Fei Yan, Krystian Mikolajczyk and Josef Kittler  
ECML-PKDD 2011

# Motivation

## ➤ Automatic analysis of visual information



Crime prevention



Visually impaired



multimedia documents



Image content search

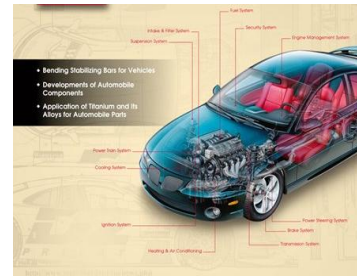
## ➤ Classification of photos



news agencies



art catalogues

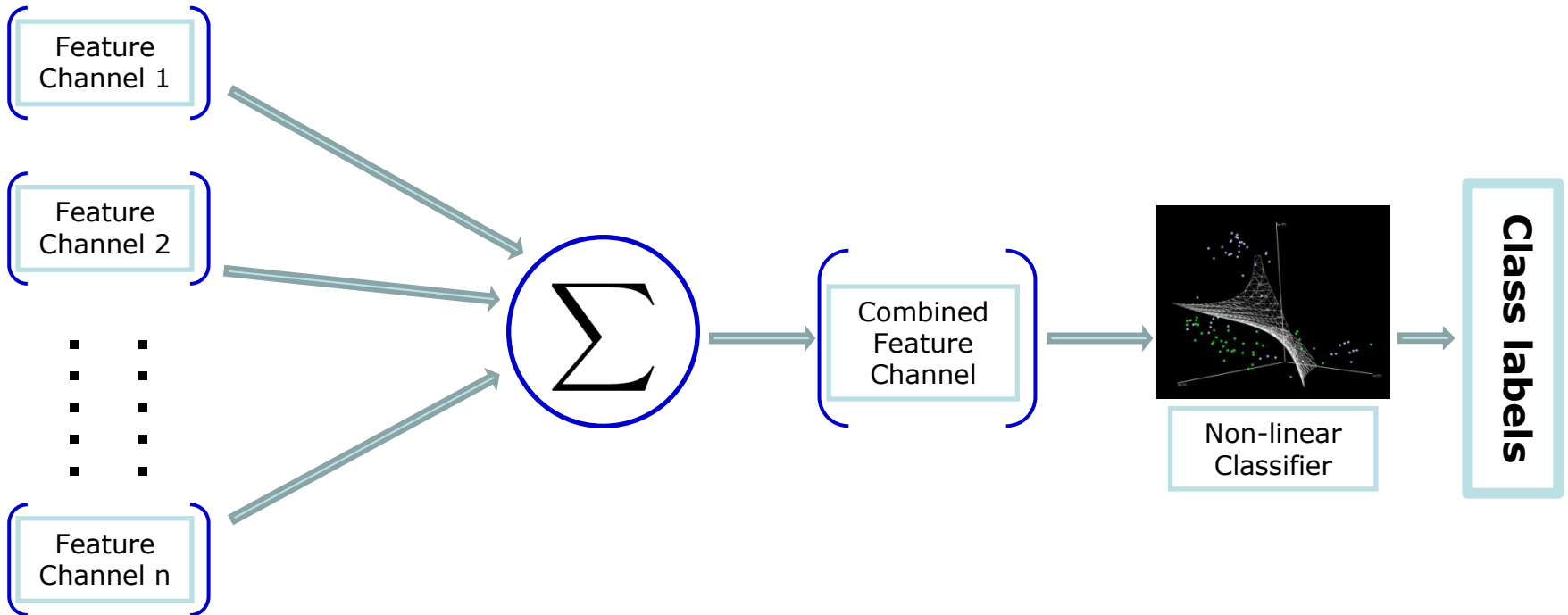


industrial components



trademarks

# Problem Statement

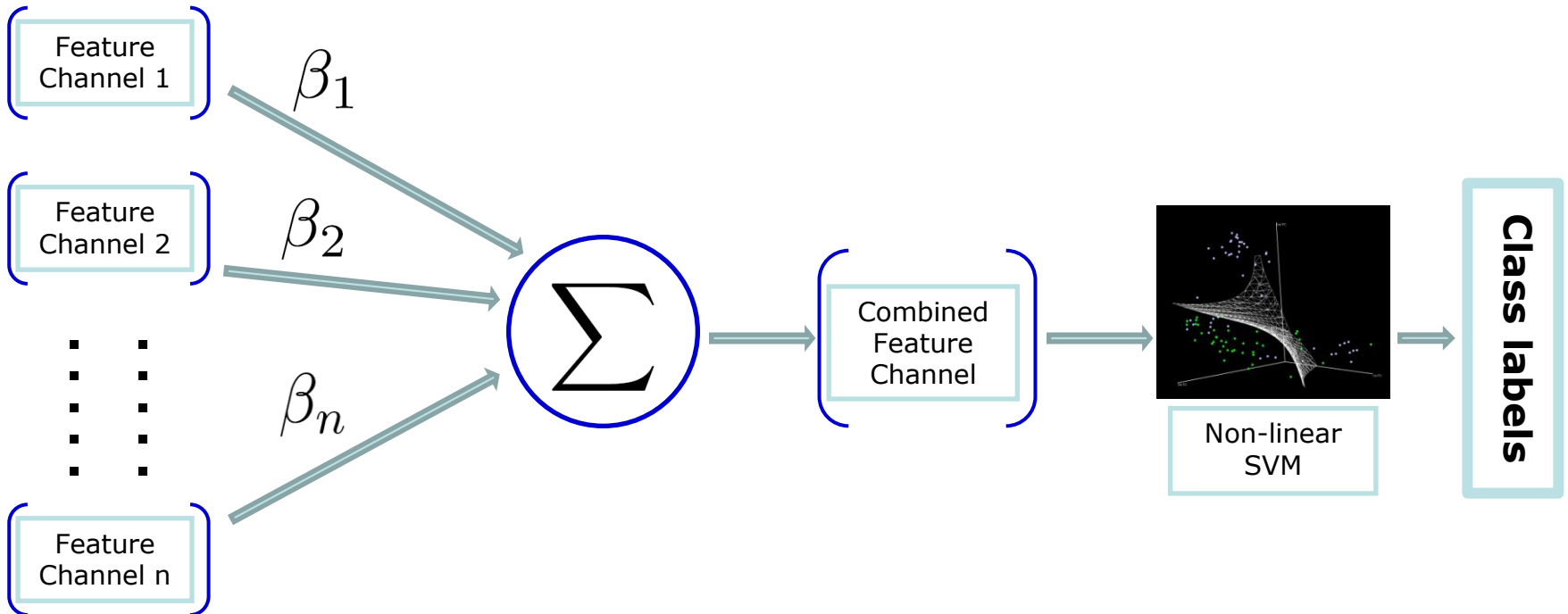


- Given a set of  $n$  features channels (kernels). Aim is to find the optimal way of combining these features channels.

# Contents

- Existing fusion techniques.
- Proposed Classifier fusion techniques.
  - Binary CLF.
  - NLP- $v$ MC.
  - NLP- $\beta$ .
  - NLP-B.
- Extended stacking.
- Evaluation on challenging datasets.
- Conclusions.

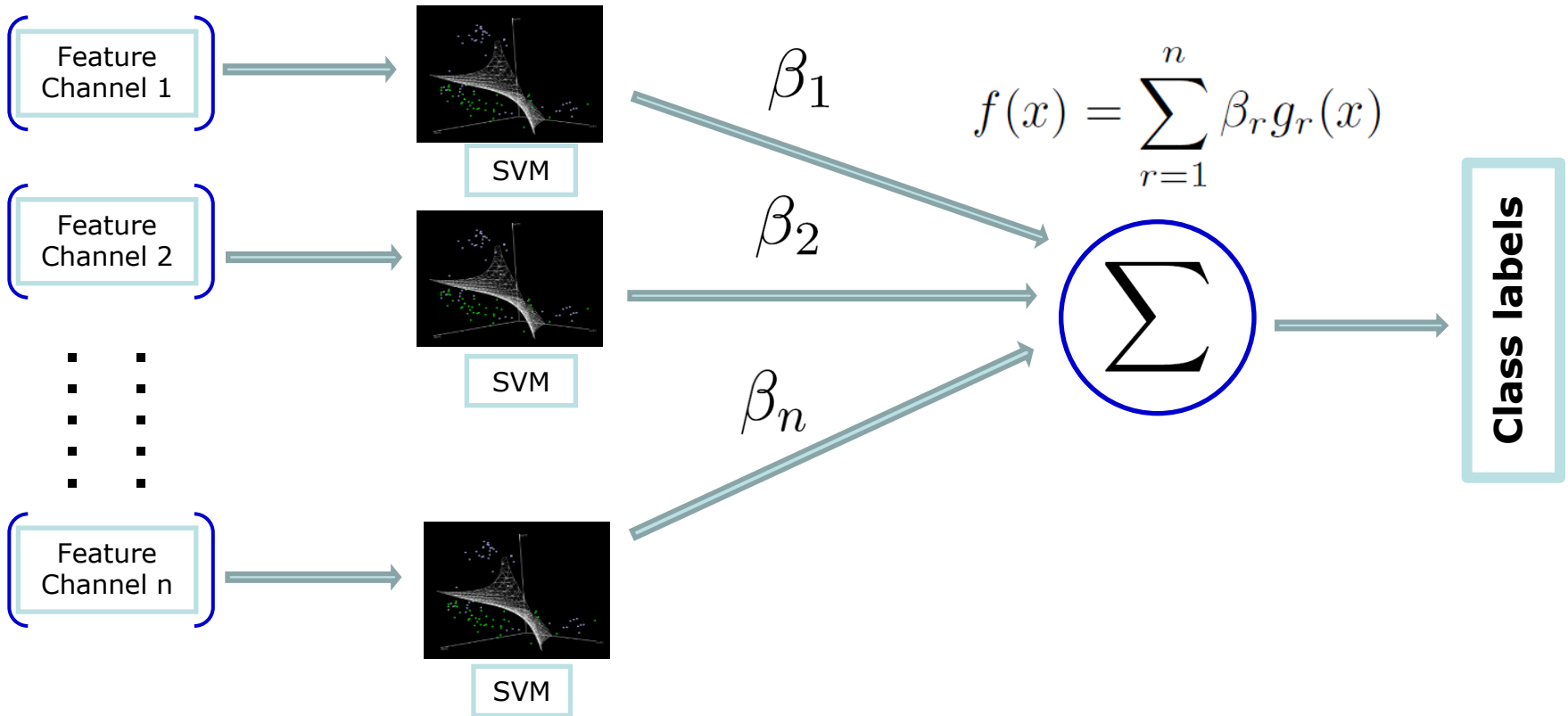
# Multiple Kernel Learning (MKL)



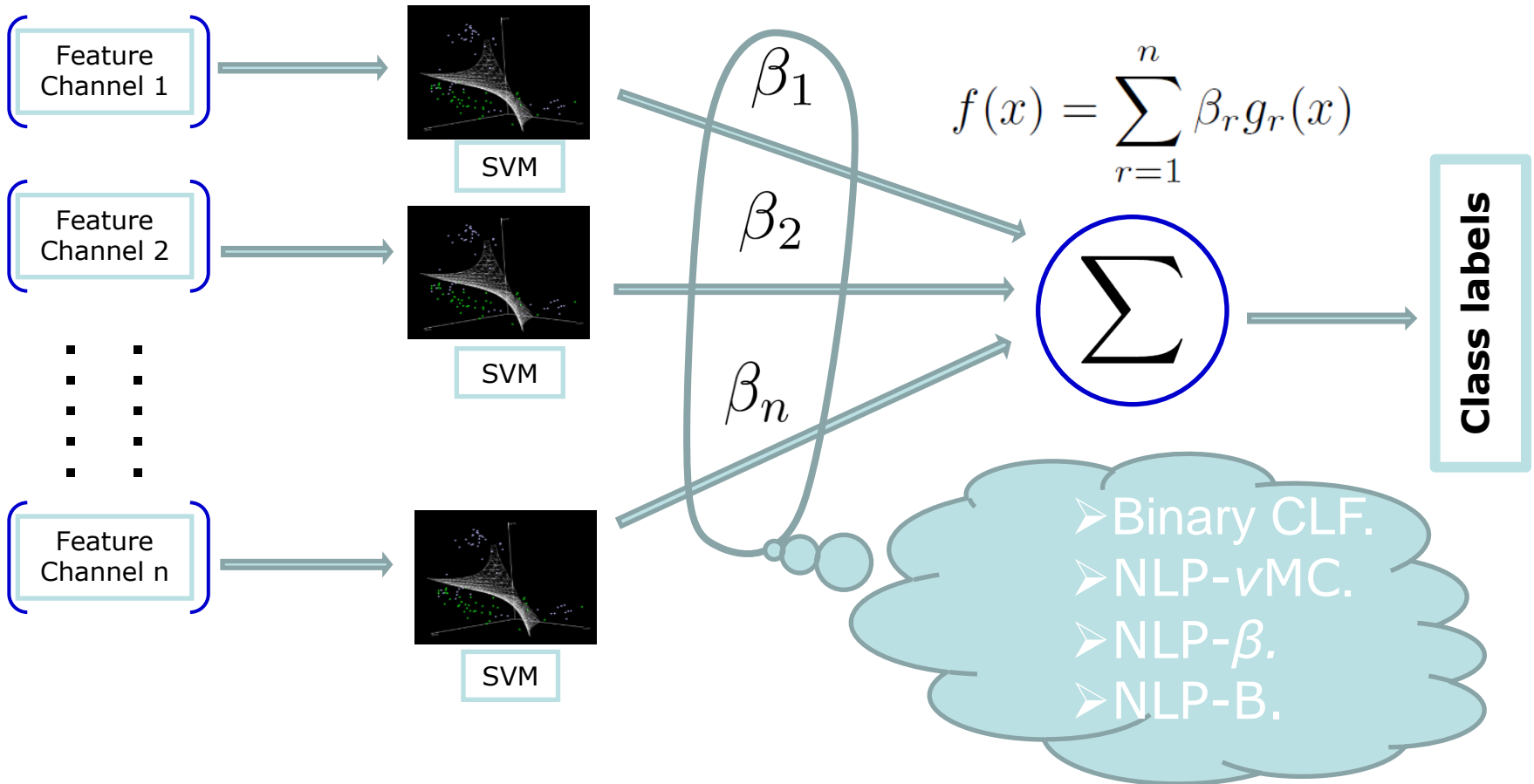
- MKL maximizes the soft margin to obtain optimal weights, for the convex combination of base kernels.

$$K = \sum_{p=1}^n \beta_p K_p$$

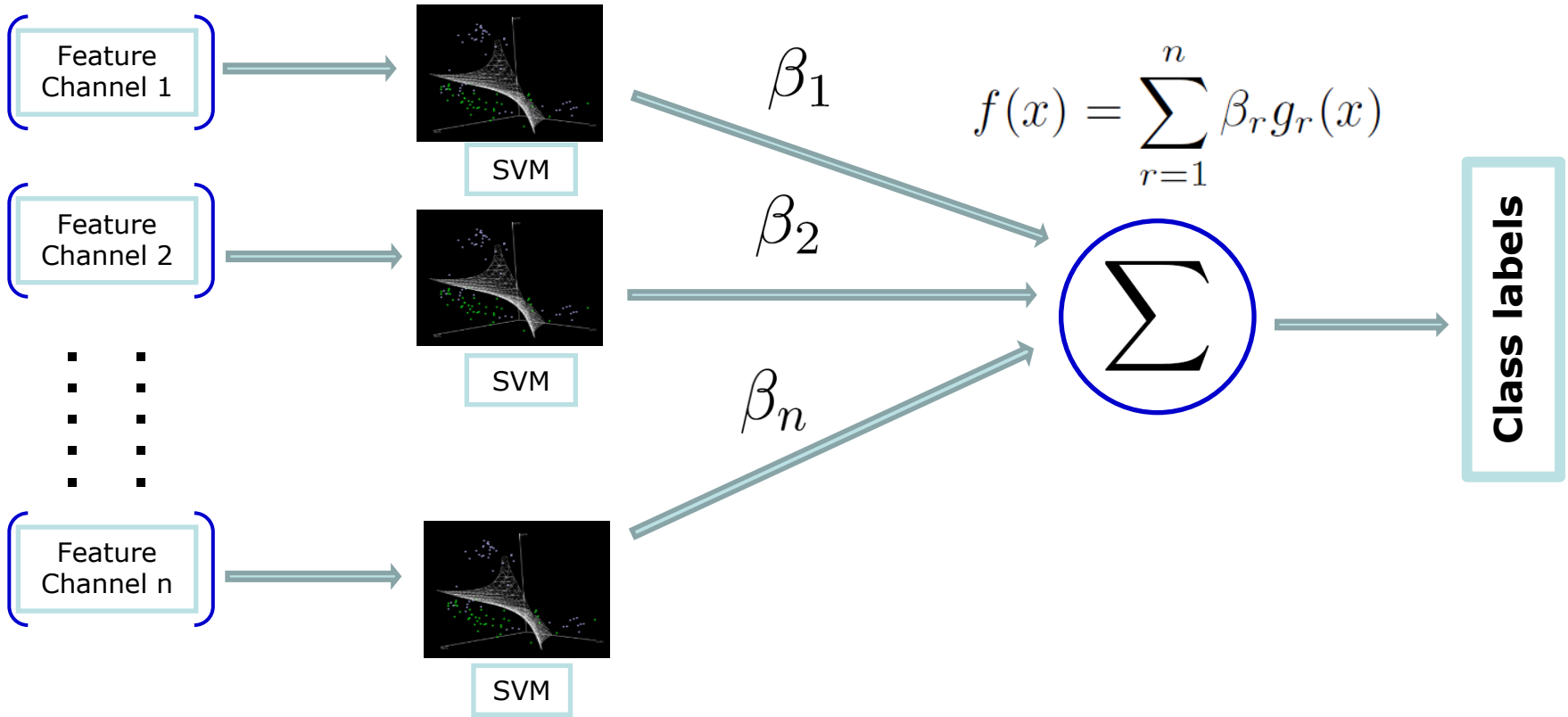
# Classifier Level Fusion (CLF)



# Classifier Level Fusion (CLF)



# Classifier Level Fusion (CLF)



$$\rho := \min_{1 \leq i \leq m} y_i f(x_i) = \min_{1 \leq i \leq m} y_i \sum_{r=1}^n \beta_r g_r(x_i)$$



# Binary CLF with Non-Linear Constraints

- Extension of  $\nu$ -LP-AdaBoost to arbitrary norms
- Normalized margin

$$\rho := \min_{1 \leq i \leq m} y_i f(x_i) = \min_{1 \leq i \leq m} y_i \sum_{r=1}^n \beta_r g_r(x_i)$$

- Nonlinear Programming CLF

$$\max_{\beta, \xi, \rho} \quad \rho - \frac{1}{\nu m} \sum_{i=1}^m \xi_i$$

$$s.t. \quad y_i \sum_{r=1}^n \beta_r f_r(x_i) \geq \rho - \xi_i \quad \forall i = 1, \dots, m$$

$$\|\beta\|_p^p \leq 1, \quad \beta \succeq 0, \xi \succeq 0, \rho \geq 0$$

# Multiclass CLF (NLP-vMC)

- Novel multiclass CLF
- Margin redefinition

$$\rho_i(x_i, \beta) := \sum_{r=1}^n \beta_{(N_C(r-1)+y_i)} g_{r,y_i}(x_i) - \sum_{r=1}^n \sum_{j=1, j \neq i}^{N_C} \beta_{(N_C(r-1)+y_j)} g_{r,y_j}(x_i)$$

- Nonlinear Programming (NLP-vMC)

$$\begin{aligned} \max_{\beta, \xi, \rho} \quad & \rho - \frac{1}{\nu m} \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & \sum_{r=1}^n \beta_{(N_C(r-1)+y_i)} g_{r,y_i}(x_i) - \sum_{r=1}^n \sum_{j=1, j \neq i}^{N_C} \beta_{(N_C(r-1)+y_j)} g_{r,y_j}(x_i) \geq \rho - \xi_i, \\ & \|\beta\|_p^p \leq 1, \quad \rho \geq 0, \beta \succeq 0 \quad \xi \succeq 0 \quad \forall i = 1, \dots, m \end{aligned}$$

# Multiclass CLF (NLP- $\beta$ )

## ➤ Nonlinear Programming- $\beta$ (NLP- $\beta$ )

$$\begin{aligned} \max_{\beta, \xi, \rho} \quad & \rho - \frac{1}{\nu m} \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & \sum_{r=1}^n \beta_r g_{r, y_i}(x_i) - \max_{y_j \neq y_i, r=1}^n \sum_{r=1}^n \beta_r g_{r, y_j}(x_i) \geq \rho - \xi_i, \quad \forall i = 1, \dots, m \\ & \|\beta\|_p^p \leq 1, \quad \beta_r \geq 0, \quad \xi_i \geq 0, \quad \rho \geq 0, \quad \forall r = 1, \dots, n, \forall i = 1, \dots, m. \end{aligned}$$

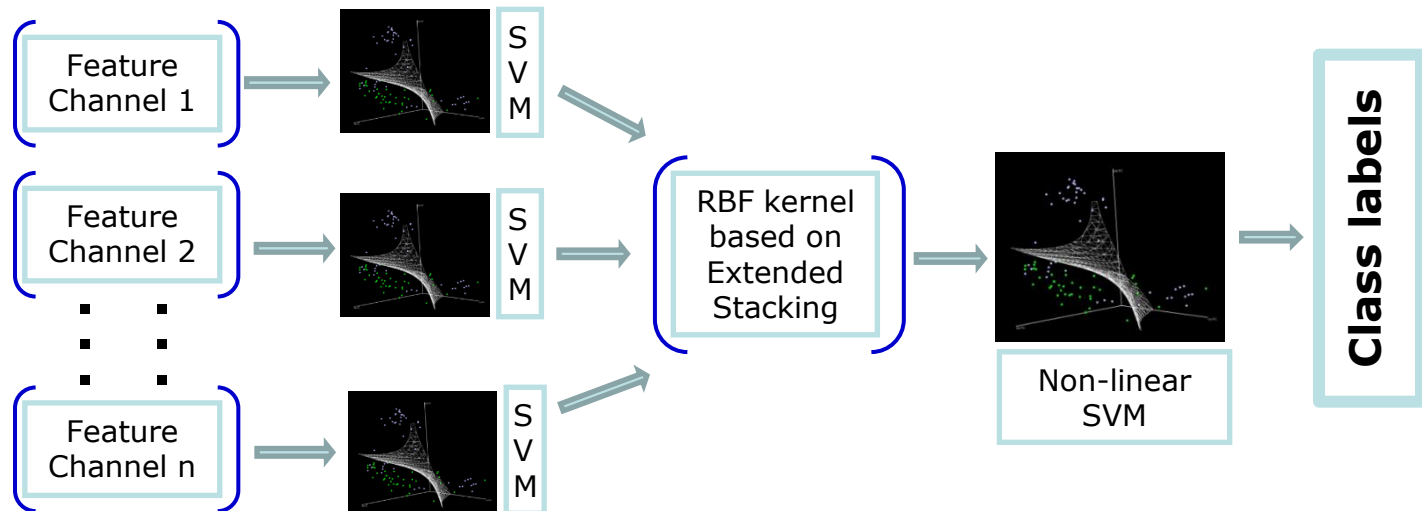
# Multiclass CLF (NLP-B)

## ➤ Nonlinear Programming-B (NLP-B)

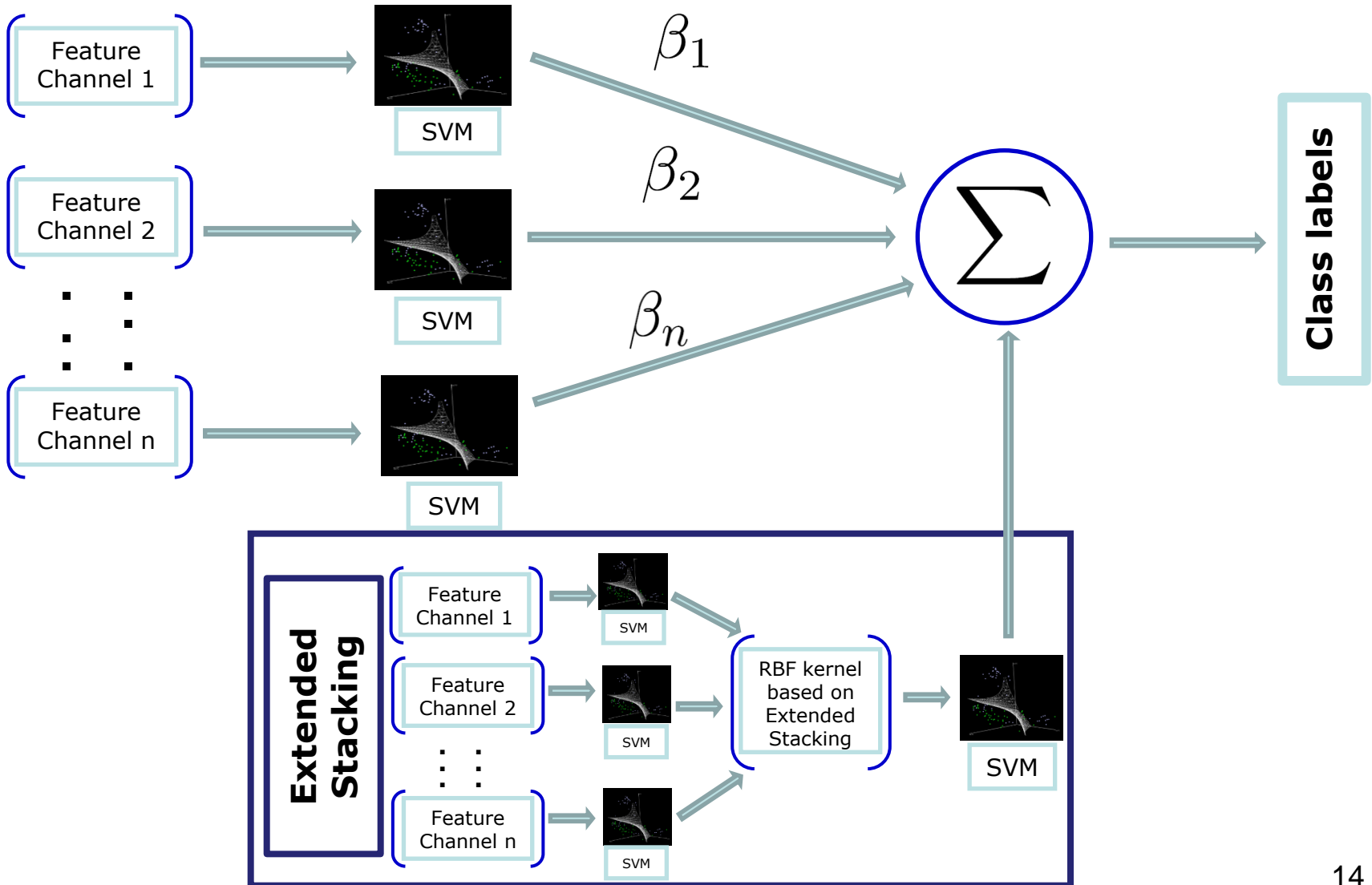
$$\begin{aligned} \max_{B, \xi, \rho} \quad & \rho - \frac{1}{\nu m} \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & \sum_{r=1}^n B_r^{y_i} g_{r, y_i}(x_i) - \sum_{y_j \neq y_i, r=1}^n B_r^{y_j} g_{m, y_j}(x_i) \geq \rho - \xi_i \quad i = 1, \dots, m, \\ & \|B\|_p^p \leq 1, \quad B_r^c \geq 0, \quad \xi \succeq 0, \rho \geq 0, \quad \forall r = 1, \dots, n, c = 1, \dots, N_C \end{aligned}$$

# Extended Stacking

- Break down multiclass problem into 1-vs-all.
  - For each sample its distances from all hyperplanes of 1-vs-all classifier is used as base feature.
- Break down multilabel problem into independent binary problem.
  - For each sample its distances from all hyperplanes of independent binary classifier is used as base feature.



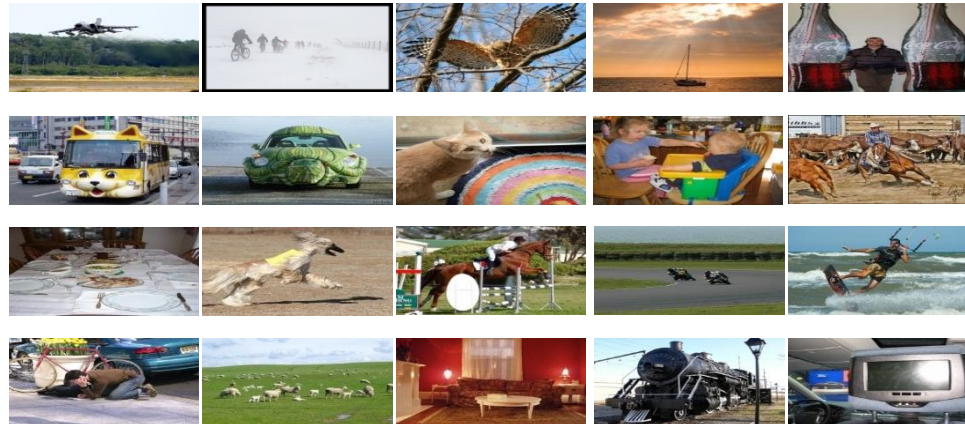
# Base plus Stacking Kernel



# Results

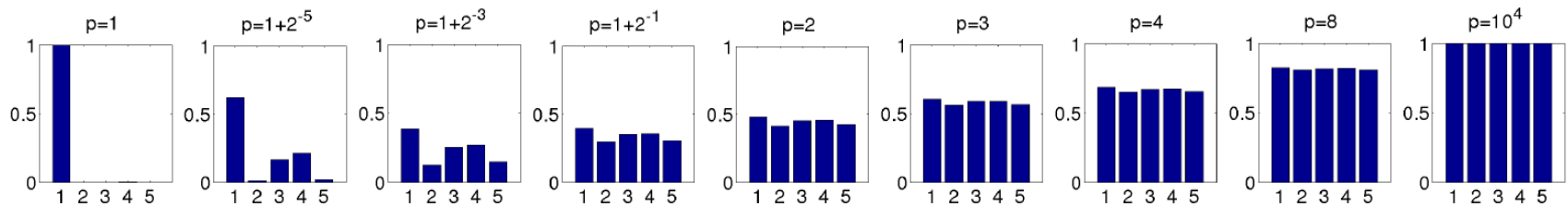
- Multi-Label Dataset (PASCAL VOC2007)
  - 20 classes, 9963 images from internet
  - Mean Average Precision (MAP)

$$MAP = \frac{1}{|R|} \sum_{k=1}^R C_k, \quad C_k = \begin{cases} \frac{|R \cap M_k|}{k} & \text{if concept true} \\ 0 & \text{if concept not true} \end{cases} \quad M_k = \{i_1, i_2, \dots, i_k\}$$



# PASCAL VOC 2007

- Feature channel weights learned with various  $l_p$ -norm for CLF.
  - Sparsity decreases with higher norms





# PASCAL VOC 2007

➤ Mean Average Precision of different fusion methods.

Fusion Methods	norms								
	1	$1 + 2^{-3}$	$1 + 2^{-2}$	$1 + 2^{-1}$	2	3	4	8	$l_{\infty}$
MKL	55.42	56.42	58.53	61.07	61.98	62.45	62.61	62.81	62.93
CLF	63.71	63.94	63.97	<b>63.98</b>	63.97	63.97	63.77	63.69	63.11
Stacking	<b>64.44</b>								
MKL (Base + Stacking)	64.39	64.55	65.06	65.75	66.06	66.23	<b>66.24</b>	66.09	65.93
CLF (Base + Stacking)	65.18	65.20	65.45	65.57	65.65	65.63	65.59	65.54	65.48

# Results

## ➤ Multi-Class Datasets

- Oxford Flower17 (17 classes)
- Oxford Flower102 (102 classes)
- Caltech101 (101 classes)
- Mean Accuracy
- Protein Subcellular Localization (4 datasets)
  - 1- MCC in percentage



# Mean Accuracy on Oxford Flower17

- Decrease in sparsity led to performance improvement.
- Best results by combining base and stacking kernels.

ML-Methods	1	$1 + 2^{-3}$	$1 + 2^{-1}$	2	3	4	8
MKL	$87.2 \pm 2.7$	$74.9 \pm 1.7$	$72.2 \pm 3.6$	$71.2 \pm 2.7$	$70.6 \pm 3.8$	$73.1 \pm 3.9$	$81.0 \pm 4.0$
NLP- $\beta$	$86.5 \pm 3.3$	$86.6 \pm 3.4$	$86.6 \pm 1.1$	$86.7 \pm 1.2$	$87.4 \pm 1.5$	<b><math>87.9 \pm 1.8</math></b>	$87.8 \pm 2.1$
NLP- $\nu$ MC	$85.5 \pm 1.3$	$86.6 \pm 2.0$	$87.6 \pm 2.2$	$87.7 \pm 2.6$	<b><math>87.8 \pm 2.1</math></b>	$87.7 \pm 2.0$	$87.8 \pm 1.9$
NLP-B	$84.6 \pm 2.5$	$84.6 \pm 2.4$	$84.8 \pm 2.6$	$84.8 \pm 2.5$	$85.5 \pm 3.7$	$86.9 \pm 2.7$	$87.3 \pm 2.7$
Stacking	<b><math>89.4 \pm 0.5</math></b>						
MKL(Base + Stacking)	$89.3 \pm 0.9$	$79.7 \pm 2.7$	$77.6 \pm 1.2$	$74.7 \pm 2.4$	$73.8 \pm 2.6$	$77.8 \pm 4.3$	$86.3 \pm 1.9$
NLP- $\beta$ (Base + Stacking)	<b><math>90.2 \pm 1.5</math></b>	$89.3 \pm 0.7$	$89.6 \pm 0.5$	$89.2 \pm 1.6$	$89.3 \pm 1.2$	$89.1 \pm 1.4$	$89.0 \pm 1.0$
NLP- $\nu$ MC(Base + Stacking)	$86.1 \pm 2.5$	$87.3 \pm 1.4$	$88.5 \pm 0.5$	$88.6 \pm 0.9$	$88.6 \pm 0.9$	$88.8 \pm 1.1$	$88.9 \pm 1.2$
Comparison with State-of-the-Art							
MKL-FDA ( $\ell_p$ ) [Yan et al. CVPR10](best state-of-the-art using 7 kernels)							<b><math>86.7 \pm 1.2</math></b>
MKL-avg ( $\ell_\infty$ )(using 7 kernels)							$84.9 \pm 1.9$
CLF ( $\ell_\infty$ ) (using 7 kernels)							$86.7 \pm 2.7$

# Summary of Mean Accuracy on Computer Vision Datasets

- Similar trend was observed for oxford flower102 and caltech101.
  - Summary of results is given in table (more details in paper).

ML-Methods	Oxford flower 17	Oxford flower 102	caltech101
state-of-the-art-best	86.7±1.2	72.8	68.6±2.2
proposed Extended Stacking	89.4±0.5	77.7	68.0±2.4
proposed-CLF-best	<b>90.2±1.5</b>	<b>80.3</b>	<b>70.7±1.9</b>

# Protein Subcellular Localization

- Results on 4 bioinformatics datasets validate our experiments.
  - Prediction error is measured as 1- Matthew Correlation Coefficient (MCC) in percentage.

ML-Methods	plant	nonpl	psortNeg	psortPos
SVM-best [kloft et al. JMLR11]	8.18±0.47	8.97±0.26	9.87±0.34	13.01±0.63
FDA-best [Yan et al. JMLR11]	10.85±2.37	10.84±1.72	9.74±2.00	12.59±4.10
proposed-best	<b>5.25±1.88</b>	<b>6.39±1.12</b>	<b>9.16±1.67</b>	<b>10.40±3.56</b>

# Conclusions

- A novel nonlinear separable convex optimization formulation for multiclass classifier fusion.
  - Learns weight of each class in each feature channel.
- Arbitrary norm in the existing CLF formulation
  - Don't reject informative channels
  - Robust against noisy and redundant channels
  - Norm can be learnt using validation set
- Extended stacking for binary and multiclass problems.
  - Stacking plus base kernels gives best results
- Extensive evaluation on challenging datasets.

*Thanks!*  
*Questions ???*