

***Algorithmes parallèles pour  
la classification de grands  
ensembles de données***



*François Poulet*



*Passage à l'échelle dans la recherche d'information multimédia*

***SVM incrémental et  
parallèle sur GPU***



# *Introduction*

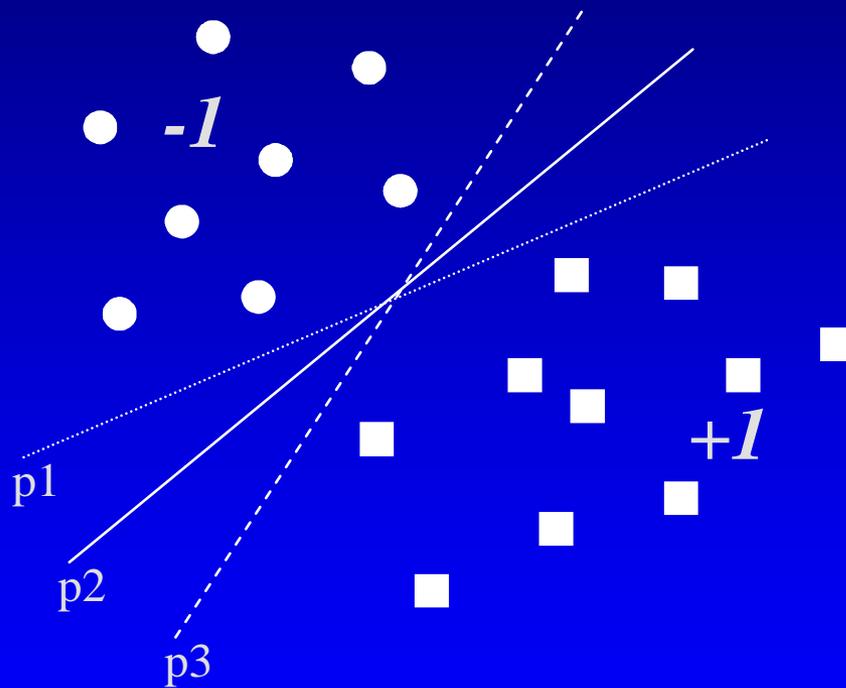
- algorithme de SVM (Vapnik, 1995)
  - hyperplan de séparation des données en 2 classes
  - SVM + fonction de noyau = méthode robuste
  - classification, régression, détection d'outlier
  - bons résultats
- programme quadratique : coûteux
- contribution
  - apprentissage incrémental et parallèle LS-SVM
  - fouille de très grands ensembles de données sur GPU



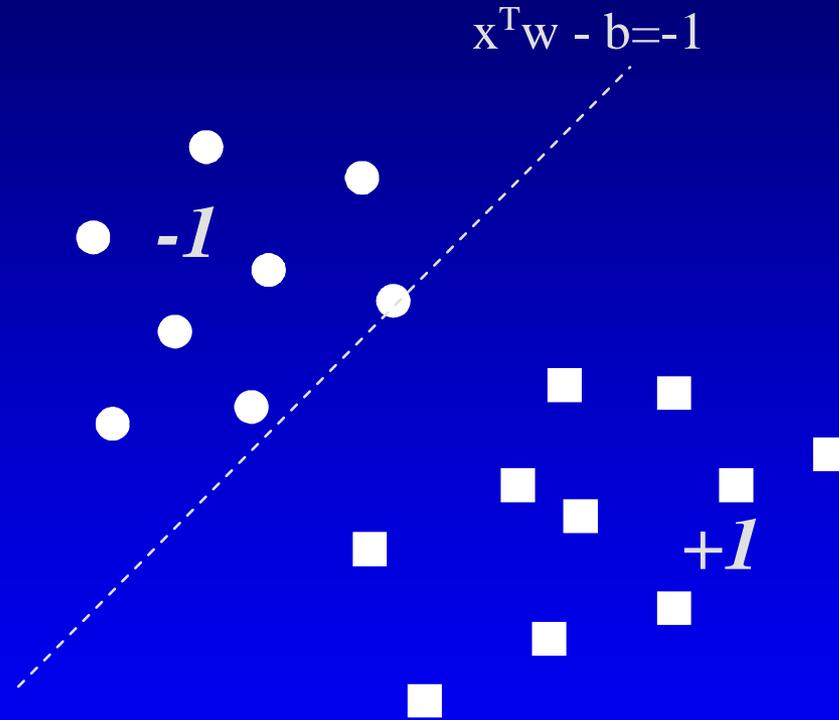
## *GPU ???*

- processeur spécialisé pour accélération graphique
- avantages / CPU pour calcul intensif :
  - meilleur taux de transfert mémoire
  - meilleurs possibilités en calculs flottants
  - centaines d'unités parallèles en SIMD
- alternative intéressante au cluster de CPUs
- GPUs récentes incluent possibilités calcul non graphiques : GP-GPU (General Purpose)

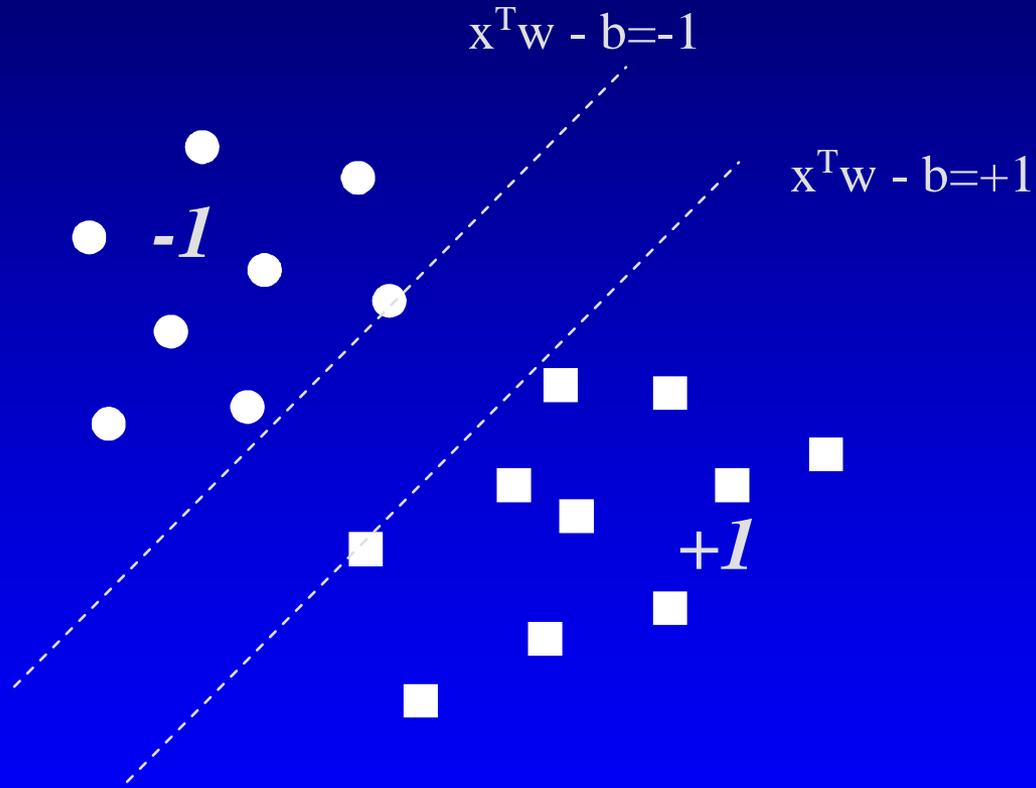
# Support vector machine



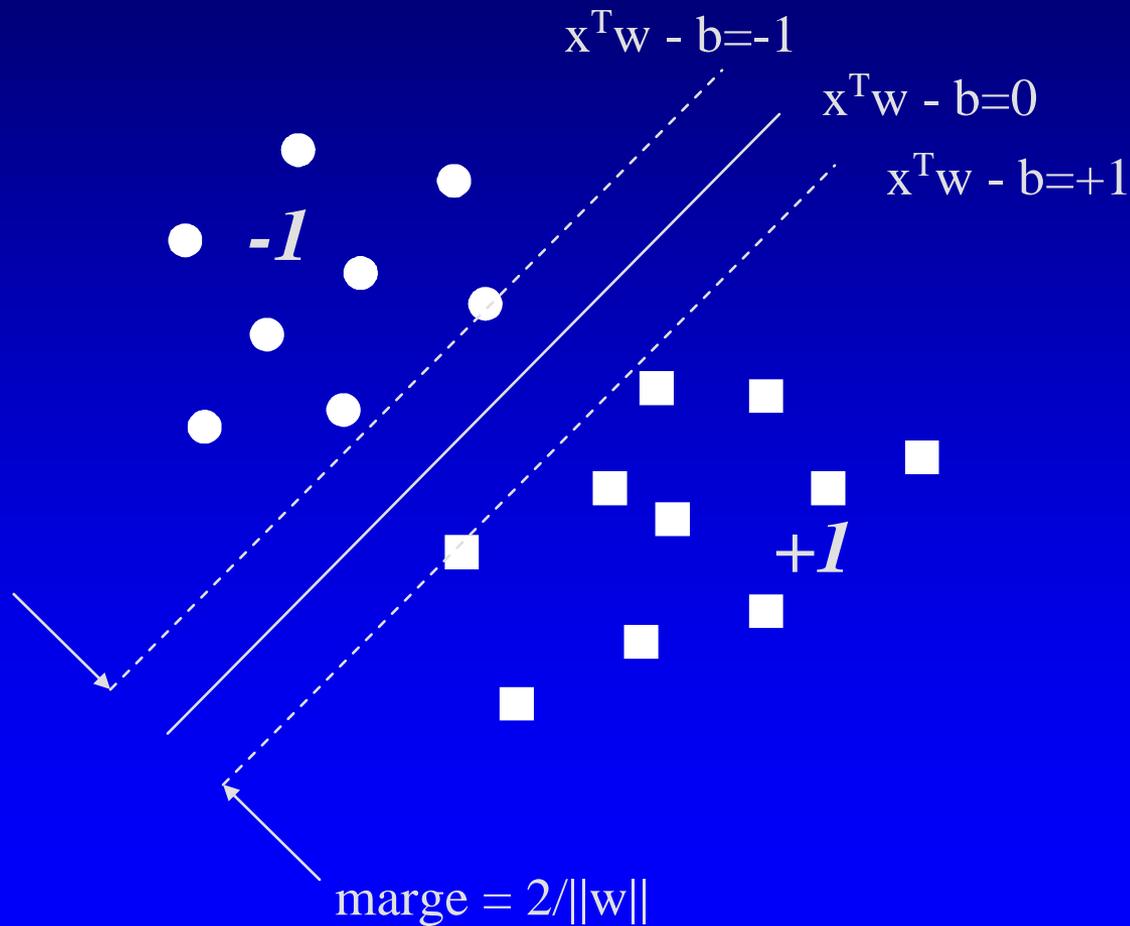
# Support vector machine



# Support vector machine

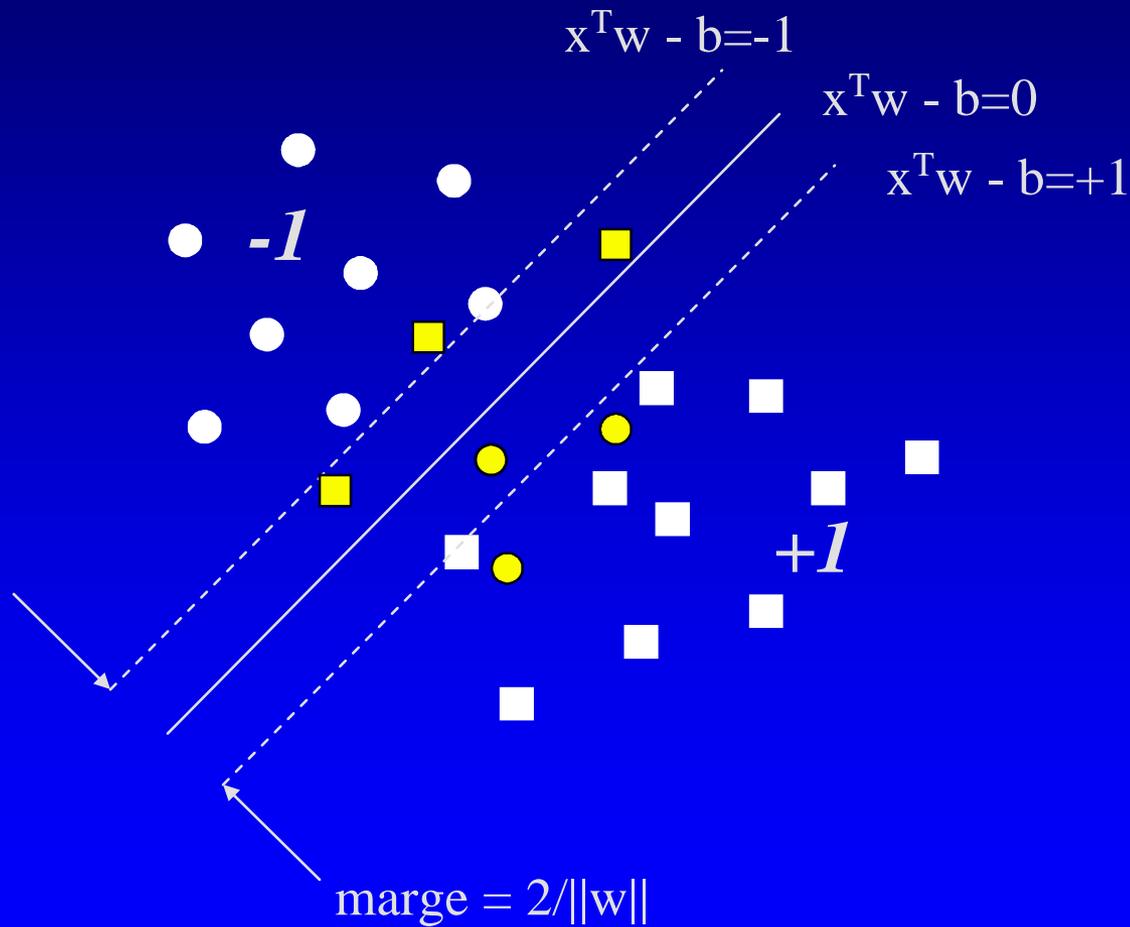


# Support vector machine





# Support vector machine





# Notations

- $A_{m \times n}$  : matrice représentant  $m$  ind. en  $n$  dim.
- $D_{m \times m}$  : matrice diagonale des classes ( $\pm 1$ ) des  $m$  ind.
- $e$  : vecteur colonne de 1
- $X^T$  : transposée de  $X$
- $w, b$  : plan
- $\|w\|$  : norme-2 du vecteur  $w$



## *Algorithme de SVM*

- maximisation marge + minimisation erreurs
  - $\min f(z,w,b) = (c/2) \|z\| + (1/2) \|w\|^2$
  - avec  $D(Aw - eb) + z \geq e$  (1)
  - où  $z \geq 0$  variable de ressort et  $c$  constante positive
- résolution du programme quadratique (1) :  $w, b$
- classification d'un nouvel individu  $x$  :  $\text{signe}(x^T w - b)$



## *LS-SVM (Suykens et al., 1999)*

- Least Squares SVM
  - remplace inégalité par égalité :  $D(Aw - eb) + z = e$
  - minimisation des erreurs :  $\min (1/2) \|z\|^2$
- LS-SVM : système linéaire à  $(n+1)$  inconnues  $(w, b)$ 
  - au lieu du programme quadratique



## *Apprentissage incrémental*

- grands ensembles de données
  - limite : mémoire => apprentissage incrémental
  - petit bloc de données en mémoire à la fois
  - même solution
- nombre individus  $\ll$  nombre dimensions
  - incrémental en colonne
- nombre individus  $\gg$  nombre dimensions
  - incrémental en ligne



## *Apprentissage incrémental*

- l'équation à résoudre est :
- $(w_1 \ w_2 \ \dots \ w_n \ b)^T = (I^{\circ}/c + E^T E)^{-1} E^T D e$
- complexité spatiale :
  - $E^T E$  et  $E^T D e$  :  $[n+1 \times n+1]$  et  $[n+1 \times 1]$
  - $E$ ,  $E^T$  et  $D$  :  $m \times (n+1)$ ,  $(n+1) \times m$  et  $m \times m$
  - carré du nombre de dimensions et du nombre d'individus
- impossible traiter grandes quantités données !

# Apprentissage incrémental et parallèle

- on découpe en blocs :

- $E^T E = \sum E_i^T E_i$  et  $E^T D e = \sum E_i^T D_i e_i$

- pour obtenir finalement :

$$[w_1 w_2 w_3 \dots w_n b]^T = \left( \frac{I^o}{c} + \sum_{i=1}^k E_i^T E_i \right)^{-1} \sum_{i=1}^k E_i^T D_i e_i$$

- complexité spatiale :

- $(n+1) \times b_i$  et  $b_i \times b_i$ ,  $b_i =$  taille bloc lignes x nombre dimensions

- complexité en temps :

- distribution des calculs de blocs sur différents processeurs



## *Déroulement de l'algorithme*

- pour  $i$  de 1 à  $k$ 
  - charge blocs de données  $E_i$  et  $D_i$  en mémoire CPU
  - copie  $E_i$  et  $D_i$  de CPU vers GPU
  - calcule  $E^T E = E^T E + E_i^T E_i$  et  $E^T D e = E^T D e + E_i^T D_i e_i$  sur GPU
- fin pour
  
- copie  $E^T E$  et  $E^T D e$  vers CPU
- résolution équation finale





## *Environnement de test*

- NVidia GeForce 8800 GTX :
  - 16 multiprocesseurs
  - chaque multi-processeur a 8 processeurs (total=128)
  - chaque processeur a un cache L1 et ALU
  - 768 Mo mémoire
  - 330 GFlops
  - 86Go/s bande passante mémoire
  - transfert CPU/GPU : 2Go/s



## *Environnement de test*

- NVidia GTX - 280 :
  - 3 blocs de 10 multiprocesseurs
  - chaque multiprocesseur a 8 processeurs :
  - 240 processeurs (128)
  - 1Go mémoire (768 Mo)
  - bande passante mémoire 142Go/s (86Go/s)
  - 933 GFlops (330 GFlops)
- nVidia CUDA, API en langage C pour GP-GPU
- cublas : implémentation CUDA de BLAS :  
Basic Linear Algebra Subprograms



## *Ensembles de données*

Dataset	# dimensions	Training set	Test set
Adult	110	32,561	16,281
Forest Cover Types	20	495,141	45,141
KDD Cup 1999	41	4,868,429	311,029
Ringnorm-1M	20	1,000,000	100,000
Ringnorm-10M	20	10,000,000	1,000,000

- attributs nominaux -> binaires
- 2 plus grandes classes Forest Cover Type (pour se comparer aux autres, idem pour le protocole de test)



# Résultats

Dataset	Time (s)		ratio		Accuracy (%)
	GPU	CPU	CPU/GPU	SVM/GPU	
Adult	0.007	2.00	285	18150 (libSVM)	85.08
Forest Cover Types	0.06	10.00	166,67	2850 (svm-perf)	76.41
KDD Cup 1999	0.58	55.67	95.98	8189 (CB-svm)	91.96
Ringnorm-1M	0.06	3.33	55.50	N/A	75.07
Ringnorm-10M	0.61	31.67	51.92	N/A	76.68

- en moyenne GPU 130 fois plus rapide que même algorithme en version CPU !!!
- comparaison algorithmes habituels :
  - 10000 fois plus rapide que libSVM, svm-perf, CB-svm



## *Conclusion*

- algorithme de SVM incrémental :
  - résoudre problème de mémoire
- algorithme de SVM parallèle :
  - résoudre problème de temps d'exécution
- moyenne :
  - 130 fois plus rapide que CPU
  - + de 3000 fois plus rapide algorithmes habituels



## *Perspectives*

- utilisé une seule GPU
- $n$  GPUs...
  - 10 GPUs : 30000 fois plus rapide que les algorithmes standards
  - 1 an de calcul  $\rightarrow$  1/2 heure de calcul
  - aucune optimisation sur GPU...
- nouvelles applications potentielles en fouille de données...



# *RFO*

- Random Forest (Breiman, 2001)
- construire n arbres de décision sur un sous-ensemble de dimensions / individus
- traiter grands nb dim / ind
- remplace chaque coupe arbre de décision par SVM linéaire, "arbre oblique"
- utilise un cluster de PCs (1 arbre oblique / cœur)



# *RFO*

- comparaison avec :
  - libSVM (noyau RBF, tuning paramètres)
  - Random Forest de C4.5
- 15 ensembles données Kent Ridge Bio Medical Dataset Repository (grand nombre de dimensions)



Dataset	#pts	#dim	Classes	Test
Colon Tumor	62	2000	tumor, normal	loo
ALL-AML-Leukemia	72	7129	ALL, AML	trn-tst
*MLL-Leukemia	72	12582	MLL, rest	trn-tst
Breast Cancer	97	24481	relapse, non-relapse	trn-tst
Duke Breast Cancer	42	7129	cancer, normal	loo
Prostate Cancer	136	12600	cancer, normal	trn-tst
Lung Cancer	181	12533	cancer, normal	trn-tst
Central Nervous System	60	7129	pos, neg	loo
Translation Initiation Site	13375	927	pos, neg	10-fold
Ovarian Cancer	253	15154	cancer, normal	loo
Diffuse Large B-Cell Lymphoma	47	4026	germinal, activated	loo
*Subtypes of Acute Lymphoblastic (Hyperdip)	327	12558	Hyperdip, rest	trn-tst
*Subtypes of Acute Lymphoblastic (TEL-AML1)	327	12558	TEL-AML1, rest	trn-tst
*Subtypes of Acute Lymphoblastic (TEL-ALL)	327	12558	TEL-ALL, rest	trn-tst
*Subtypes of Acute Lymphoblastic (Others)	327	12558	others, diagnostic groups	trn-tst


**RFO**

Dataset	precision			recall		
	lib-SVM	RF-C4.5	RF-ODT	lib-SVM	RF-C4.5	RF-ODT
Colon Tumor	68.18	76.19	<b>82.61</b>	75.00	72.73	<b>86.36</b>
ALL-AML-Leukemia	<b>100</b>	95.24	95.24	95.00	<b>100</b>	<b>100</b>
*MLL-Leukemia	75.00	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>
Breast Cancer	69.23	83.33	<b>84.62</b>	75.00	83.33	<b>91.67</b>
Duke Breast Cancer	85.00	<b>94.12</b>	90.00	<b>94.44</b>	80.00	90.00
Prostate Cancer	73.53	75.76	<b>100</b>	<b>100</b>	<b>100</b>	96.00
Lung Cancer	88.26	<b>93.75</b>	<b>93.75</b>	<b>100</b>	<b>100</b>	<b>100</b>
Central Nervous System	47.62	45.46	<b>61.91</b>	55.56	23.81	<b>61.91</b>
Translation Initiation Site	83.13	<b>92.58</b>	90.78	<b>84.42</b>	73.83	79.75
Ovarian Cancer	<b>100</b>	98.78	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>
Diffuse Large B-Cell Lymphoma	91.30	<b>95.65</b>	92.00	87.50	91.67	<b>95.83</b>
*Subtypes (Hyperdip)	95.46	95.24	<b>100</b>	<b>95.46</b>	90.91	<b>95.46</b>
*Subtypes (TEL-AML1)	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	96.30	96.30
*Subtypes (TEL-ALL)	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>
*Subtypes (Others)	92.59	<b>100</b>	<b>100</b>	39.68	29.63	<b>55.56</b>

Dataset	F1 Measure			accuracy		
	lib-SVM	RF-C4.5	RF-ODT	lib-SVM	RF-C4.5	RF-ODT
Colon Tumor	71.43	74.42	<b>84.44</b>	80.65	82.26	<b>88.71</b>
ALL-AML-Leukemia	97.44	<b>97.56</b>	<b>97.56</b>	<b>97.06</b>	<b>97.06</b>	<b>97.06</b>
*MLL-Leukemia	85.71	<b>100</b>	<b>100</b>	93.33	<b>100</b>	<b>100</b>
Breast Cancer	72.00	83.33	<b>88.00</b>	63.16	78.94	<b>84.21</b>
Duke Breast Cancer	89.47	86.49	<b>90.00</b>	90.48	88.10	<b>90.48</b>
Prostate Cancer	84.75	86.21	<b>97.96</b>	73.53	76.47	<b>97.06</b>
Lung Cancer	93.75	<b>96.77</b>	<b>96.77</b>	98.66	<b>99.33</b>	<b>99.33</b>
Central Nervous System	51.28	31.25	<b>61.91</b>	68.33	63.33	<b>73.33</b>
Translation Initiation Site	83.77	82.15	<b>84.91</b>	92.15	92.30	<b>93.20</b>
Ovarian Cancer	<b>100</b>	99.39	<b>100</b>	<b>100</b>	99.21	<b>100</b>
Diffuse Large B-Cell Lymphoma	89.36	93.62	<b>93.88</b>	89.36	93.62	<b>93.62</b>
*Subtypes (Hyperdip)	95.46	93.02	<b>97.67</b>	98.21	97.32	<b>99.11</b>
*Subtypes (TEL-AML1)	<b>100</b>	98.11	98.11	<b>100</b>	99.11	99.11
*Subtypes (TEL-ALL)	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>	<b>100</b>
*Subtypes (Others)	55.56	45.71	<b>71.43</b>	64.29	83.93	<b>89.29</b>



## *RFO - Conclusion*

- mesure F1 et précision : RFO meilleur partout sauf dans un cas (libSVM)
- plus gros gain (/libSVM) dans le cas de Breast Cancer +20% (le + gd nb de dimensions = 24481)
- coût : 10 000 mots visuels 6500 images => 75 jours cpu (100 cœurs -> - d'1 jour cpu)



## Références

- F.Poulet, T-N.Do, V-H.Nguyen : *SVM incrémental et parallèle sur GPU*, in Revue des Nouvelles Technologies de l'Information, Série Extraction et Gestion des Connaissances, RNTI-E-15, Cépaduès, 2009, 103-114.
- T-N. Do, V-H. Nguyen, F. Poulet: *Speed up SVM Algorithm for Massive Classification Tasks* in proc. of ADMA'08, the 4th International Conference on Advanced Data Mining and Applications, 2008, Chengdu, China, in Advanced Data Mining and Applications, Lecture Notes in Computer Science, LNCS 5139, 147-157.
- T-N. Do, V-H. Nguyen, F. Poulet: *A Fast Parallel SVM Algorithm for Massive Classification Tasks* in proc. of MCO'08, Modelling, Computation and Optimization in Information Systems and Management Sciences, 2008, Metz, France, Communications in Computer and Information Science, CCIS-14, Springer-Verlag, 419-428.
- CUDA Programming Guide : [http://www.nvidia.com/object/cuda\\_develop.html](http://www.nvidia.com/object/cuda_develop.html)
- CUBLAS :  
[http://developer.download.nvidia.com/compute/cuda/1\\_0/CUBLAS\\_Library\\_1.0.pdf](http://developer.download.nvidia.com/compute/cuda/1_0/CUBLAS_Library_1.0.pdf)
- nvidia 280 GTX : [http://www.nvidia.com/object/product\\_geforce\\_gtx\\_280\\_us.html](http://www.nvidia.com/object/product_geforce_gtx_280_us.html)