

Classification générique d'images : Approches aléatoires et convolutionnelles

Begon Jean-Michel

Mémoire de fin d'études réalisé en vue de l'obtention du grade de Master en sciences informatiques

UNIVERSITÉ DE LIÈGE

Faculté des sciences appliquées
Année académique 2013-2014

Plan

- ▶ Introduction
- ▶ La méthode RandConv
- ▶ Les modes de classification
- ▶ Base de données d'évaluation
- ▶ Résultats
- ▶ Conclusion et perspectives

Introduction

- ▶ **Classification générique d'images**

- ▶ Apprentissage automatique

Objet	X1	X2	X3	...	Xn	Classe
#1	45	10	11	...	255	Chien
#2	255	27	200	...	14	Bateau
...
#m	16	0	4	...	144	Avion

- ▶ **Images (données structurées)**

- ▶ Transformer la collection d'objets en matrice d'apprentissage

- ▶ Descripteur d'images

La méthode RandConv

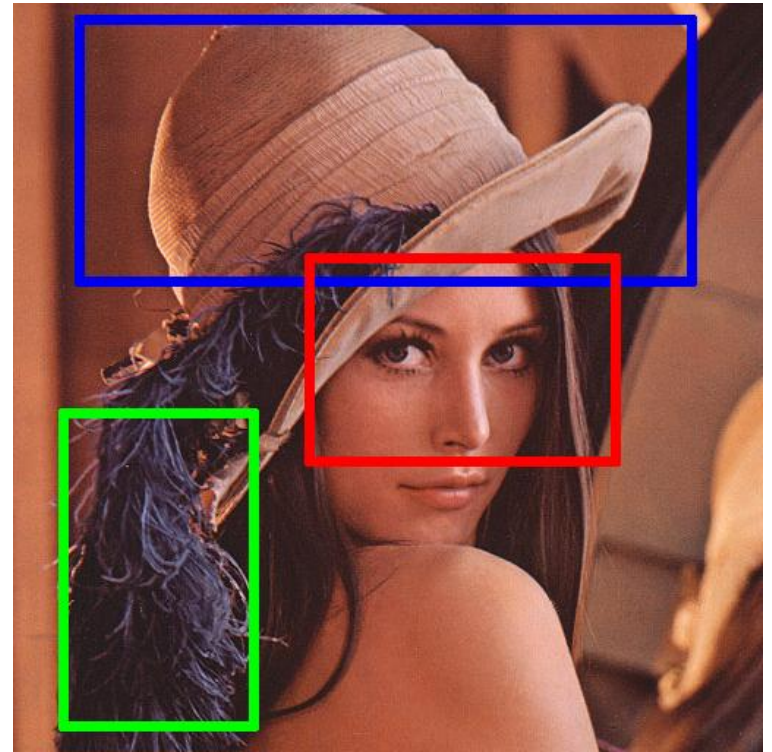
▶ Objectif :

- ▶ Combiner les avantages de la méthode Pixit
 - ▶ Simplicité
 - ▶ Temps de calcul (et d'optimisation)
- ▶ à ceux des réseaux de neurones à convolution (ConvNets)
 - ▶ Exactitude

▶ Principe :

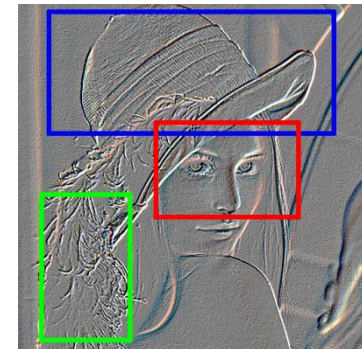
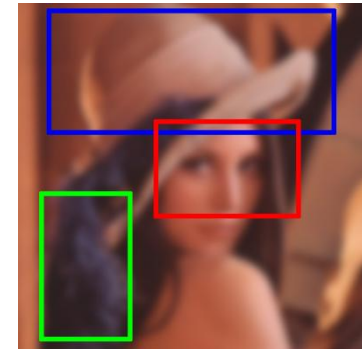
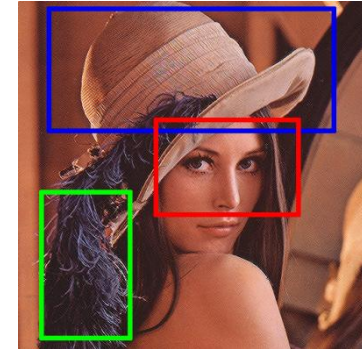
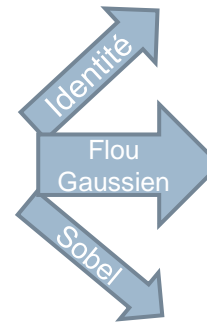
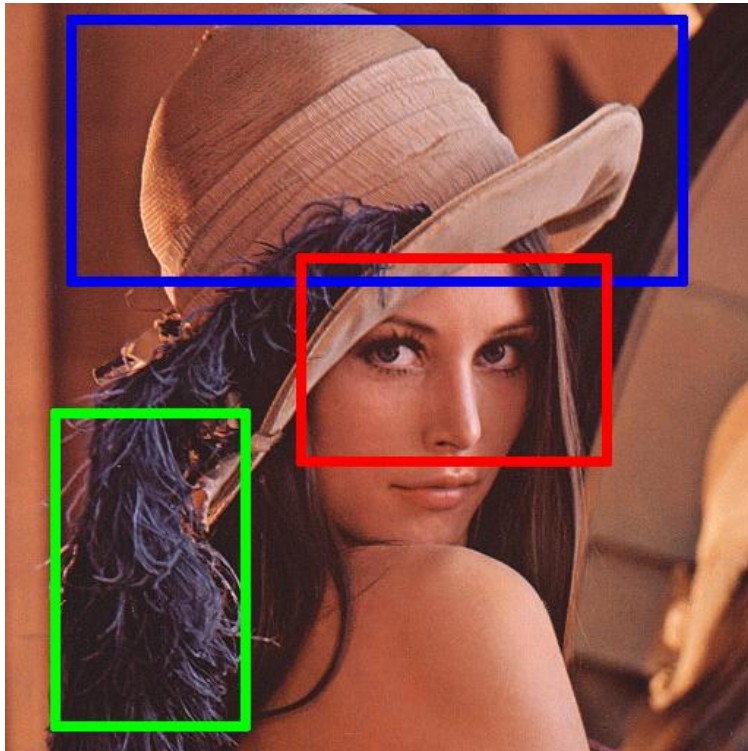
- ▶ Application de filtres aléatoires
 - ▶ Échantillonnage spatial
 - ▶ Extraction de sous-fenêtres et redimensionnement
 - ▶ Description par les pixels bruts
- } ConvNets
- } Pixit

La méthode RandConv : Sélection des sous-fenêtres (Pixit)

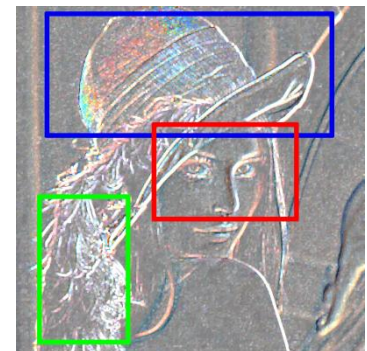
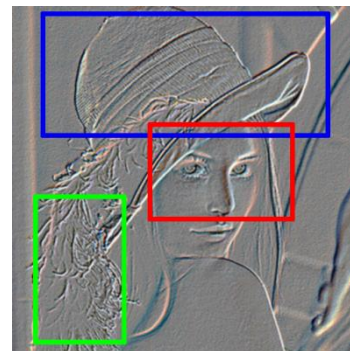
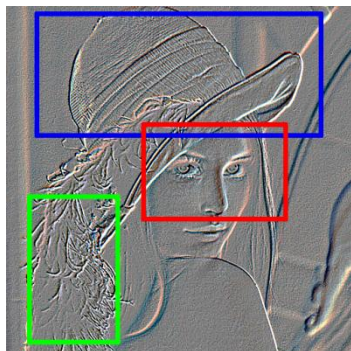
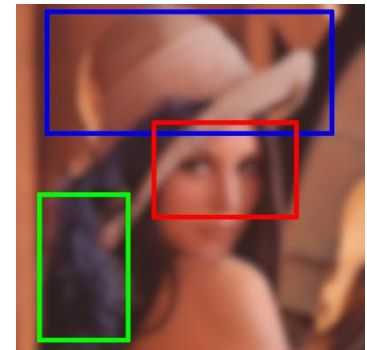
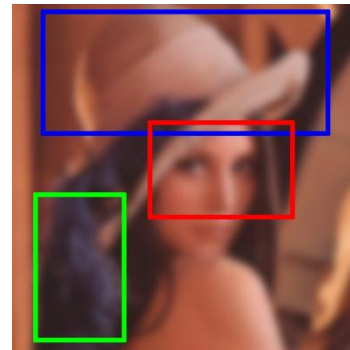
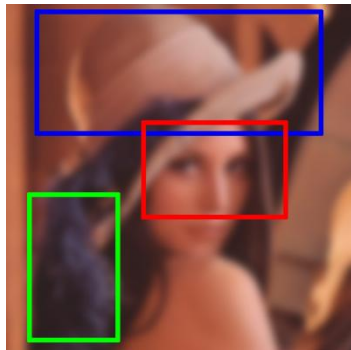
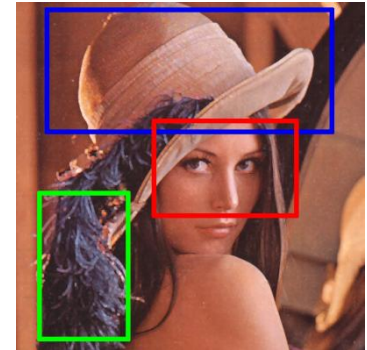
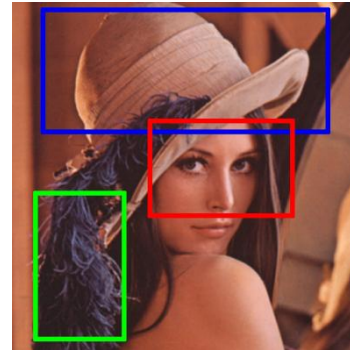
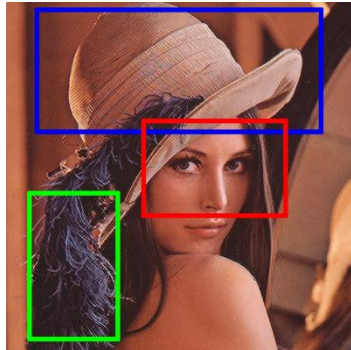


- ▶ Tirage aléatoire uniforme

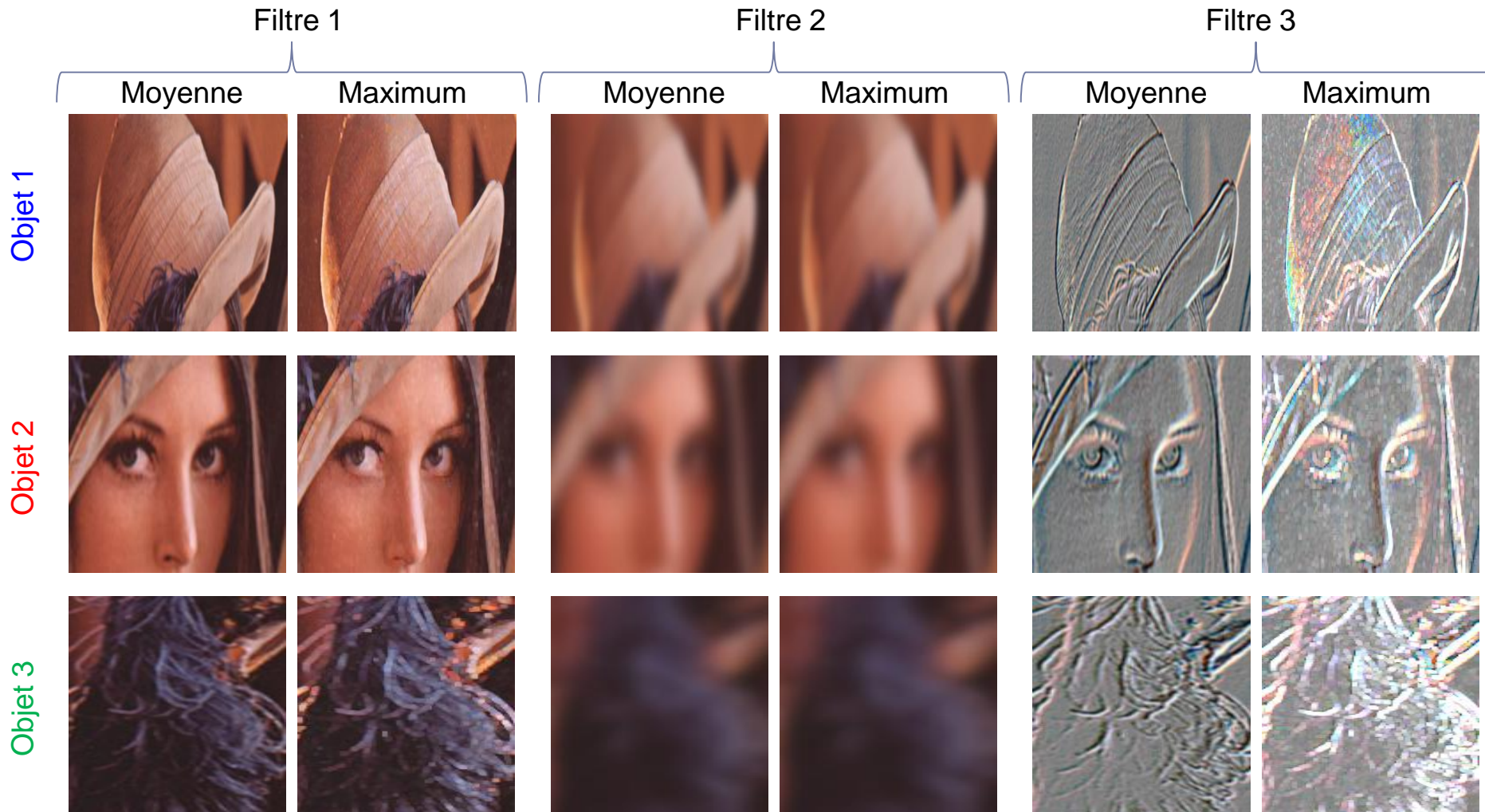
La méthode RandConv : Filtrage (ConvNets)



La méthode RandConv : Échantillonnage spatial (ConvNets)

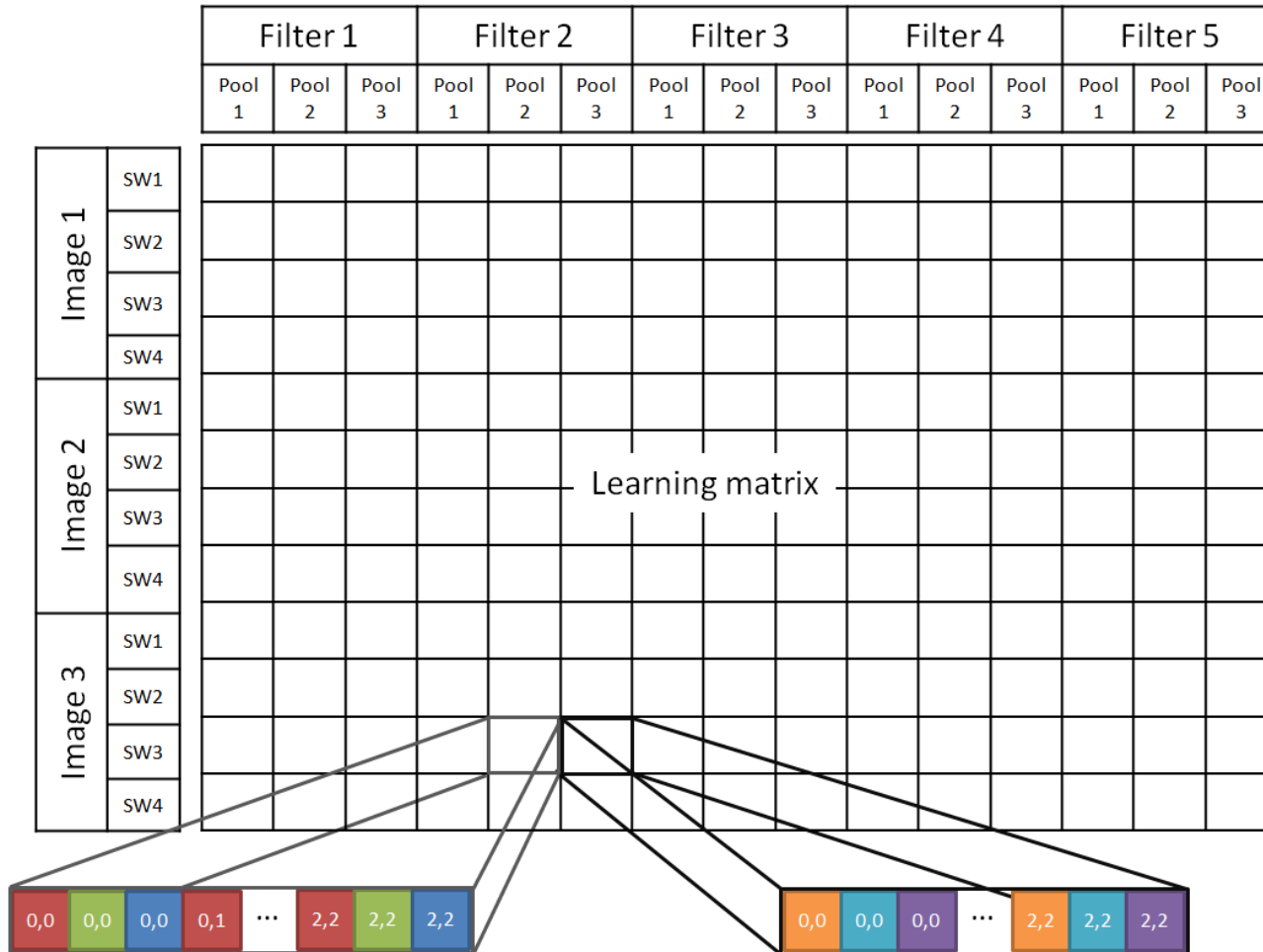


La méthode RandConv : Extraction et redimensionnement des sous-fenêtres (Pixit)



La méthode RandConv : Construction de la matrice d'apprentissage

Learning matrix layout



Les modes de classification

▶ Deux modes

▶ ET-DIC (*ExtraTrees for Direct Image Classification*)

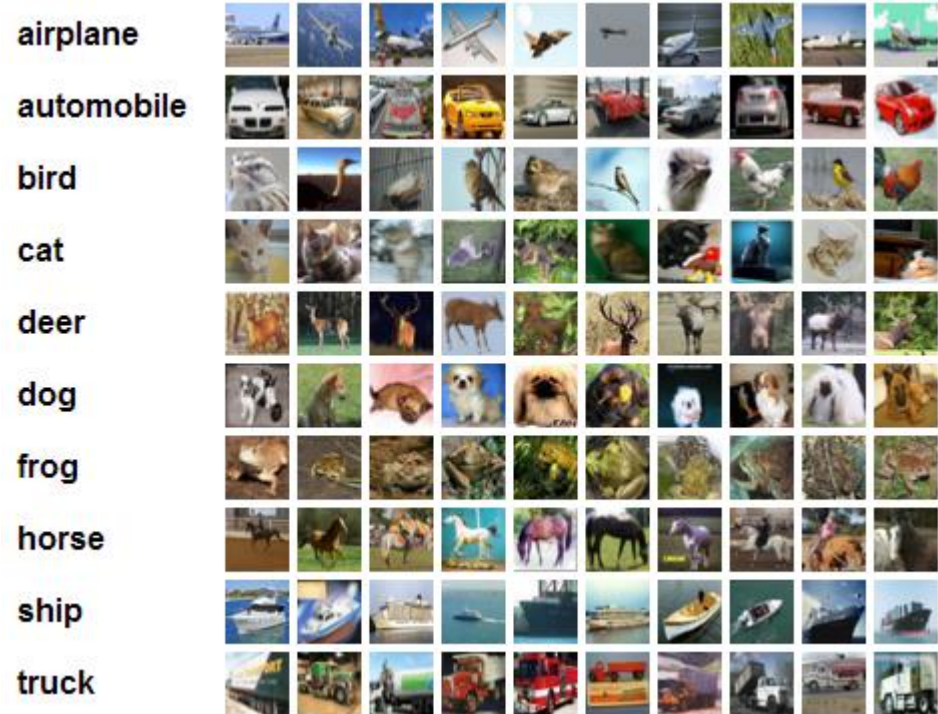
- ▶ Soumission de la matrice d'apprentissage aux *ExtraTrees*
- ▶ Agrégation des prédictions des sous-fenêtres d'une même image (comme Pixit)

▶ ET-FL (*ExtraTrees for Feature Learning*)

- ▶ Création d'un dictionnaire (*ExtraTrees*)
- ▶ Histogramme de « mots visuels » sur base des sous-fenêtres d'une même image
- ▶ Classification par une machine à vecteurs de support (SVM)

Base de données d'évaluation : CIFAR-10

- ▶ 50 000 images d'apprentissage
- ▶ 10 000 images de test
- ▶ 32x32 RGB
- ▶ 10 catégories
- ▶ Difficile pour Pixit
 - ▶ 53,67%
- ▶ Meilleur résultat :
 - ▶ 91.2% (ConvNets)



Résultats du mode ET-DIC

▶ Meilleurs résultats de Pixit (ET-DIC)

- ▶ 53.67% (20 sous-fenêtres, 10 arbres, $n_{\min} = 10$)
- ▶ 49.77% (10 sous-fenêtres, 30 arbres complètement développés)

▶ RandConv

- ▶ 51.51% (10 sous-fenêtres, 30 arbres complètement développés, 100 filtres aléatoires, fenêtre glissante moyennante 3x3)
- ▶ Faible variabilité de l'exactitude
- ▶ Grande stabilité de l'importance des filtres

Peu important	Important
Filtres (nombre et type)	Paramètres des arbres Nombre de sous-fenêtres Échantillonnage spatial (fonction , échelle)

Résultats du mode ET-DIC

- ▶ RandConv
 - ▶ Meilleur résultat : 63,30% (~24h de calculs)
 - ▶ Optimisation des paramètres
 - ▶ Présélection des filtres
 - ▶ Inutile
 - ▶ Combinaison de forêts
 - ▶ Gain lié au nombre d'arbres

Résultats du mode ET-FL

- ▶ **Pixit : 58.79%**
 - ▶ 750 *Totally randomized trees*, $n_{\min} = 500$, fenêtre glissante moyennante 3x3
- ▶ **RandConv :**
 - ▶ 63.55%
 - ▶ Filtres prédéfinis, 750 *Totally randomized trees*, $n_{\min} = 500$, fenêtre glissante moyennante 3x3
 - ▶ 55,07 – 59.37%
 - ▶ Filtres aléatoires, 750 *Totally randomized trees*, $n_{\min} = 500$, fenêtre glissante moyennante 3x3
- ▶ **Rôle important des filtres**
 - ▶ Disparition de la sélection des « bons » filtres
 - ▶ → Chute d'exactitude
 - ▶ → Grande variabilité
- ▶ **Exactitude – taille du dictionnaire**
 - ▶ Tendance linéaire

Résultats du mode ET-FL

▶ Améliorations

- ▶ Présélection des filtres importants
 - ▶ 55,07 – 59.37 → 62.11%
- ▶ Combinaison de plusieurs ensembles de filtres (63.98%, 63.69%, 61.61%)
 - ▶ Méthode d'ensemble : 67.10%
 - ▶ Construction de plusieurs forêts et agrégation
 - Dictionnaire de même ordre de taille : 64.59%
 - Dictionnaire complet : 68.99%

Résultats du mode ET-FL

▶ Meilleur résultat

- ▶ Présélection de 60 filtres aléatoires parmi 500 avec fenêtres glissantes maximisantes 3x3 et 7x7
 - ▶ 70.74%
- ▶ Filtres prédéfinis avec fenêtres maximum glissantes 3x3 et 7x7
 - ▶ 70.62%, 70.57%, 70.87%, 70.70%, 70.97%
- ▶ Combinaison « *leave-one-out* » pour la création de dictionnaires complets
 - ▶ 73.97% - 74.27%
- ▶ Méthode d'ensemble sur les résultats
 - ▶ 74.40%

Conclusion

▶ **Objectif :**

- ▶ Combiner les avantages de la méthode Pixit à ceux des ConvNets

Sous-objectifs	ET-DIC	ET-FL
Simplicité	Oui (par rapport aux ConvNets)	
Temps	Pixit << RandConv < ConvNets	
Exactitude	Pixit < RC-DIC << ConvNets	Pixit < RC-FL < ConvNets

Perspectives

- ▶ Approfondir le mode ET-FL
- ▶ Confirmer les observations
- ▶ Introduire de la non-linéarité
- ▶ Nouveaux générateurs de filtres linéaires basés sur le domaine fréquentiel

- ▶ D'autres approches sont envisageables pour combiner les deux méthodes

Merci pour votre attention

Slides additionels:

La méthode RandConv : Extraction et redimensionnement des sous-fenêtres

Objet 1

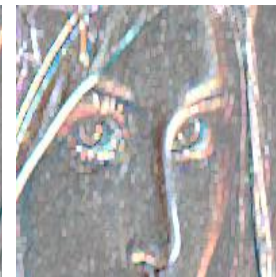
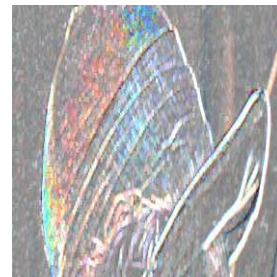
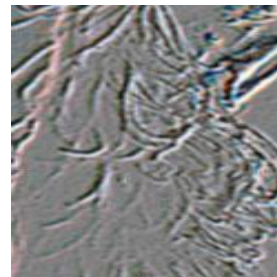
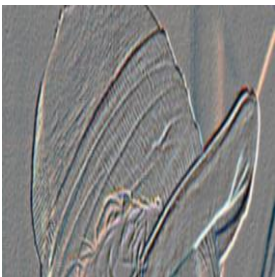
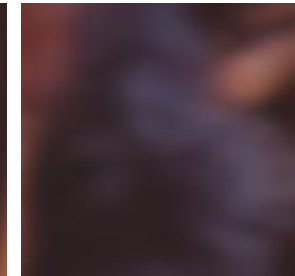
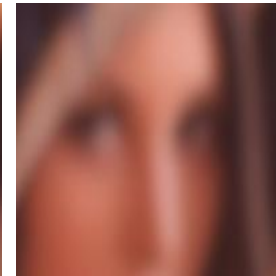
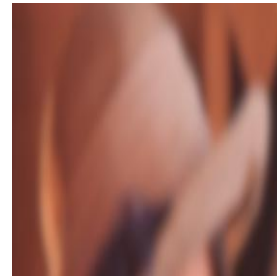
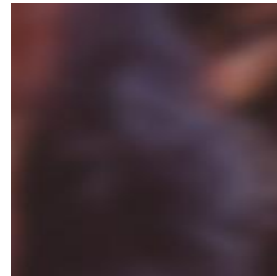
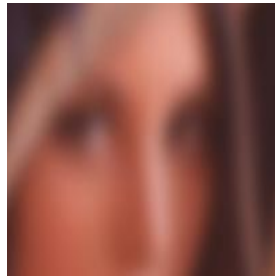
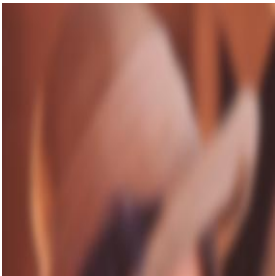
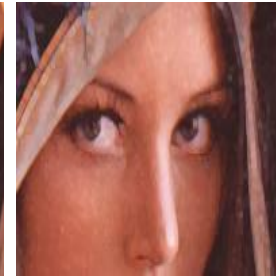
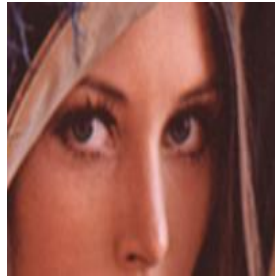
Objet 2

Objet 3

Objet 1

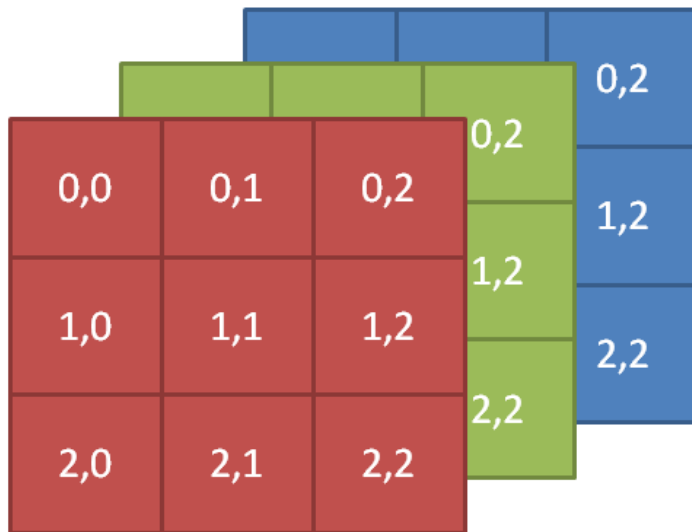
Objet 2

Objet 3

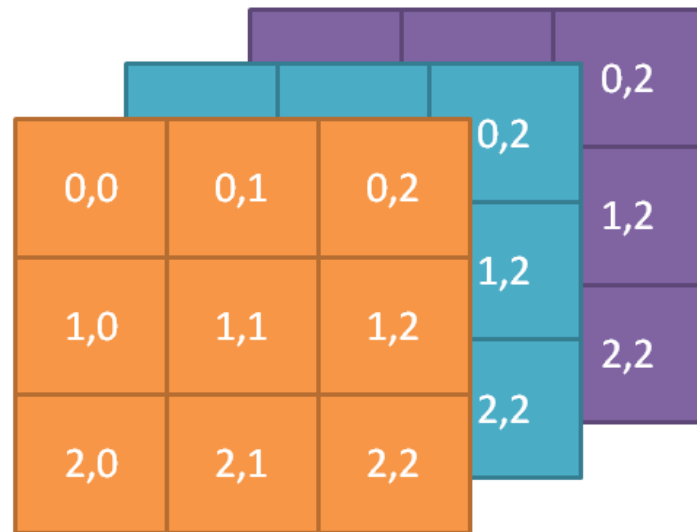


La méthode RandConv : Construction de la matrice d'apprentissage (Pixit)

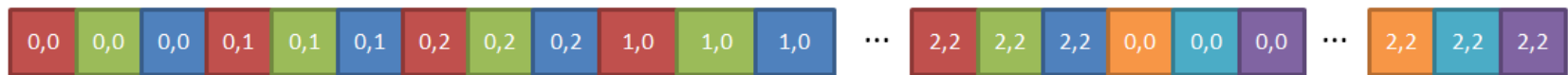
Subwindow description by raw pixels



Subwindow X of image Y
after a given filter and pooling



Subwindow X of image Y
after the next filter and pooling



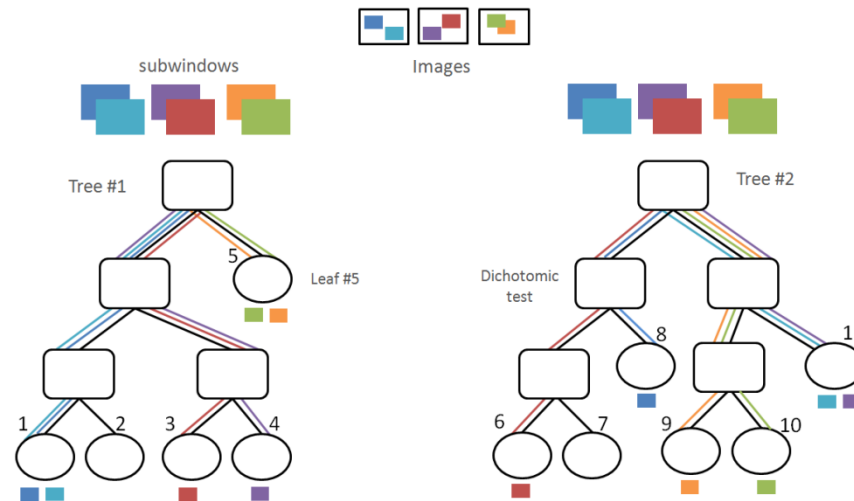
Corresponding feature vector

La méthode RandConv : remarques

- ▶ **Génération des filtres**
 - ▶ Prédéfinis
 - ▶ Centrés sur le filtre nul (perturbation)
 - ▶ Centrés sur le filtre identité
 - ▶ Perturbation
 - ▶ Distance
 - ▶ Stratifiés
- ▶ **Échantillonnage spatial**
 - ▶ Agrégation (voisinage sans chevauchement)
 - ▶ Fenêtre glissante (chevauchement)
 - ▶ Minimum, moyenne, maximum
- ▶ **Limitation**
 - ▶ Coût en mémoire
 - ▶ Coût des temps de calcul

Les modes de classification : ET-FL

ET-FL : construction of the histogram of visual words

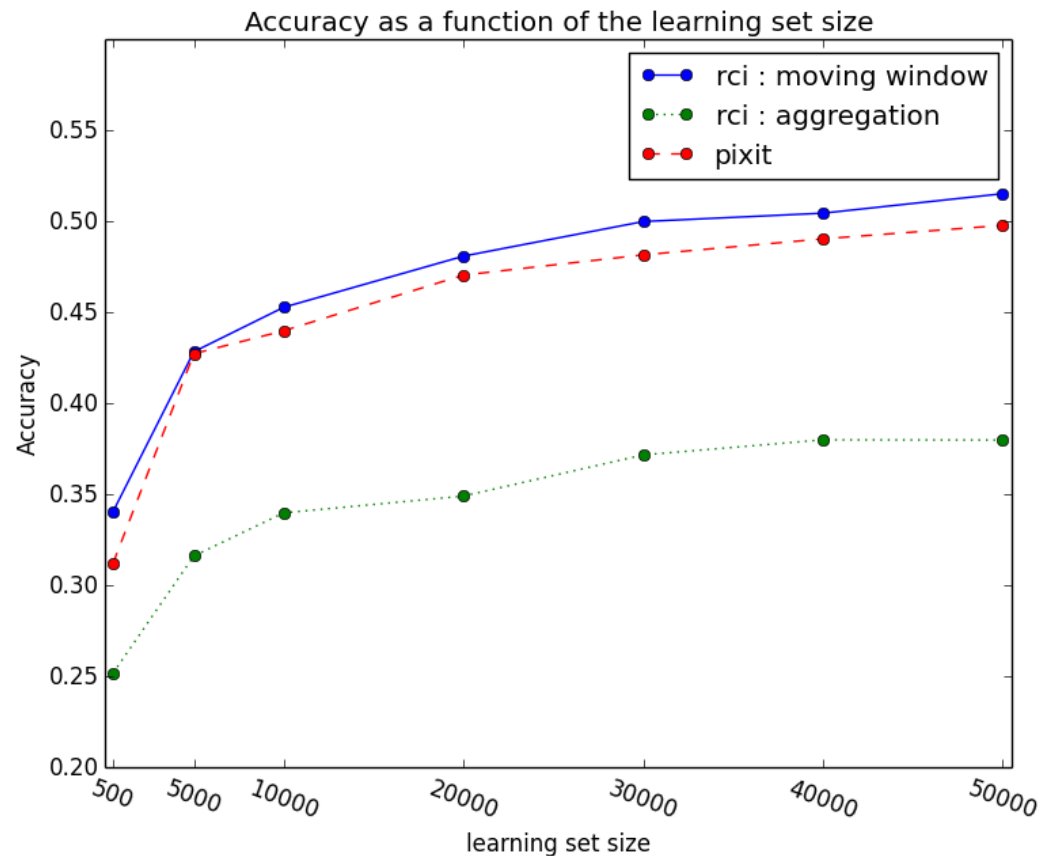


		Leaves										
		1	2	3	4	5	6	7	8	9	10	11
Image 1	Sw. 1											
	Sw. 2											
	Final word	2	0	0	0	0	0	0	1	0	0	1
Image 2	Sw. 1											
	Sw. 2											
	Final word	0	0	1	1	0	1	0	0	0	0	1
Image 3	Sw. 1											
	Sw. 2											
	Final word	0	0	0	0	2	0	0	0	1	1	0

Résultats du mode ET-DIC

- ▶ Influence de la taille de l'ensemble d'apprentissage
 - ▶ Agrégation : inefficace
 - ▶ Légère dominance de RandConv par rapport à Pixit

-30 arbres non-élagués
-10 sous-fenêtres
-100 filtres centrés sur le filtre nul et perturbés uniformément
-Fenêtre mobile moyennante 3x3



Résultats du mode ET-DIC

▶ Variabilité

▶ Liée aux arbres

Statistique	Exactitude
Minimum	51.14%
Moyenne	51.36%
Maximum	51.56%
Écart-type	0.00148

▶ Grande stabilité des importances de filtres

▶ Liée à la matrice d'apprentissage (pour un générateur fixé)

Statistique	Exactitude
Minimum	50.93%
Moyenne	51.23%
Maximum	51.83%
Écart-type	0.00257

Résultats du mode ET-DIC

- ▶ Influence des paramètres des arbres
 - ▶ Nombre d'arbres
 - ▶ Accroissement monotone de l'exactitude
 - ▶ Stabilité de l'importance des filtres
 - ▶ Taille du sous espace aléatoire de variables (k)
 - ▶ Accroissement de l'exactitude avec la taille
 - ▶ Stabilité de l'importance des filtres
 - ▶ Profondeur des arbres (n_{min})
 - ▶ Accroissement de l'exactitude avec la profondeur
 - ▶ Stabilité mitigée de l'importance des filtres
- ▶ Influence du nombre de sous-fenêtres
 - ▶ Accroissement de l'exactitude
 - ▶ Stabilité des filtres

N	Exact.
10	0.492
30 (déf.)	0.515
500	0.529

K	Exact.
1	0.447
277 (déf.)	0.515
10,000	0.528

N_{min}	Exact.
2 (déf.)	0.515
50	0.498
500	0.453

N	Exact.
5	0.48
10 (déf.)	0.515
18	0.521

Résultats du mode ET-DIC

- ▶ Influence des paramètres convolutionnels
 - ▶ Nombre de filtres
 - ▶ Influence faible sur l'exactitude
 - ▶ Taille des filtres
 - ▶ Préférence des petites tailles (3x3 – 9x9) pour CIFAR-10
 - ▶ Influence du générateur
 - ▶ Faible
 - ▶ Générateurs basés sur la distance : moins bons
 - ▶ Influence de l'échantillonnage spatial
 - ▶ Importante
 - ▶ Fonction maximum
 - ▶ Meilleur résultat : 57.41% fenêtre maximum glissante 3 à échelles (3x3, 5x5 et 7x7)
 - ▶ Importances des filtres varient avec la fonction mais peu avec l'échelle

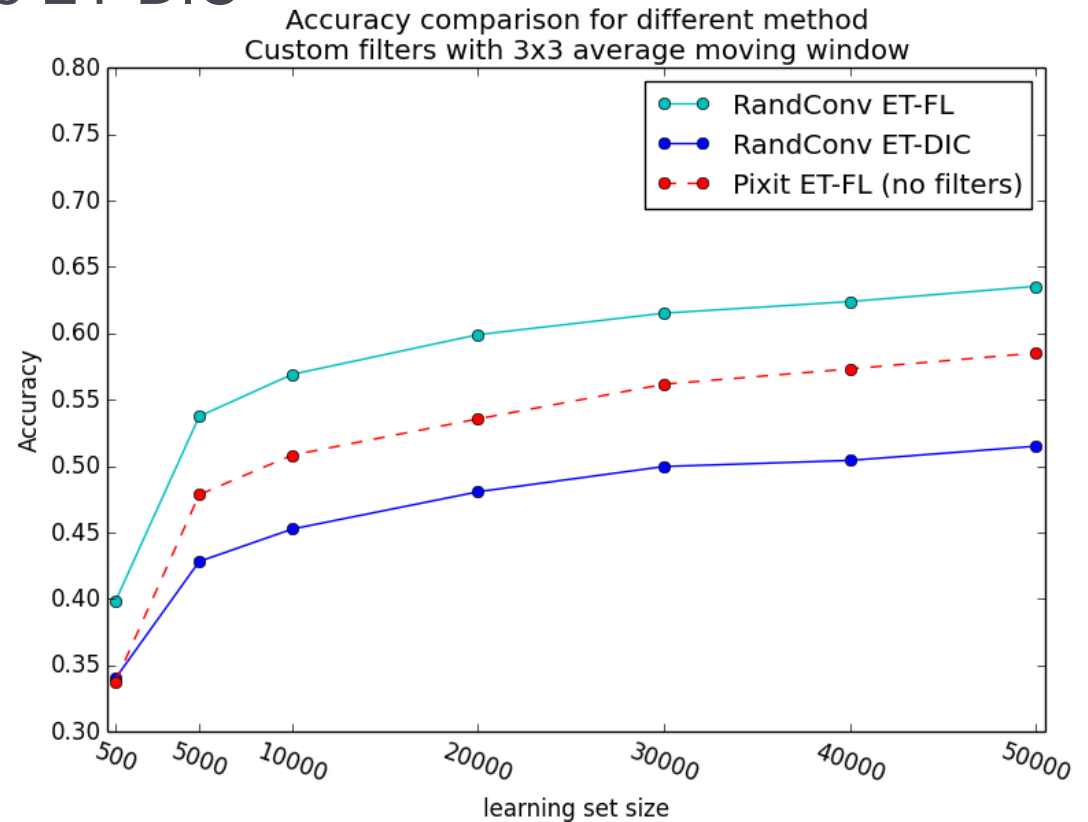
Résultats du mode ET-DIC

- ▶ Borne maximum en optimisant les paramètres
 - ▶ 63,30% (~24h de calculs)
- ▶ Améliorations
 - ▶ Présélection des filtres importants
 - ▶ 51,37%
 - ▶ Création de plusieurs forêts et agrégation des résultats
 - ▶ Combinaisons des 3 générateurs de filtres différents (51.28%, 49.99% et 50.63%) :
 - 52.41% au final
 - A relativiser avec le nombre d'arbres

Résultats du mode ET-FL

- ▶ Influence de la taille de l'ensemble d'apprentissage
 - ▶ Filtres prédéfinis
 - ▶ Meilleur que le mode ET-DIC
 - ▶ Domine le Pixit en ET-FL
 - ▶ 63.55%

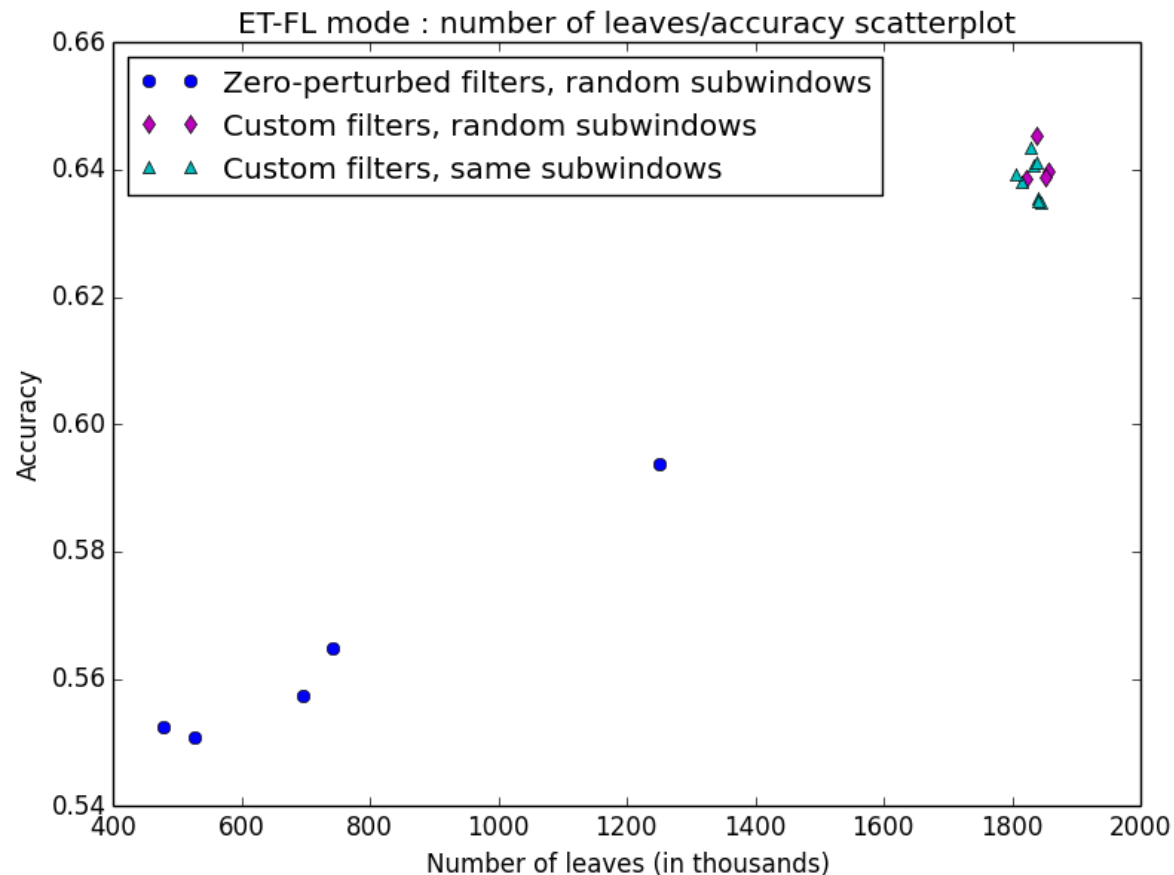
-750 arbres
-n_min = 500
-20 sous-fenêtres
-38 filtres prédéfinis
-Fenêtre mobile
moyennante 3x3



Résultats du mode ET-FL

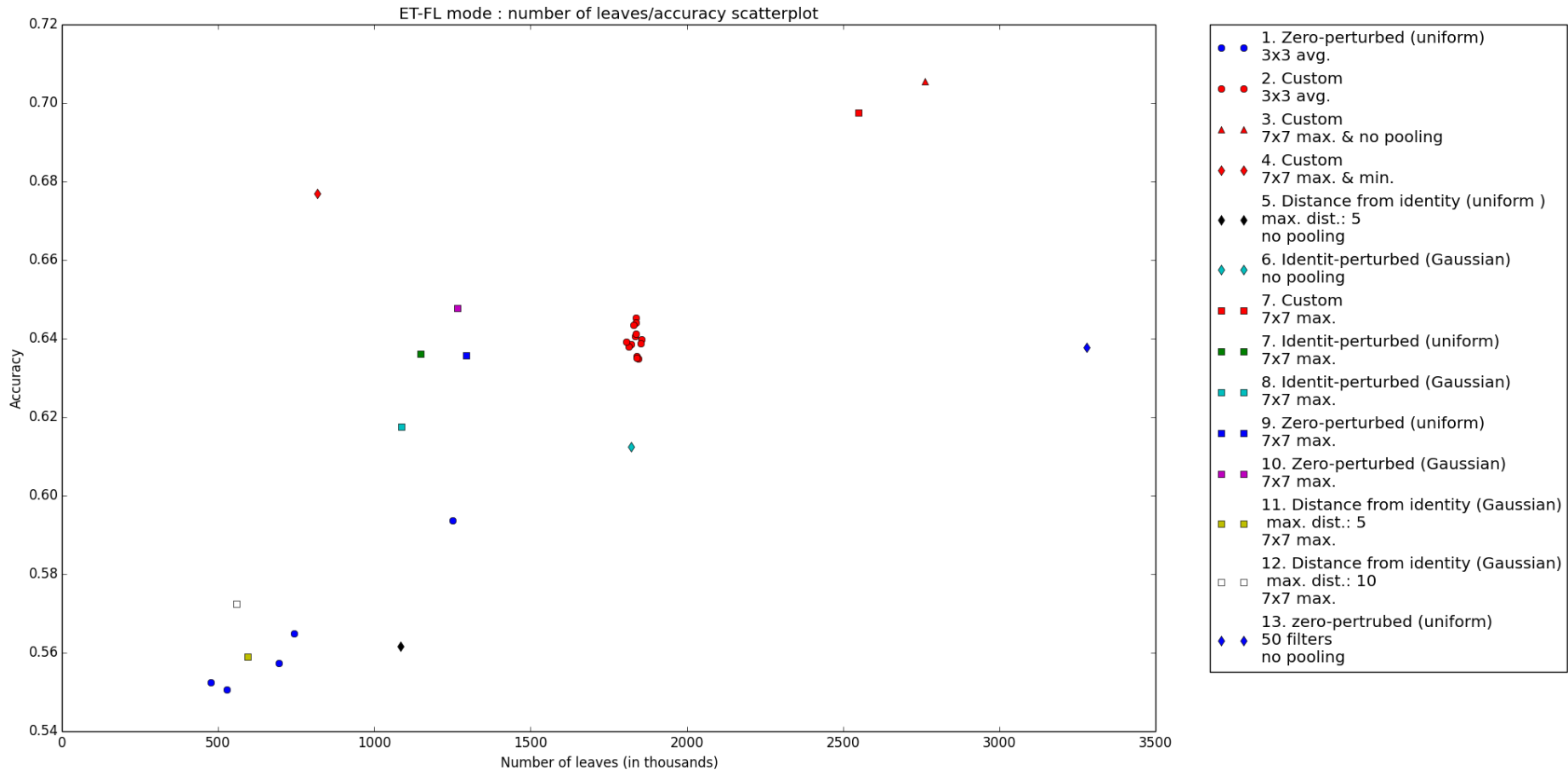
▶ Variabilité

- ▶ Plus importante qu'en ET-DIC
- ▶ Faible influence des sous-fenêtres
- ▶ Forte influence des filtres



Résultats du mode ET-FL

► Relation exactitude – taille du dictionnaire



Illustrations

Bias/variance illustration

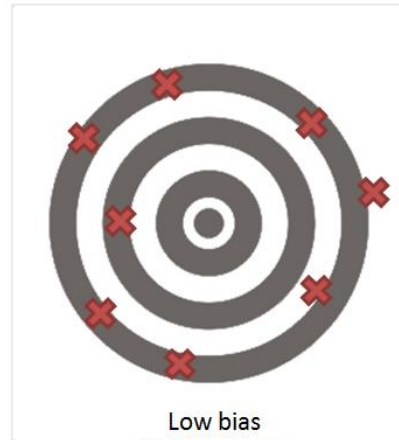


Low bias
Low variance



Bayes model

High bias
Low variance



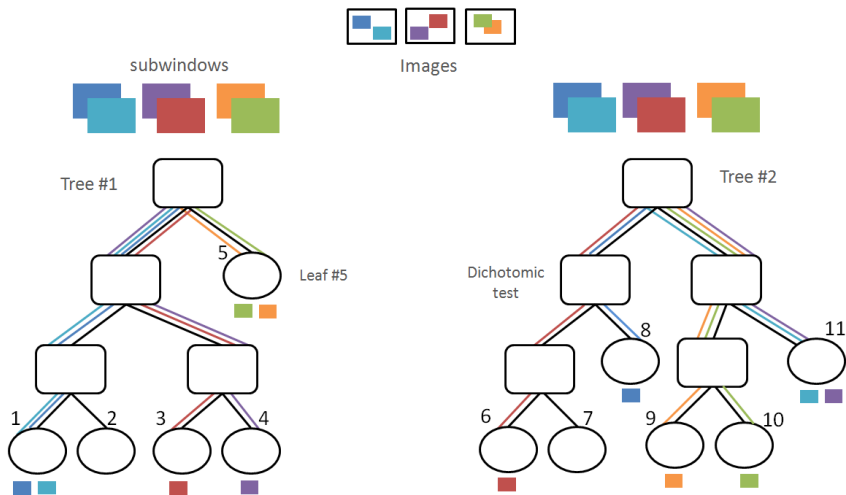
Low bias
High variance



High bias
High variance



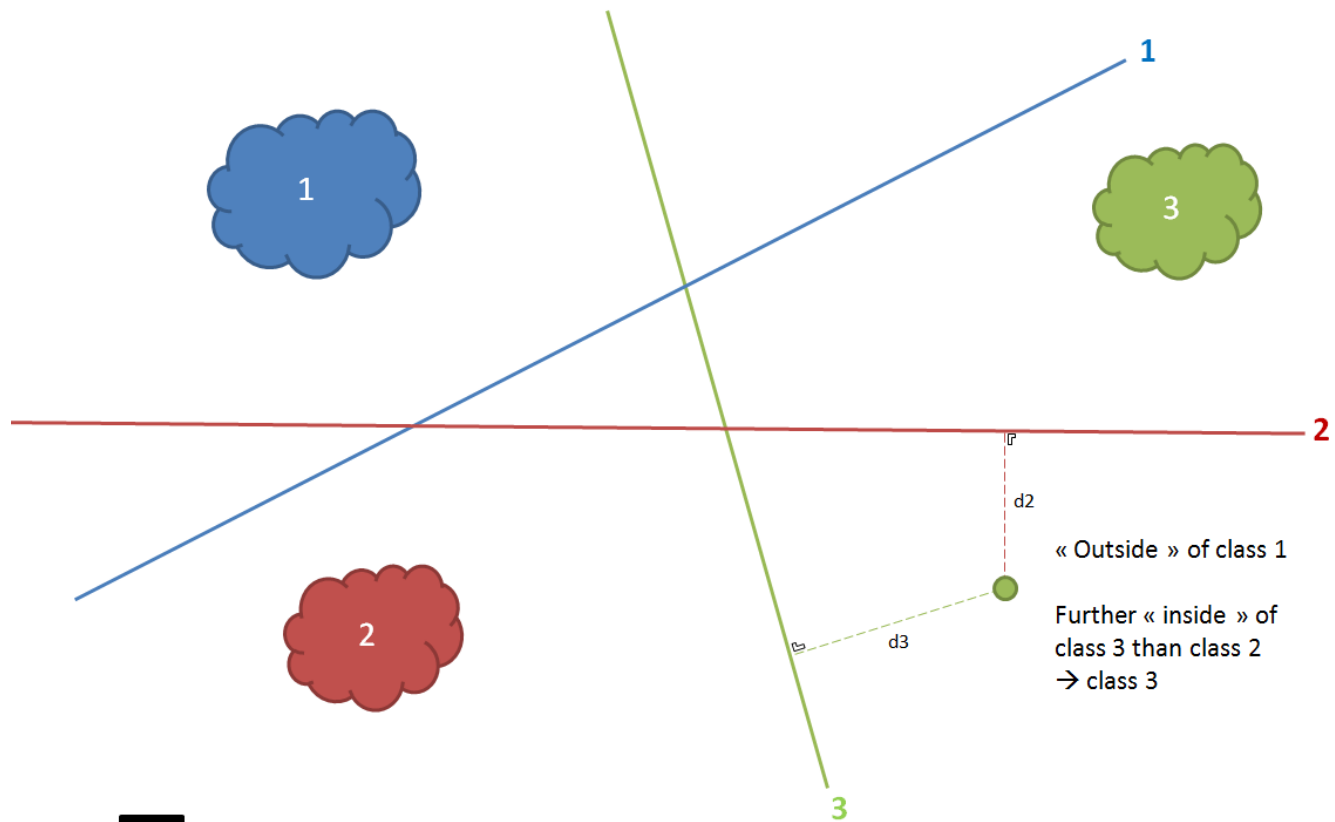
ET-FL : construction of the histogram of visual words :
Forest part



ET-FL : construction of the histogram of visual words :
Histogram part

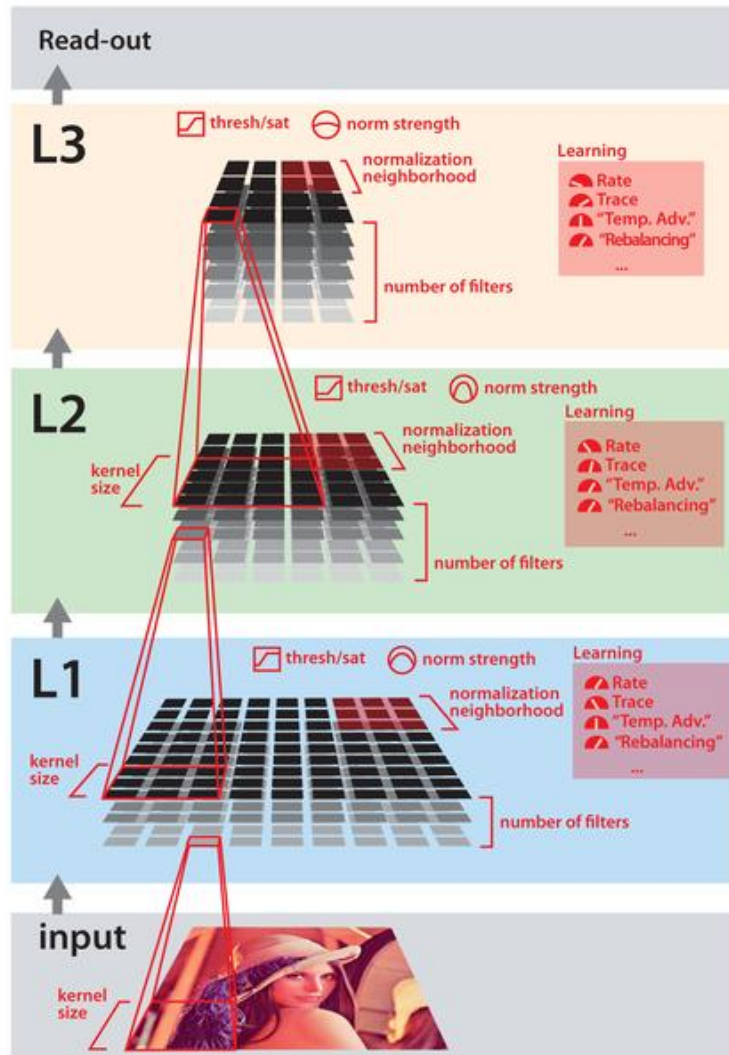
		Leaves										
		1	2	3	4	5	6	7	8	9	10	11
Image 1	Sw. 1											
	Sw. 2											
	Final word	2	0	0	0	0	0	0	1	0	0	1
Image 2	Sw. 1											
	Sw. 2											
	Final word	0	0	1	1	0	1	0	0	0	0	1
Image 3	Sw. 1											
	Sw. 2											
	Final word	0	0	0	0	2	0	0	0	1	1	0

SVM : One-versus-all multiclass scheme



X Cluster « X »

— x Hyperplane separating cluster « X » from the other clusters



Convolutional product

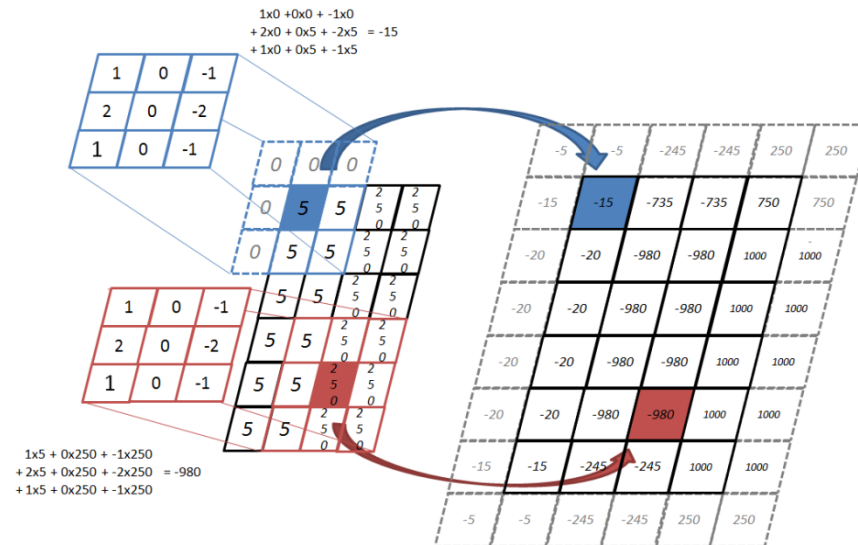
5	5	250	250
5	5	250	250
5	5	250	250
5	5	250	250
5	5	250	250
5	5	250	250

*

1	0	-1
2	0	-2
1	0	-1

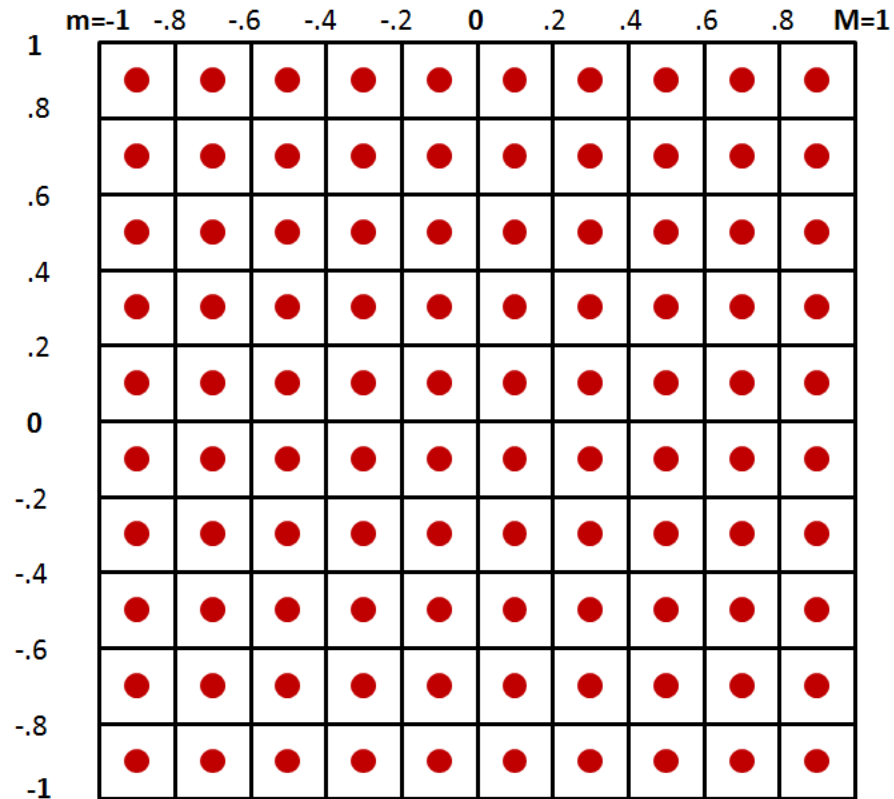
=

-15	-735	-735	750
-20	-980	-980	1000
-20	-980	-980	1000
-20	-980	-980	1000
-20	-980	-980	1000
-15	-735	-735	750



2D subspace of the stratified perturbed generator

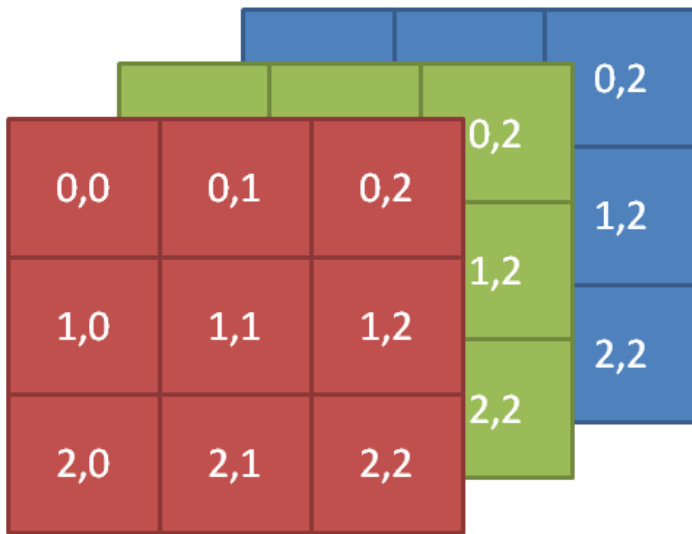
$n = 10$



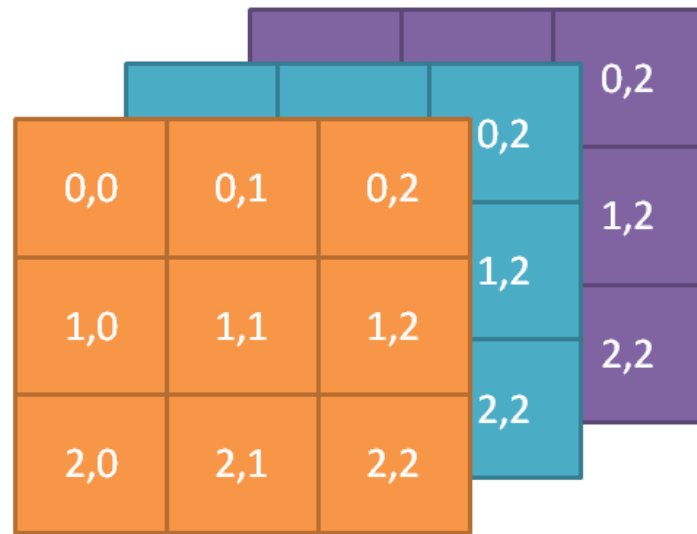
□ cell

● Candidate value
before perturbation

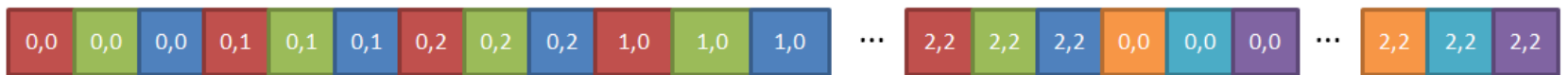
Subwindow description by raw pixels



Subwindow X of image Y
after a given filter and pooling

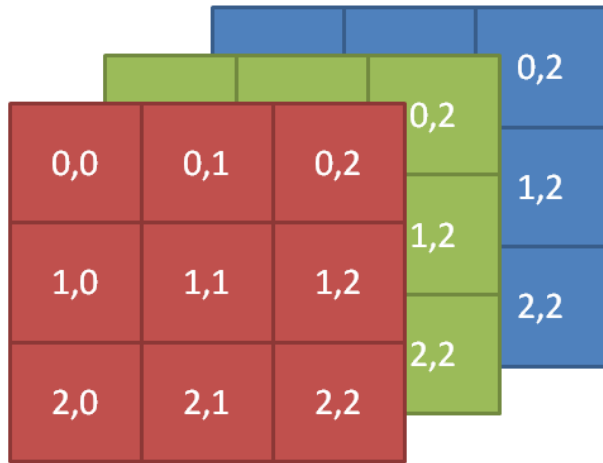


Subwindow X of image Y
after the next filter and pooling

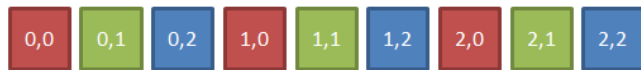


Corresponding feature vector

Image description by raw pixels with compression
1 colors by pixel

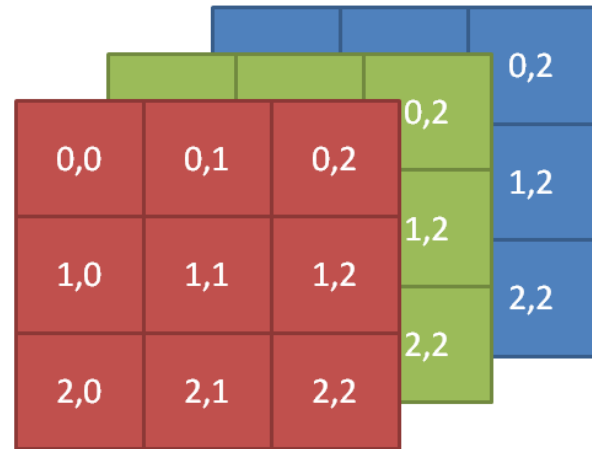


Original image



Corresponding feature vector

Image description by raw pixels with compression
2 colors by pixel

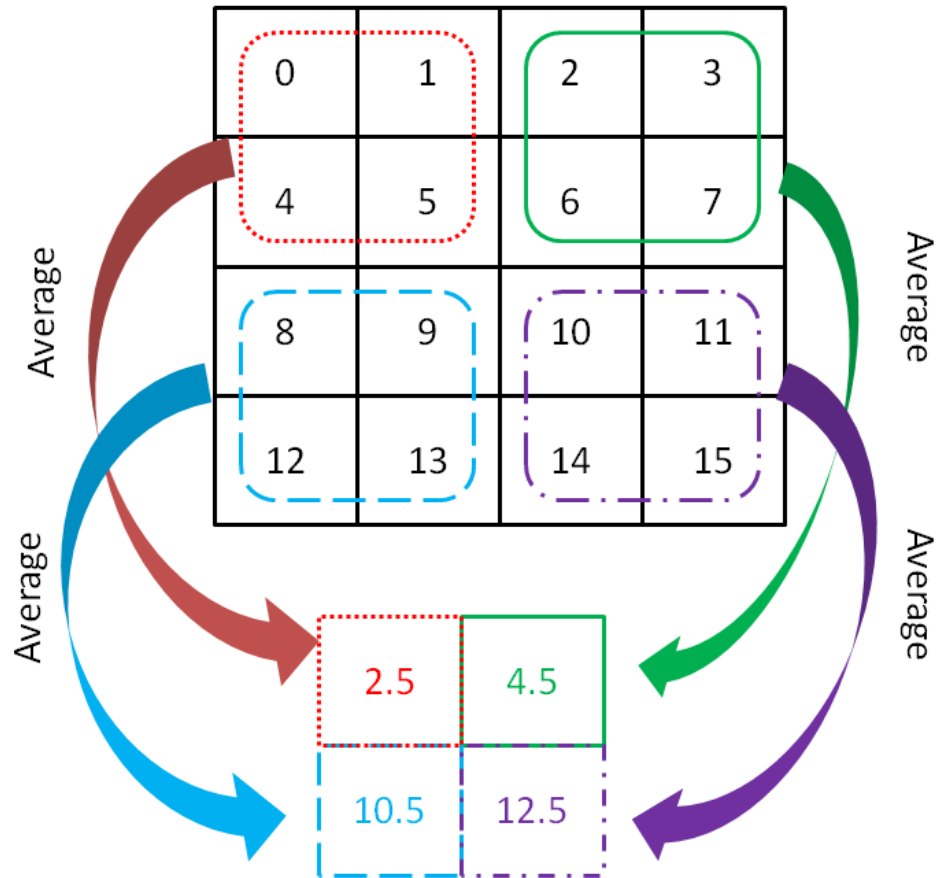


Original image

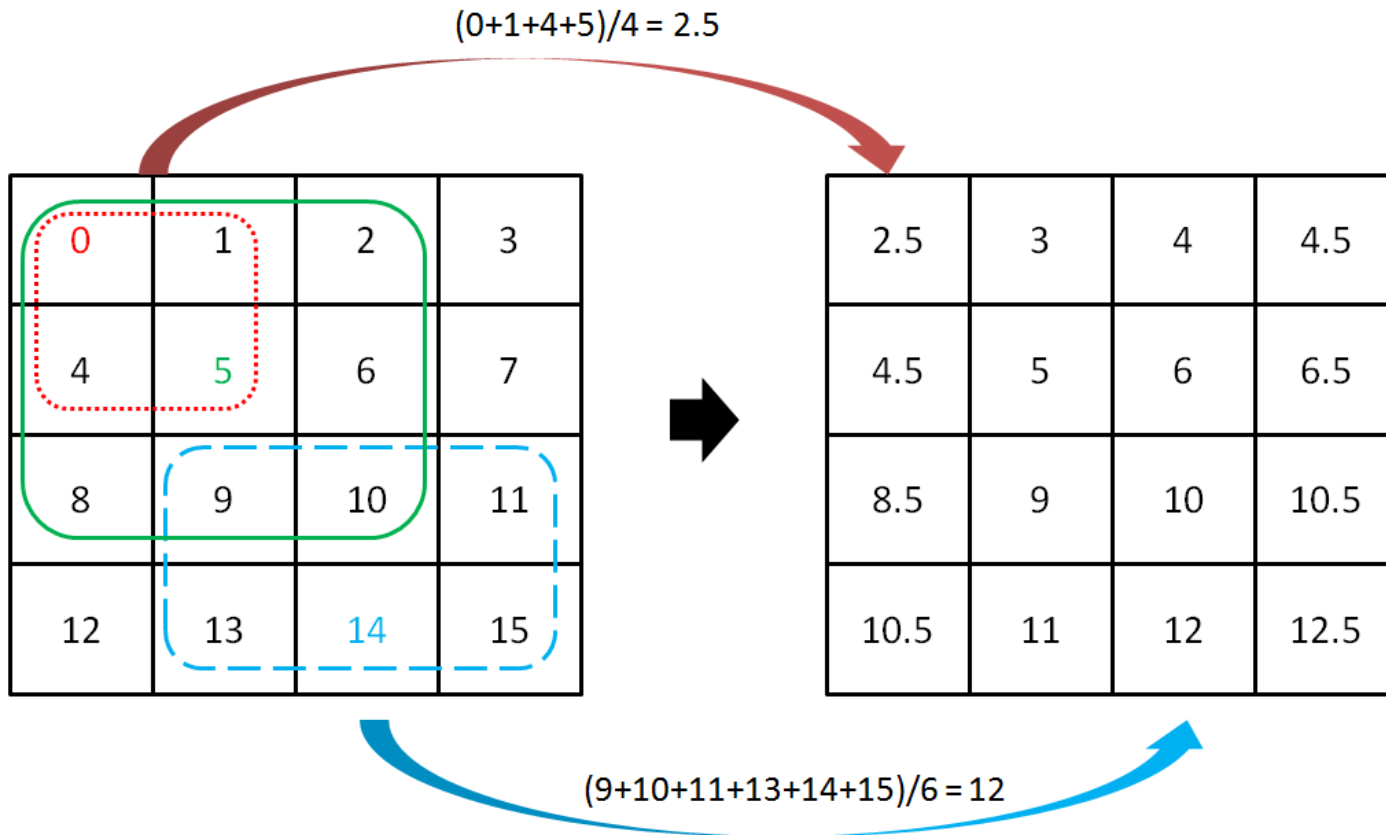


Corresponding feature vector

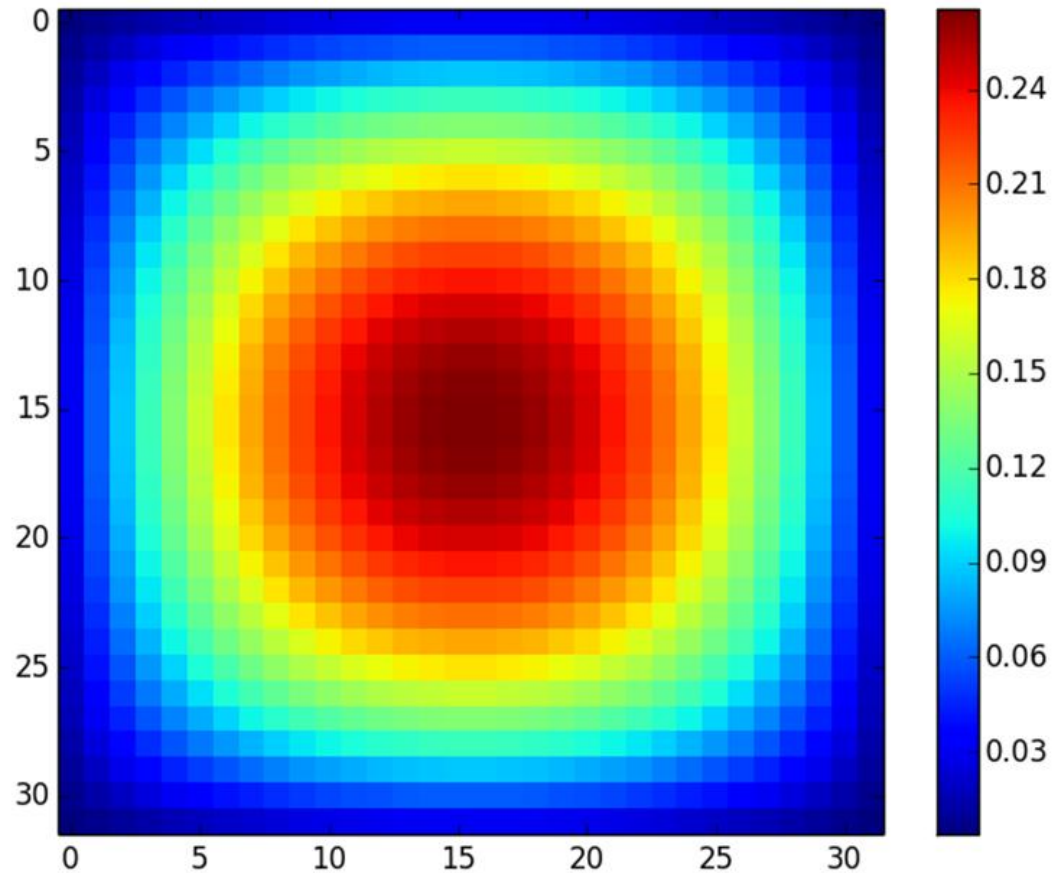
Spatial pooling by aggregation
2x2 averaging neighborhood



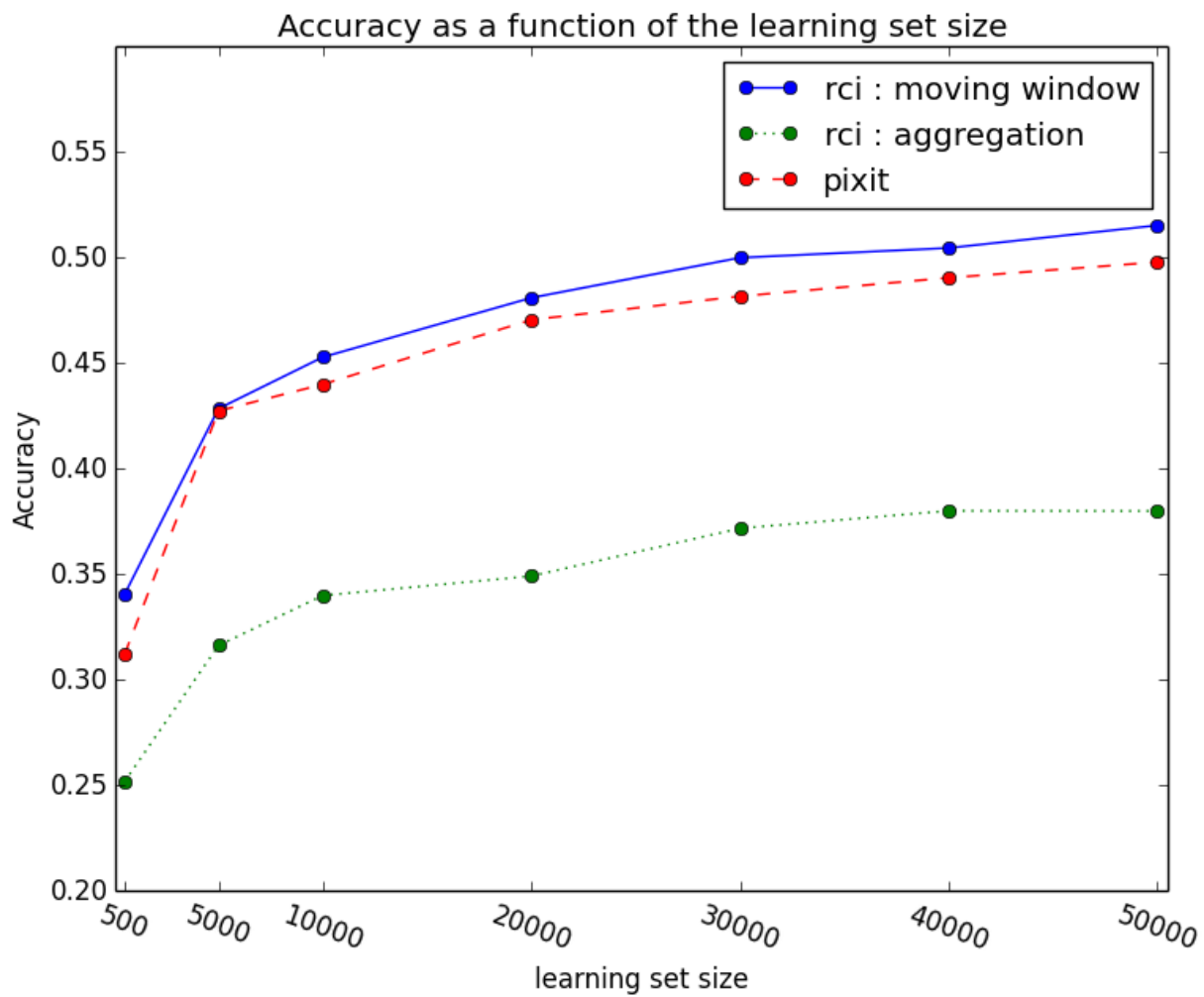
Spatial pooling by moving windows
3x3 averaging windows

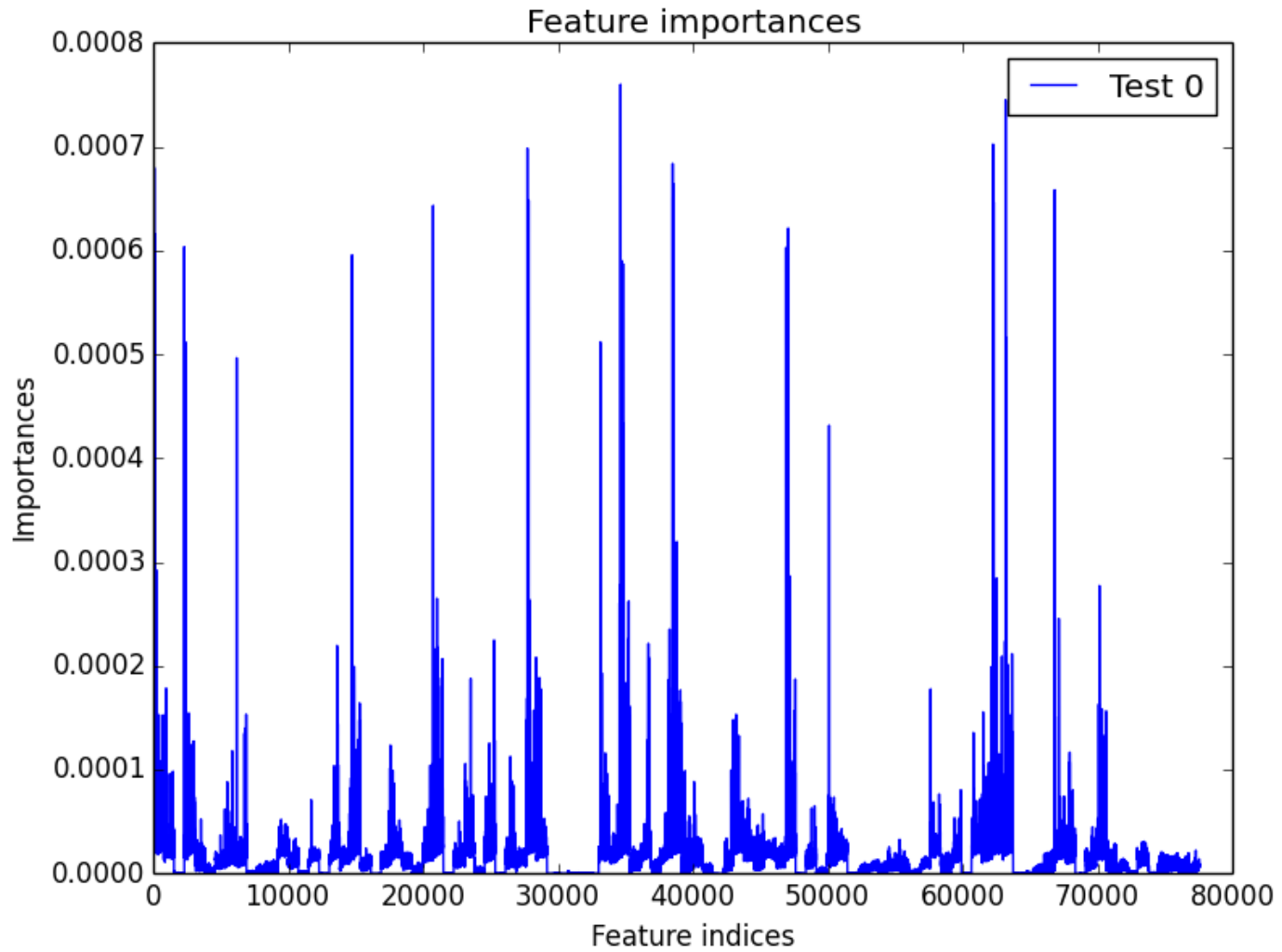


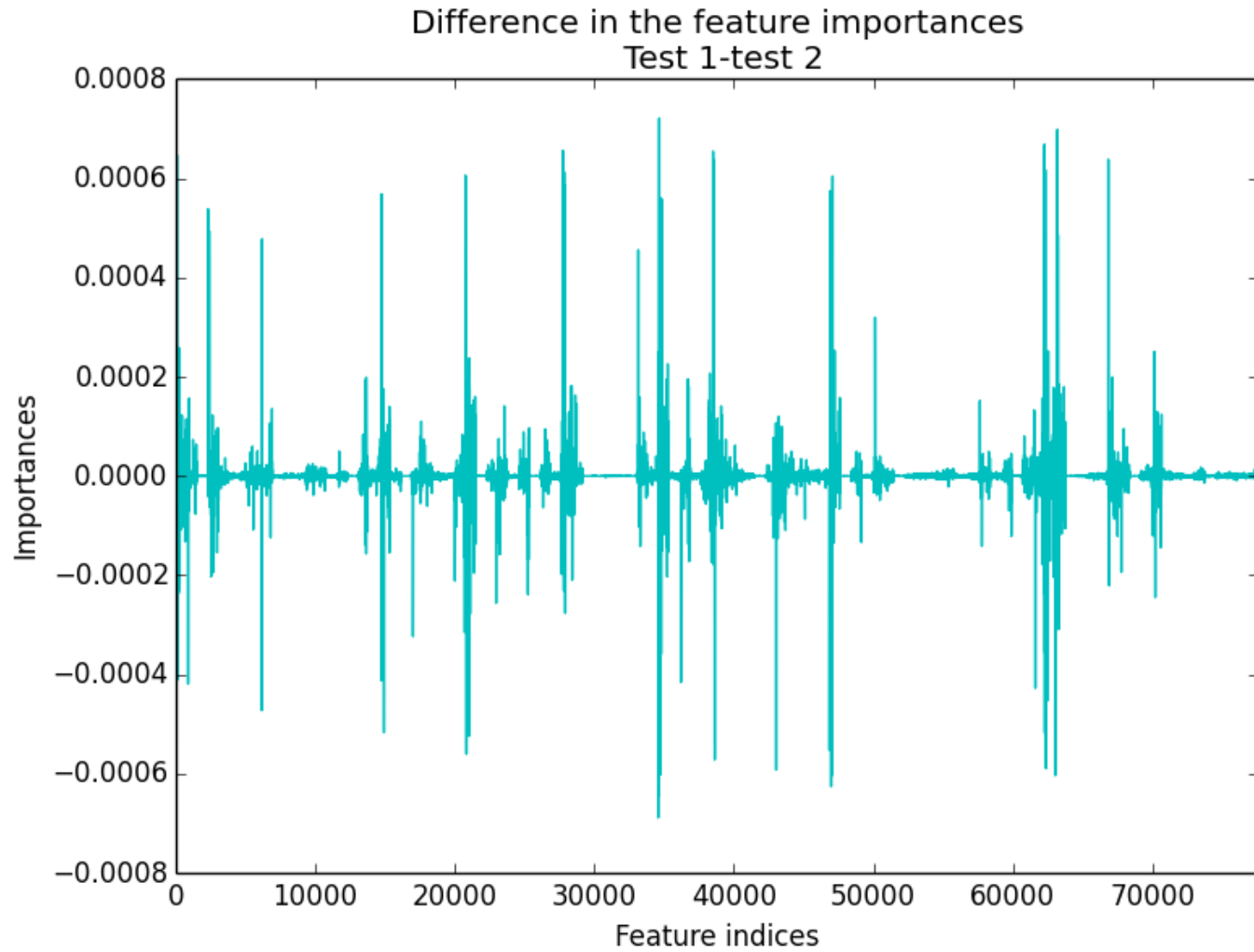
Probability per pixel of being included in a subwindow

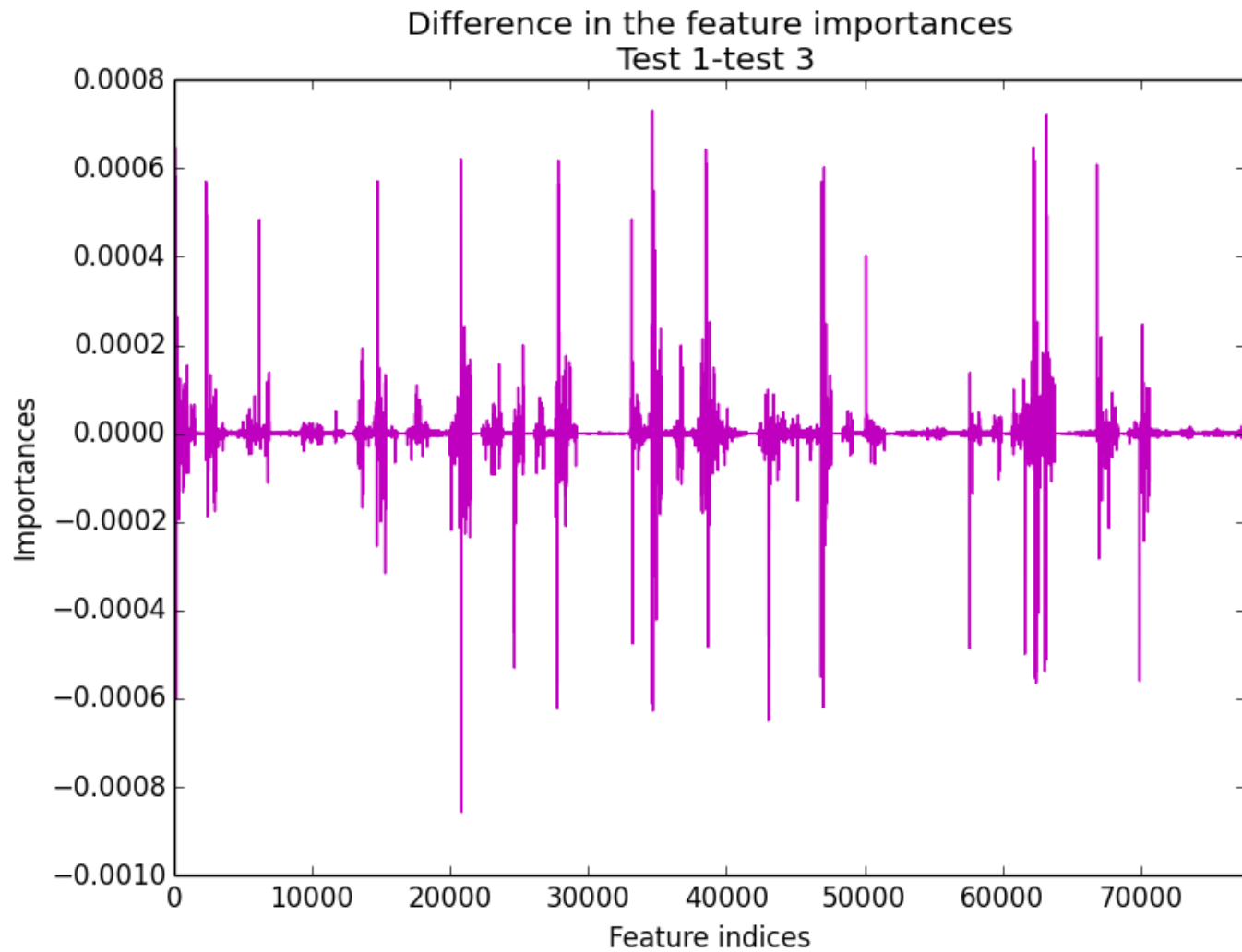


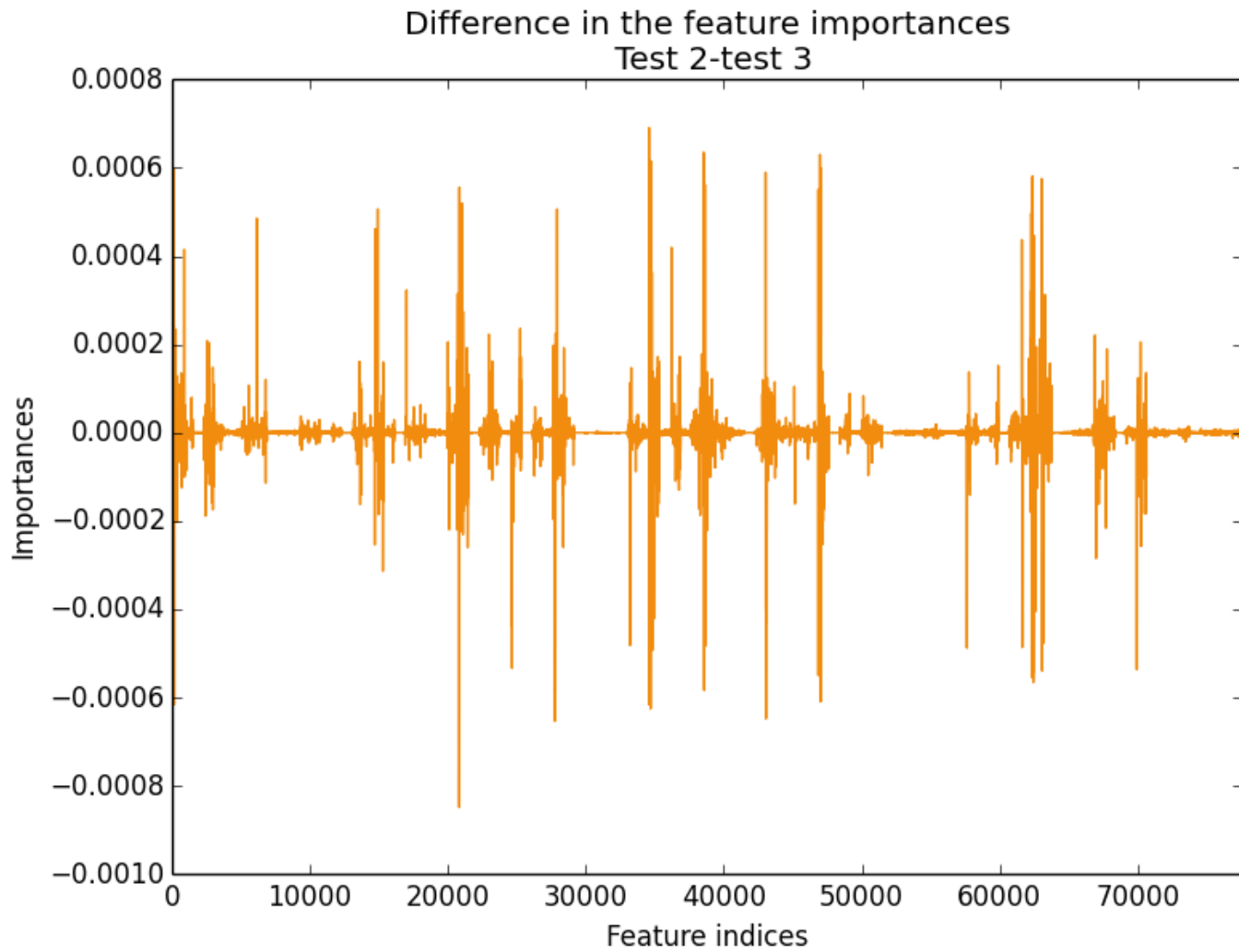
Graphes

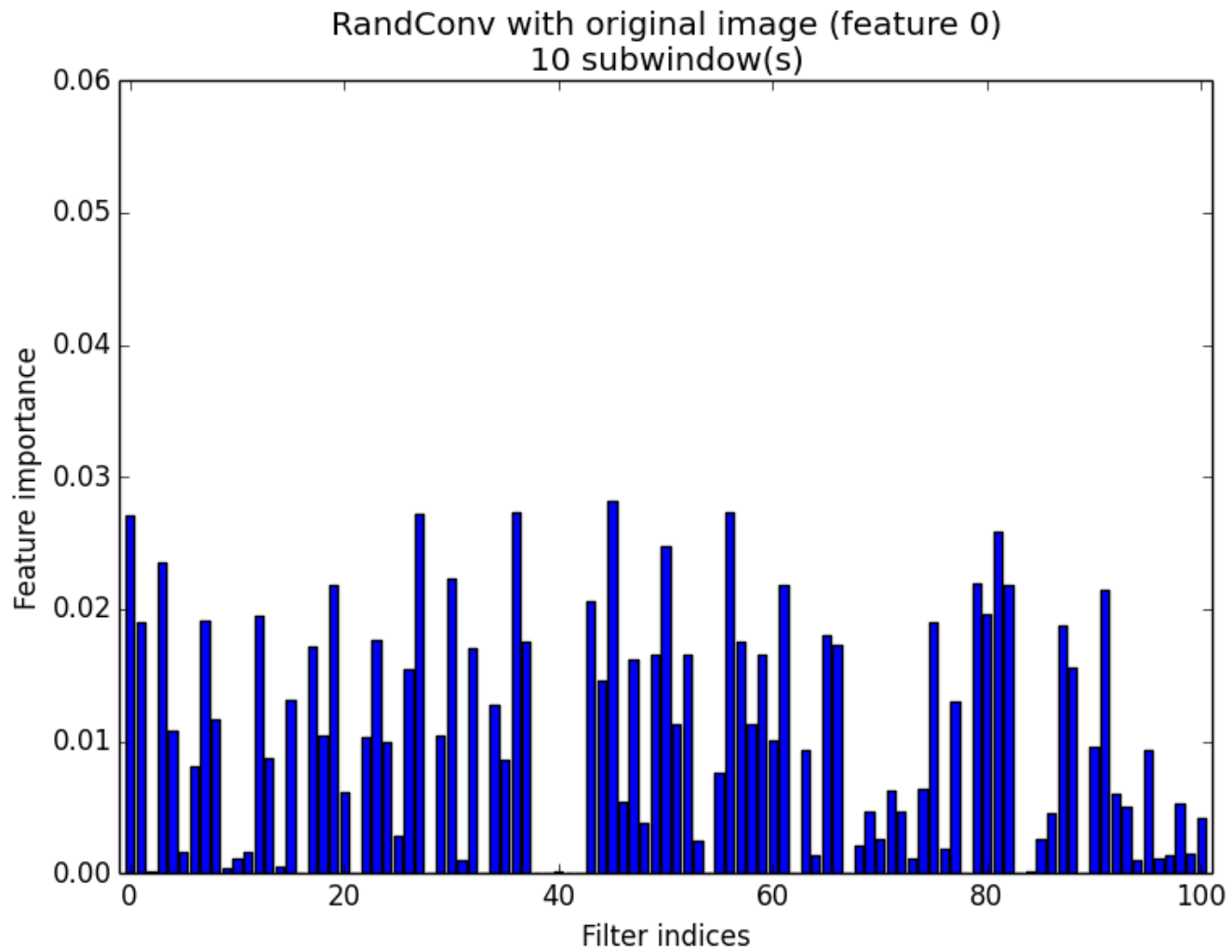


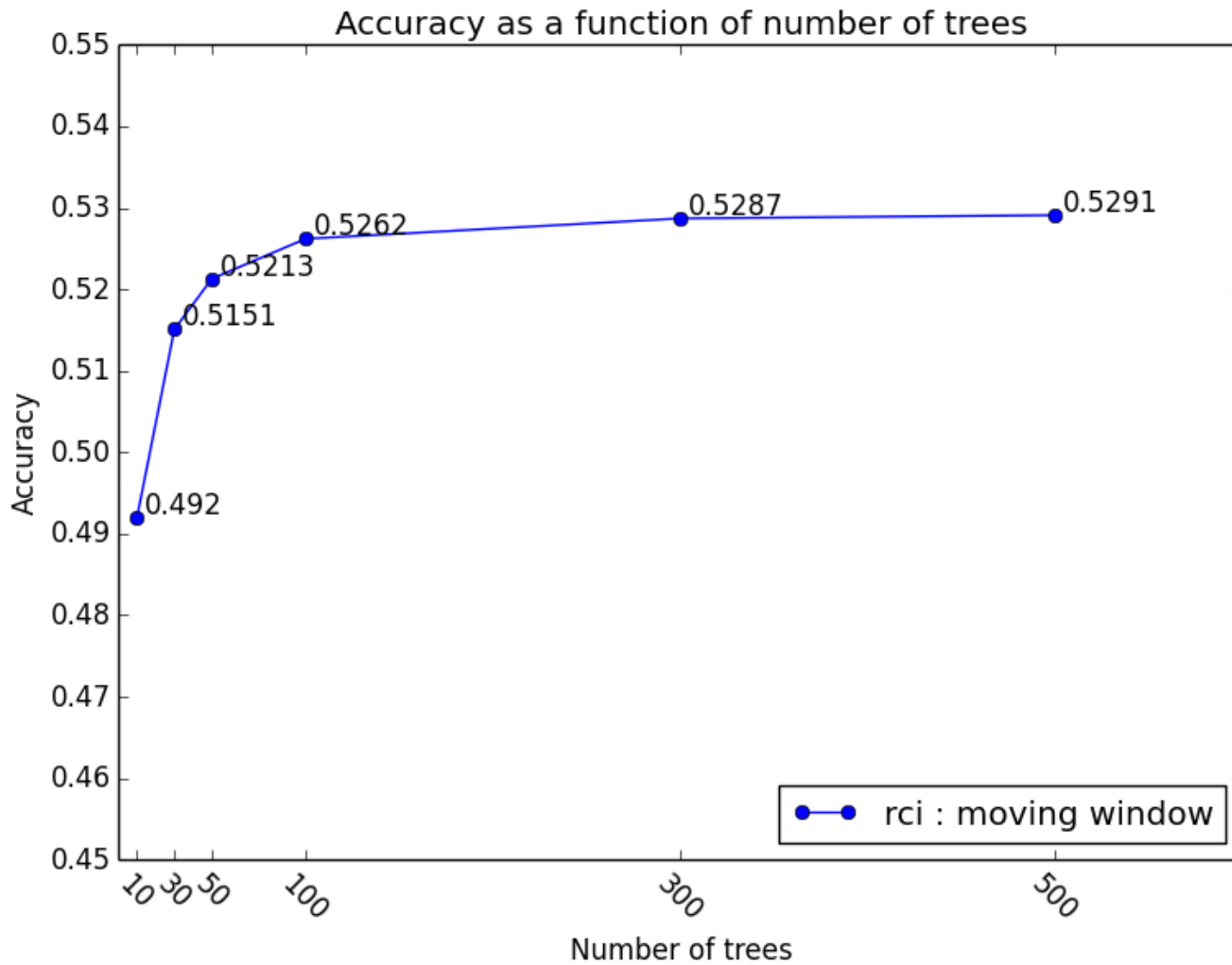




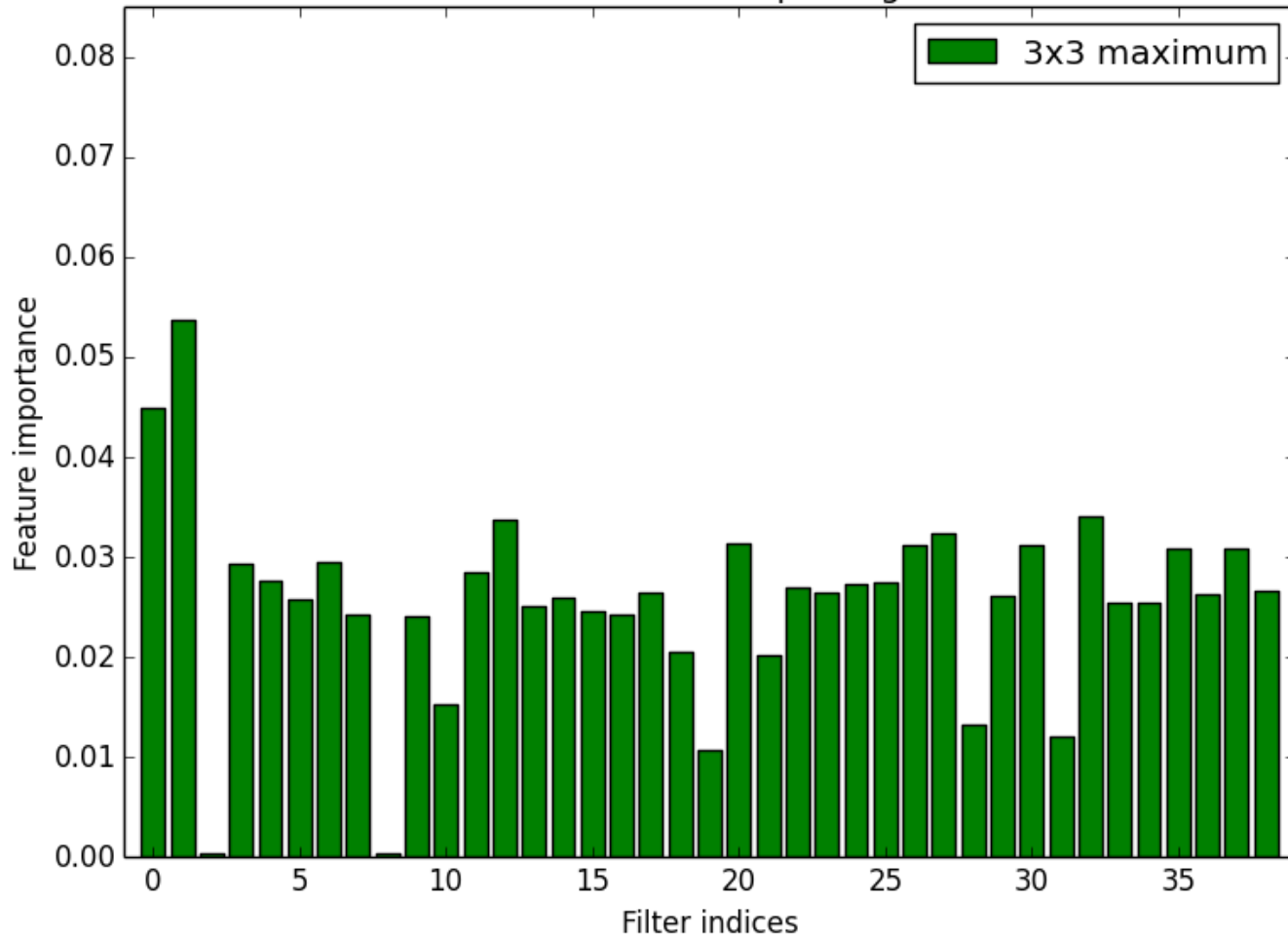




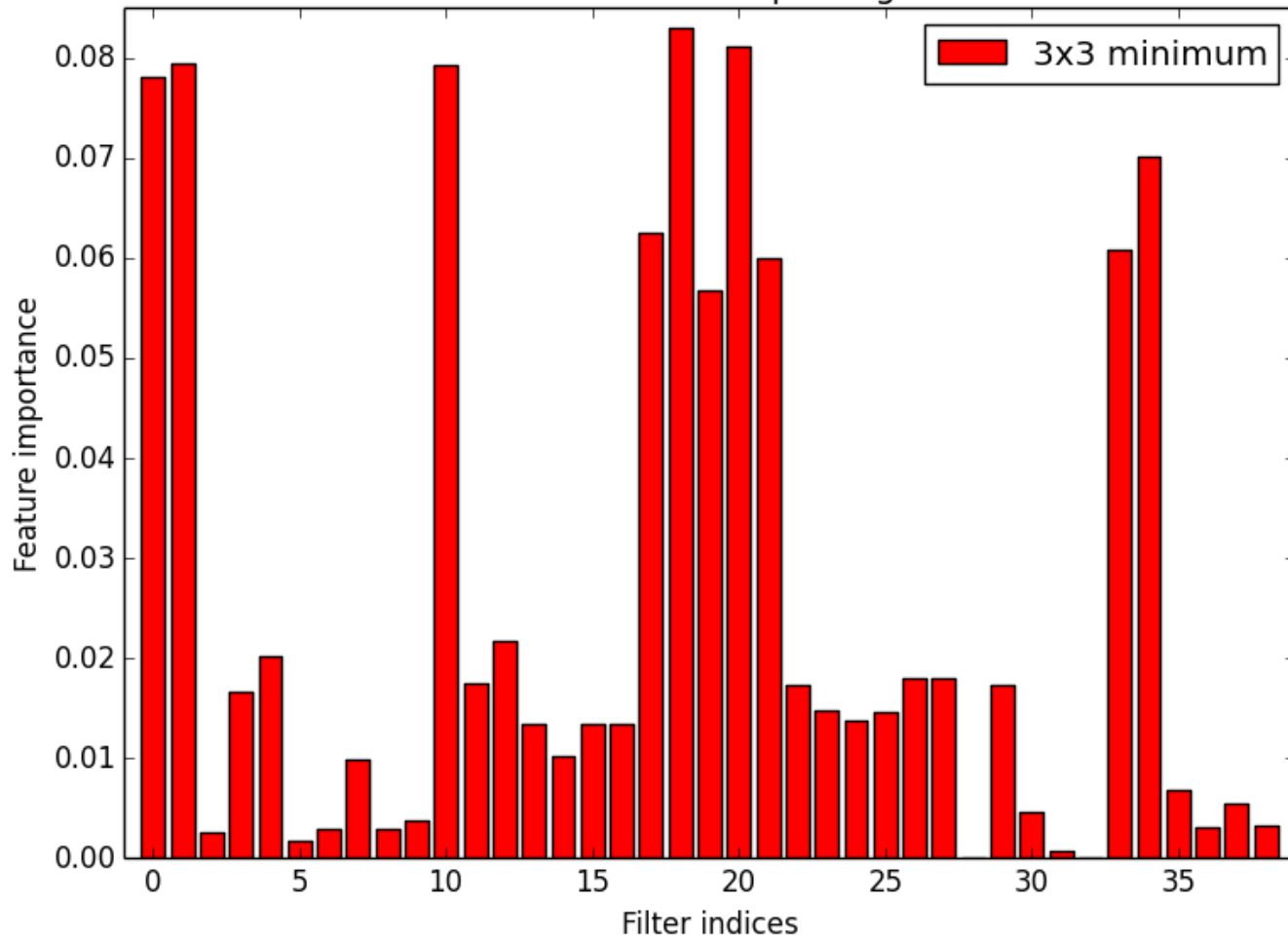




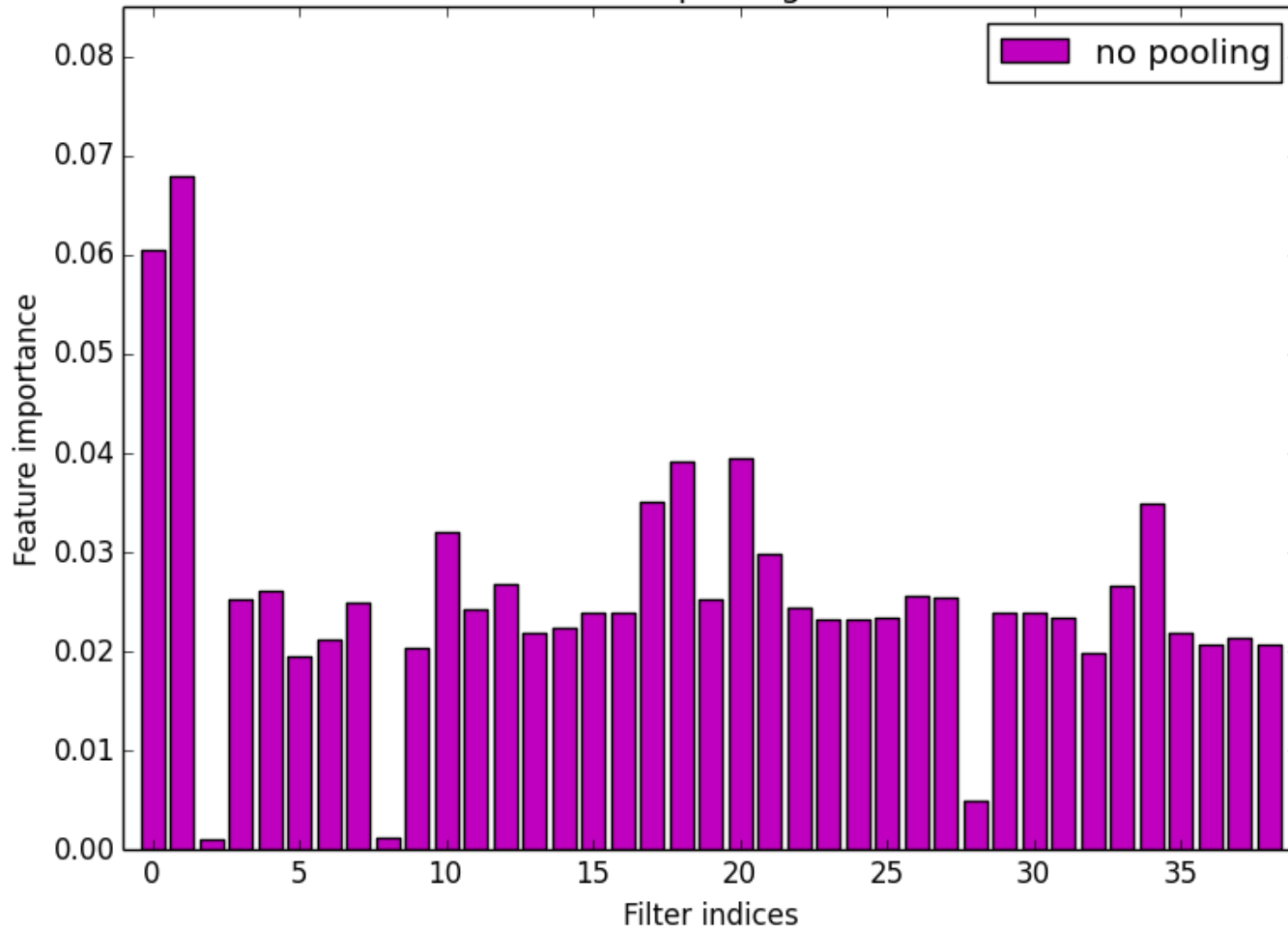
RandConv with original image (feature 0) and custom filters
3x3 maximum pooling

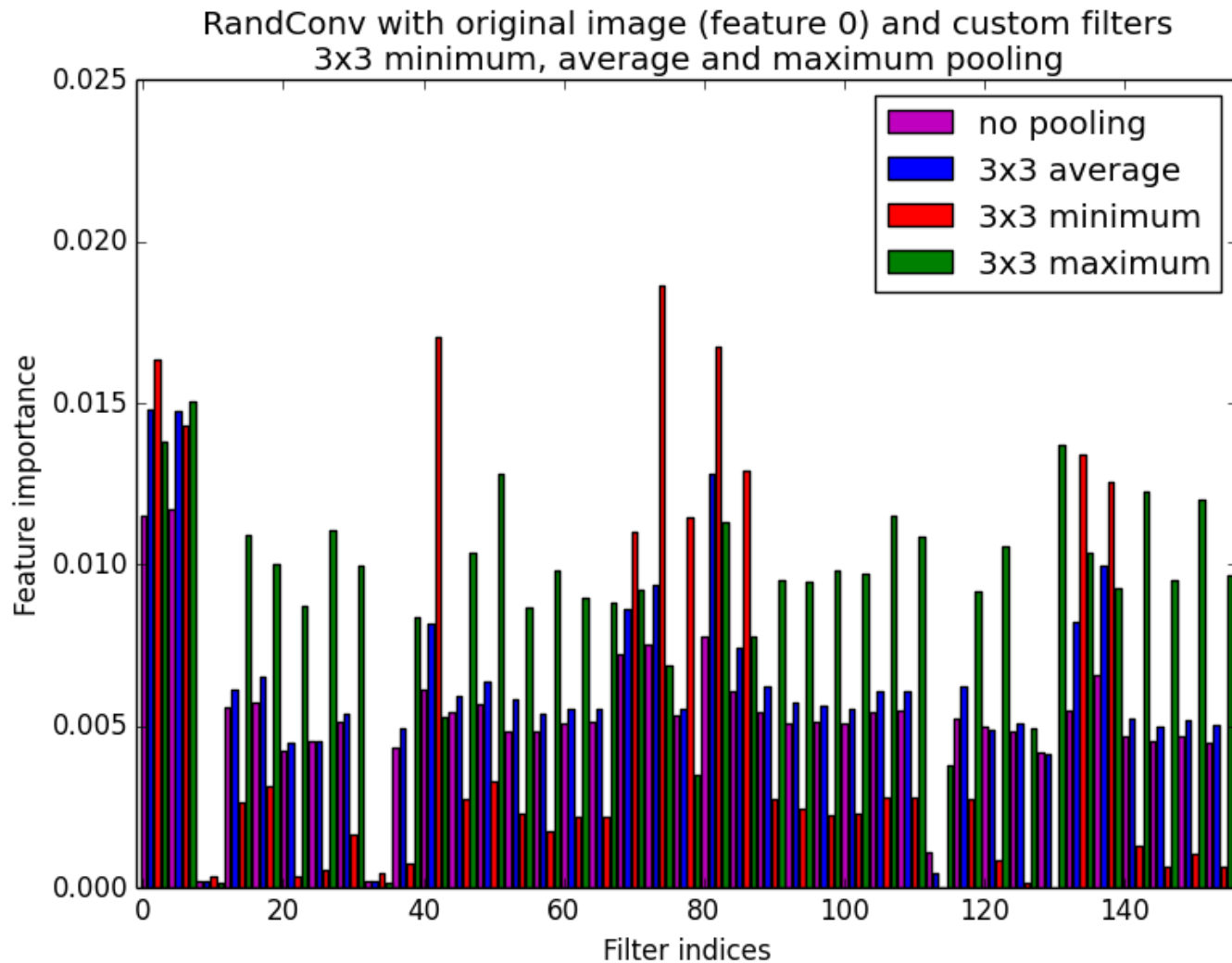


RandConv with original image (feature 0) and custom filters
3x3 minimum pooling

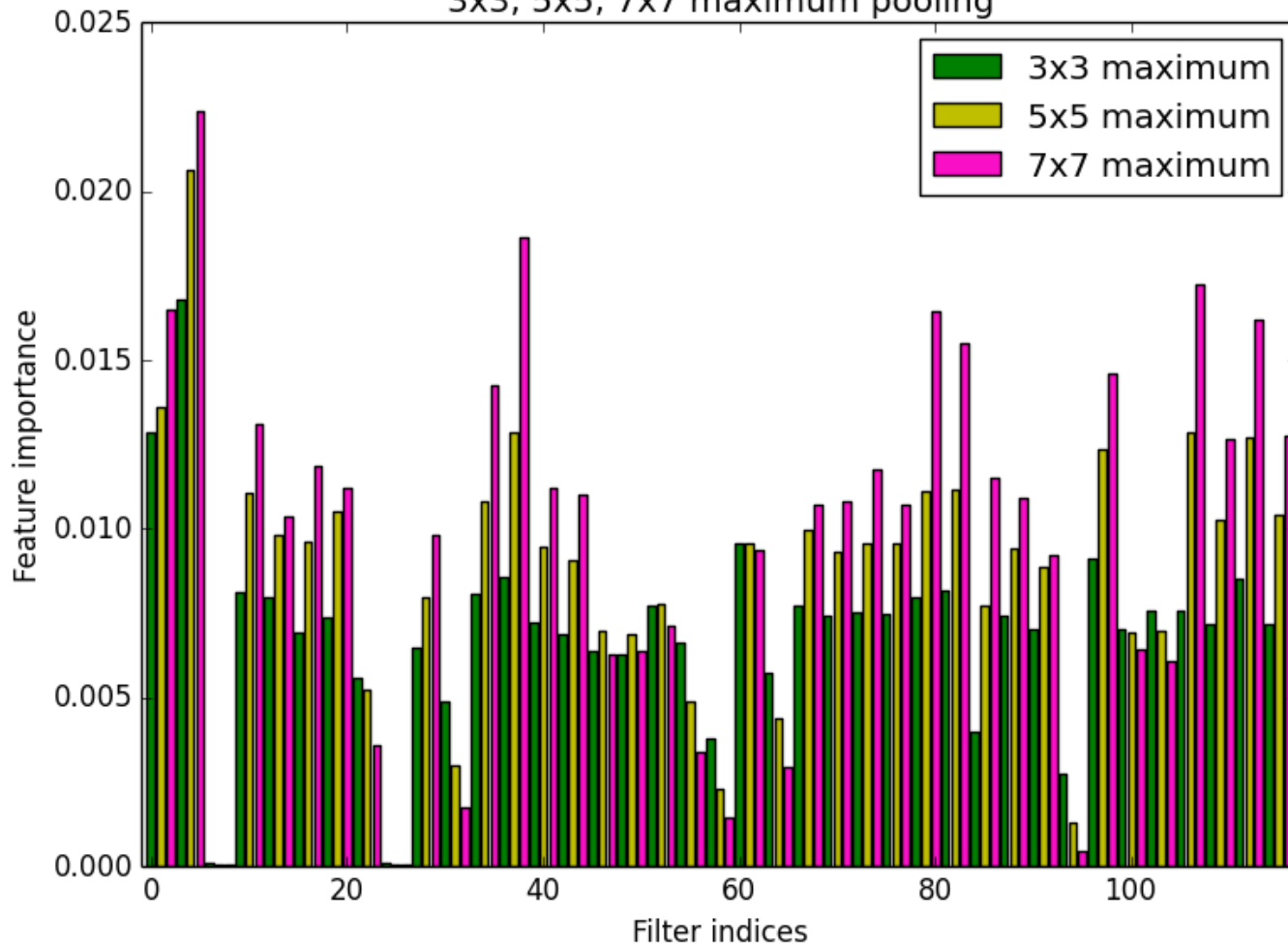


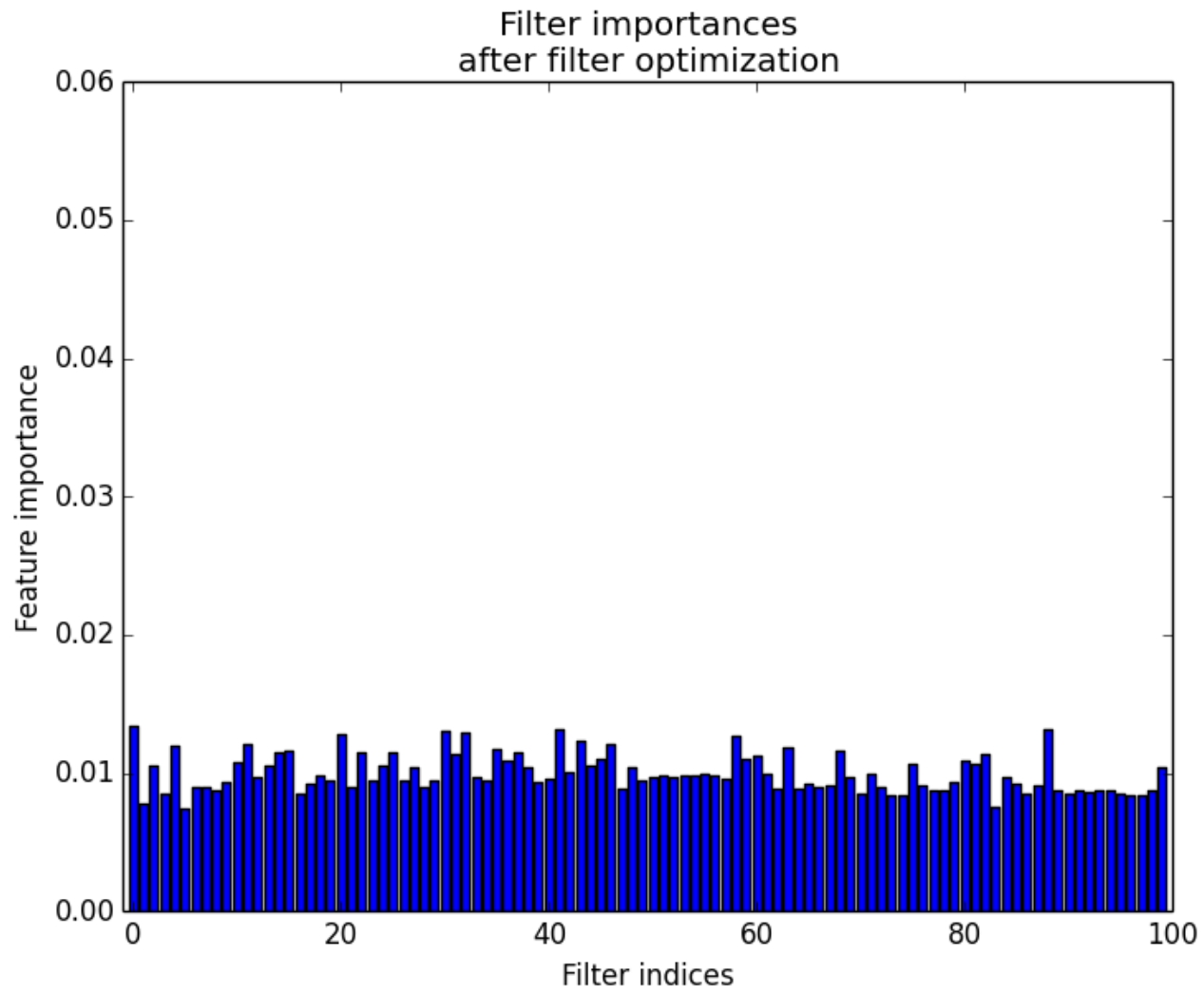
RandConv with original image (feature 0) and custom filters
No pooling



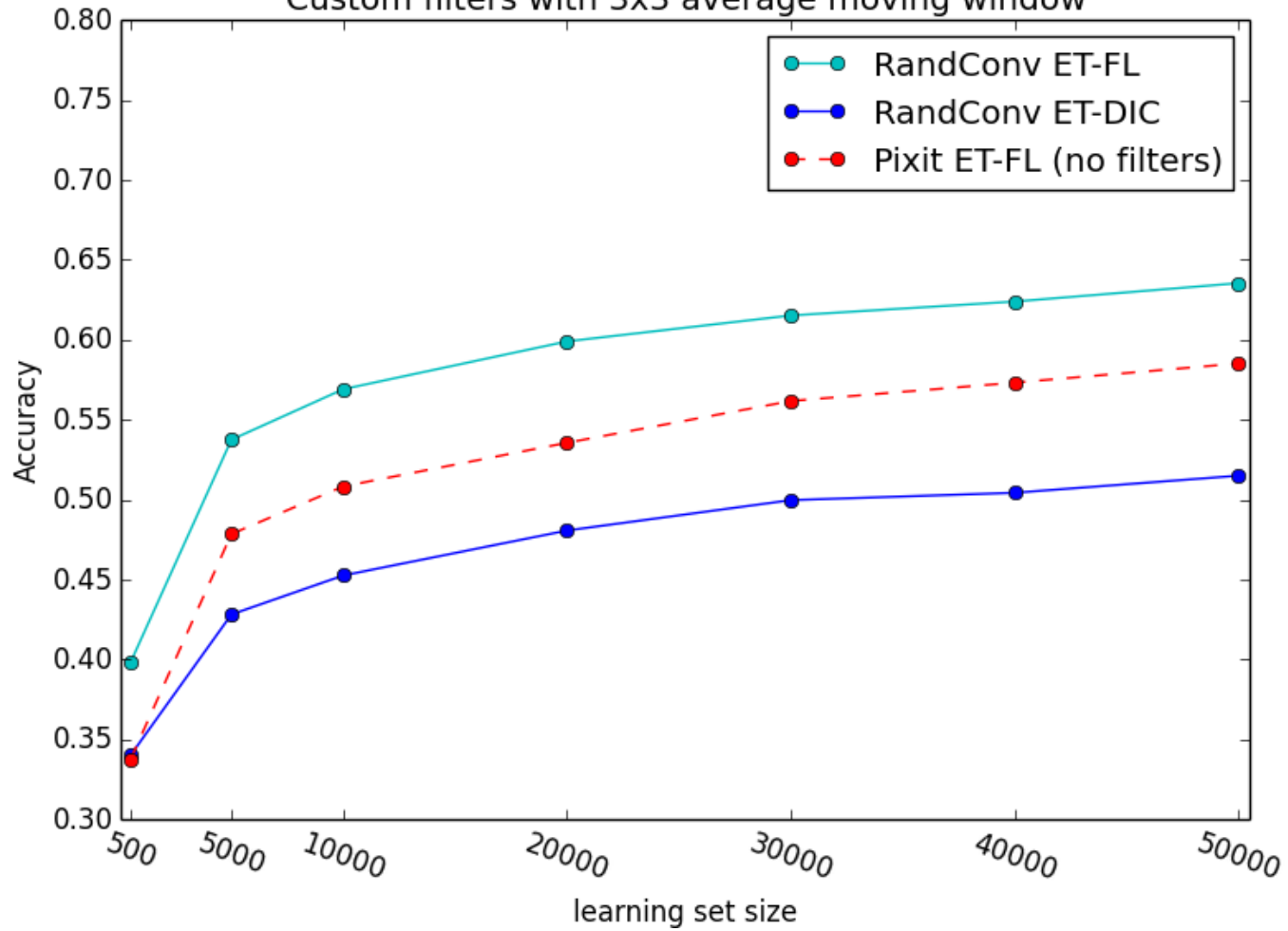


RandConv with original image (feature 0) and custom filters
3x3, 5x5, 7x7 maximum pooling

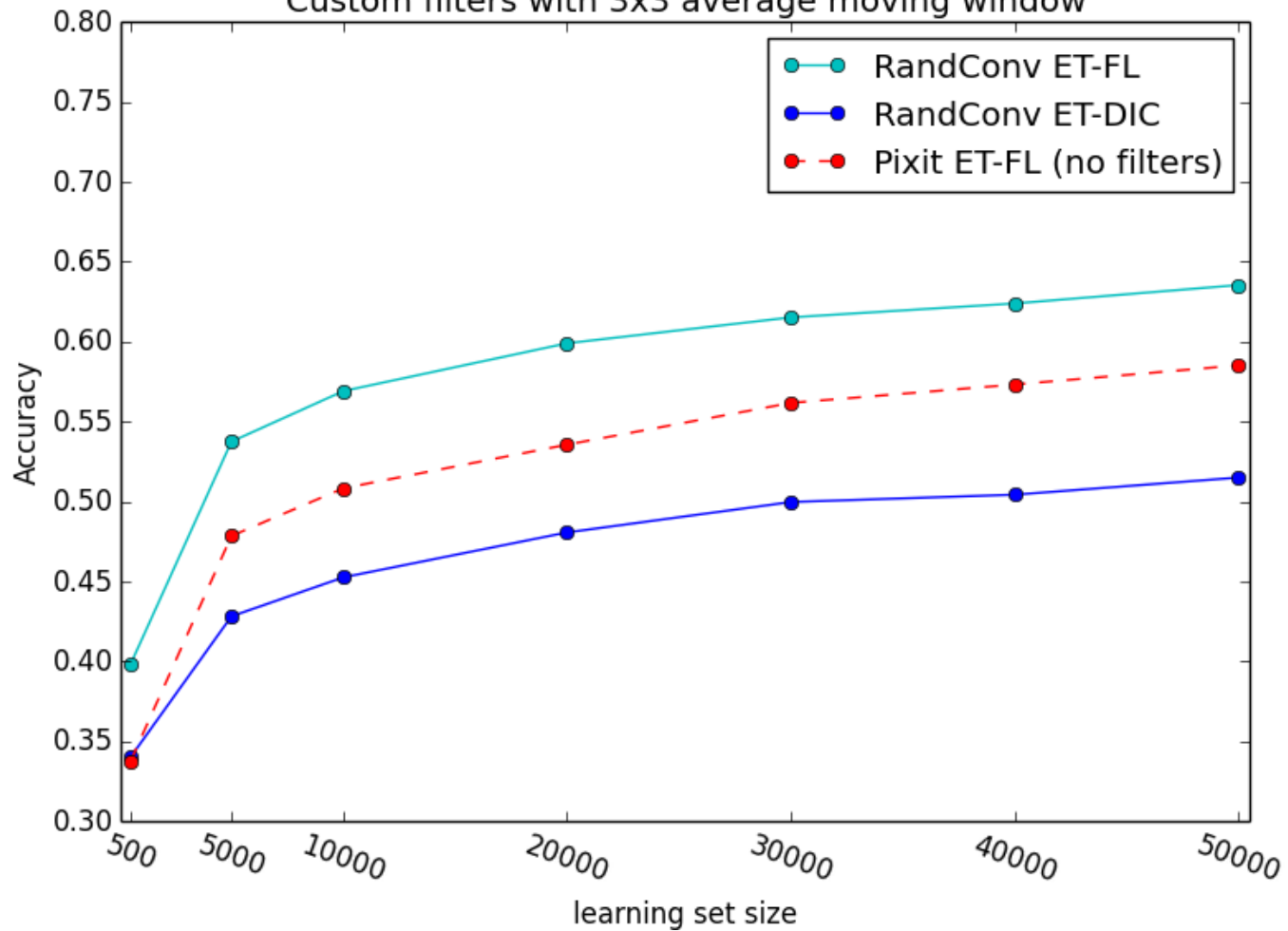


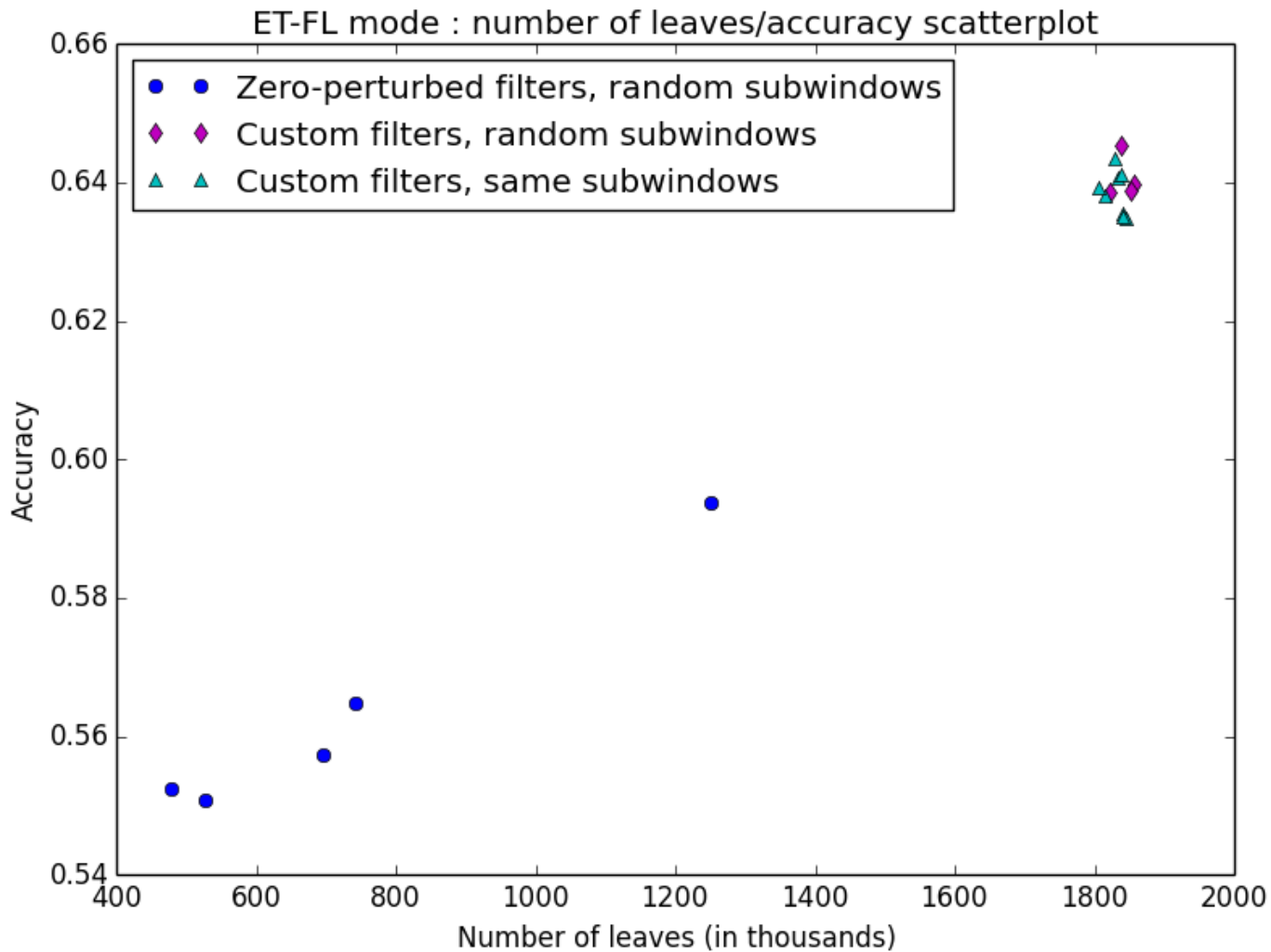


Accuracy comparison for different method
Custom filters with 3x3 average moving window



Accuracy comparison for different method
Custom filters with 3x3 average moving window





ET-FL mode : number of leaves/accuracy scatterplot

