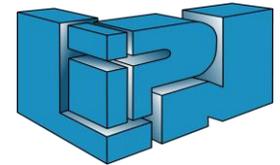


Université Paris 13 , Sorbonne Paris Cité,
LIPN, UMR 7030 du CNRS
99 Avenue J-B. Clément - 93430 Villetaneuse - France



CLUSTERING DE FLUX DE DONNÉES

Mustapha LEBBAH

MCF, HdR

LIPN – Univ. Paris 13

H. Azzag, T. Duong, **M. Ghesmoune (doctorant)**

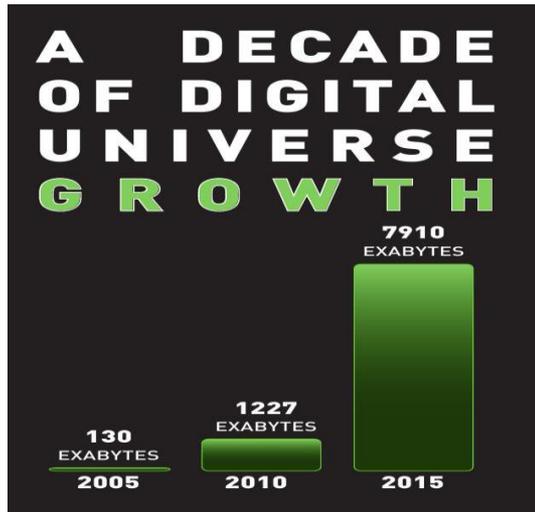
Journée Clustering, Issy Les Moulineaux le 20 Oct 2015



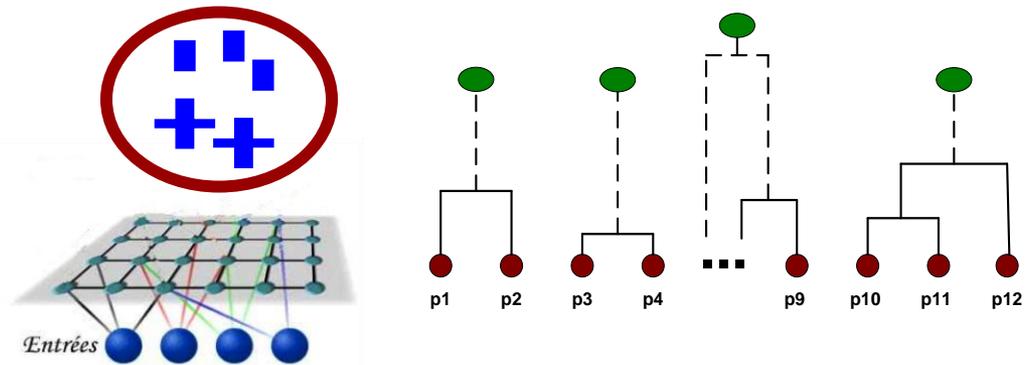
PLAN

- Contexte
- Le Big Data comme du «Data Stream»
- Stratégie du clustering de flux de données
- Approches utilisées
- Conclusion & perspectives

PARTITIONNEMENT «CLUSTERING» ET BIG DATA

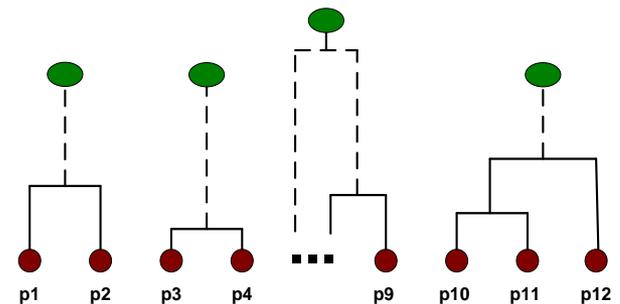
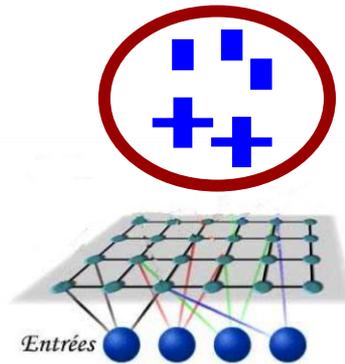
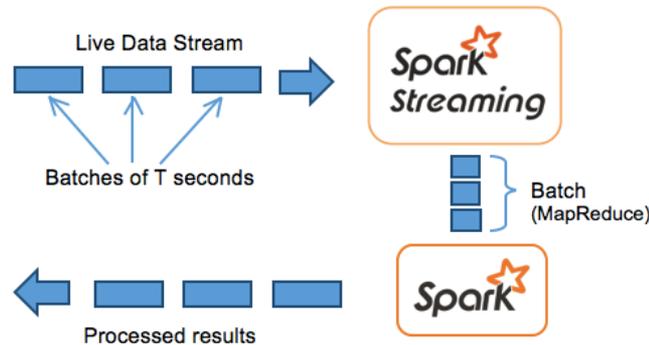
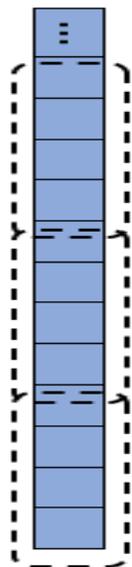
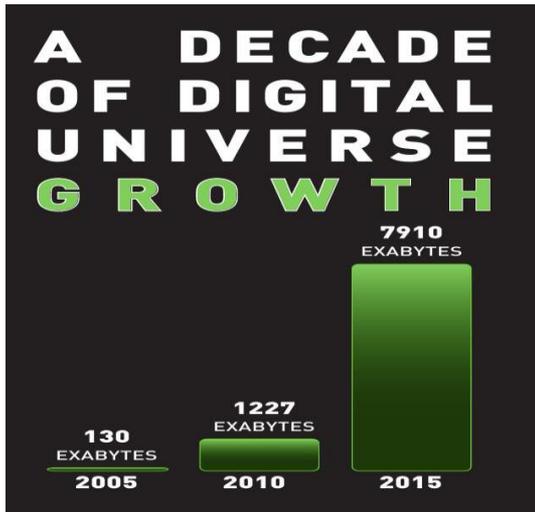


- Approche hiérarchique
- Approche par partitionnement
- Approche spectrale
- Approche à base de modèles auto-organisés
- Versions probabilistes



PARTITIONNEMENT «CLUSTERING» ET BIG DATA

- Approche hiérarchique
- Approche par partitionnement
- Approche spectrale
- Approche à base de modèles auto-organisés
- Versions probabilistes



FLUX DE DONNÉES / DATA STREAM

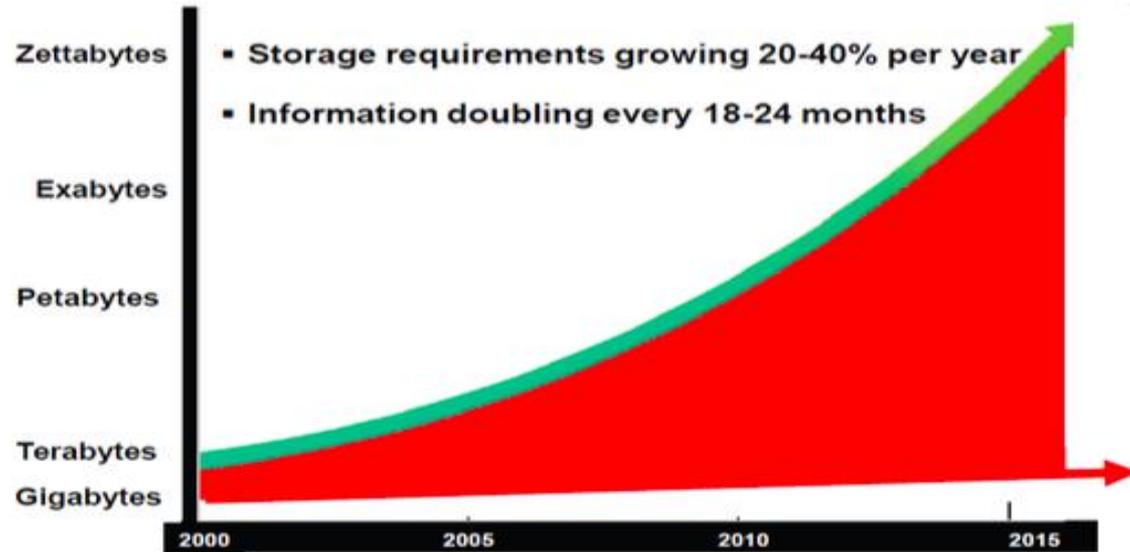
- Les données massives sont collectées à un rythme exceptionnel.
- Besoin d'une réponse rapide
- Exemples:
 - Les données de suivi de comportement pour étudier les changements dans les habitudes (achat/de circulation/intrusion)
 - Suivi des données météorologiques,
 - La taille des bases de données et des échantillons est énorme



LA COLLECTE DES DONNÉES NE S'ARRÊTE JAMAIS !



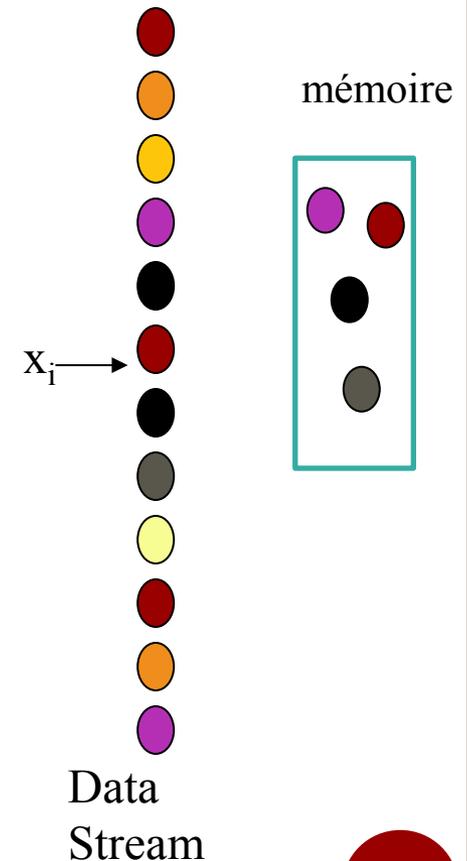
<http://www.flightradar24.com>



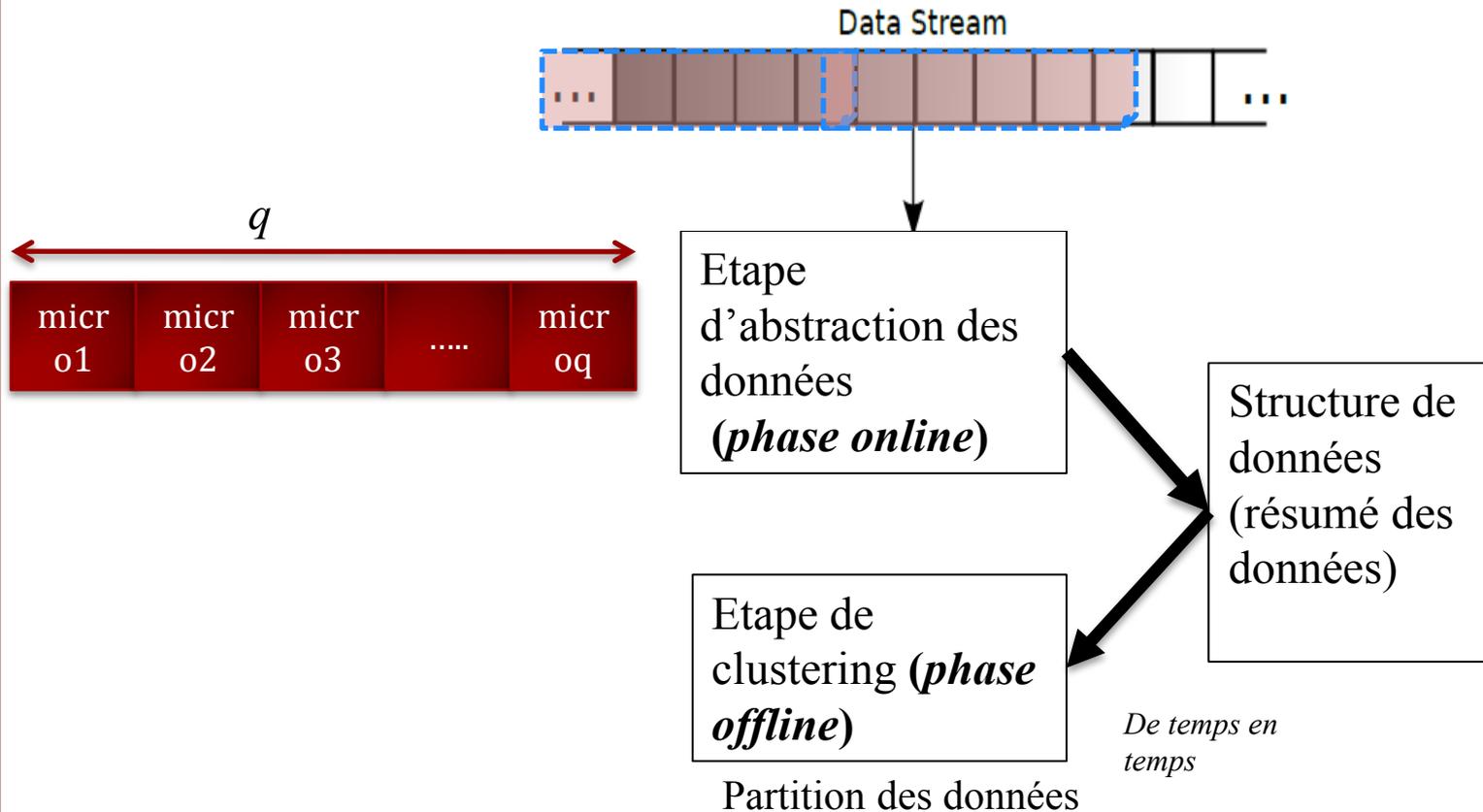
CHALLENGES

TRAITER LES « BIG DATA » COMME DU DATA STREAM

- Ensemble x_1, \dots, x_n d'observations
 - Potentiellement infini
 - Non-stationnaire
- Rapide, un seul passage
- Doit être traité dans cet ordre en une passe
- Ressources CPU et mémoire utilisées par élément faibles
- L'accès aléatoire n'est pas permis



LE BIG DATA COMME LE DATA STREAM



Framework de clustering de flux de données online-offline

CLUSTERING DE FLUX DE DONNÉES

PRINCIPE DE BASE

○ Stratégie

- **Fenêtre de temps**: diviser de manière indépendante le flux en fenêtres temporelles.

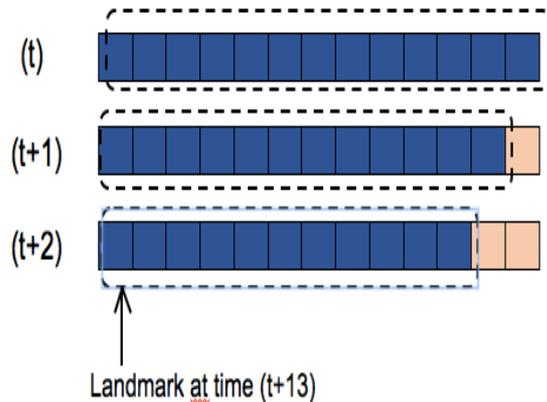


CLUSTERING DE FLUX DE DONNÉES

PRINCIPE DE BASE

o Stratégie

- **Fenêtre de temps**: diviser de manière indépendante le flux en fenêtres temporelles.
- Pondération des fenêtres temporelles (données) -> « fading »

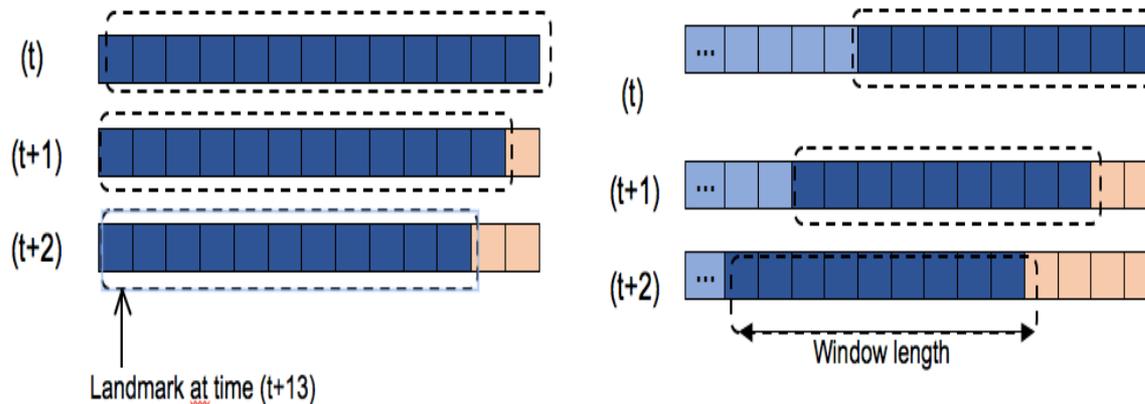


CLUSTERING DE FLUX DE DONNÉES

PRINCIPE DE BASE

o Stratégie

- **Fenêtre de temps**: diviser de manière indépendante le flux en fenêtres temporelles.
- Pondération des fenêtres temporelles (données) -> « fading »

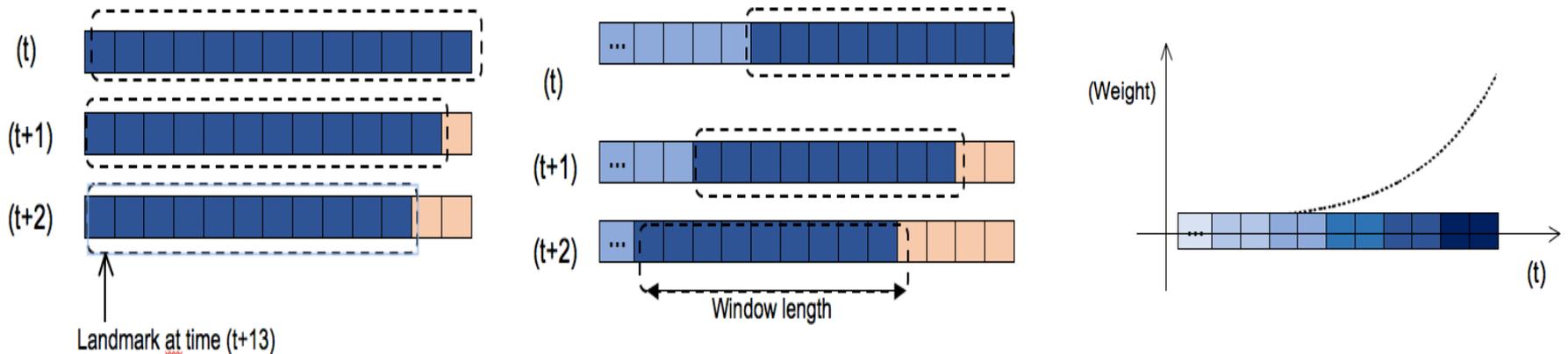
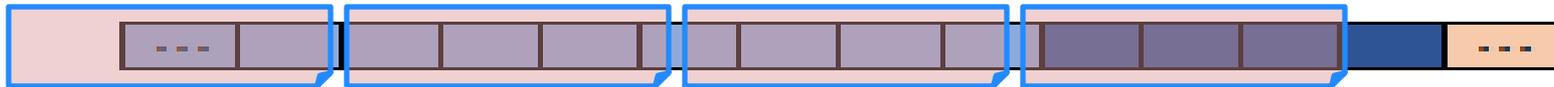


CLUSTERING DE FLUX DE DONNÉES

PRINCIPE DE BASE

o Stratégie

- **Fenêtre de temps**: diviser de manière indépendante le flux en fenêtres temporelles.
- Pondération des fenêtres temporelles (données) -> « fading »



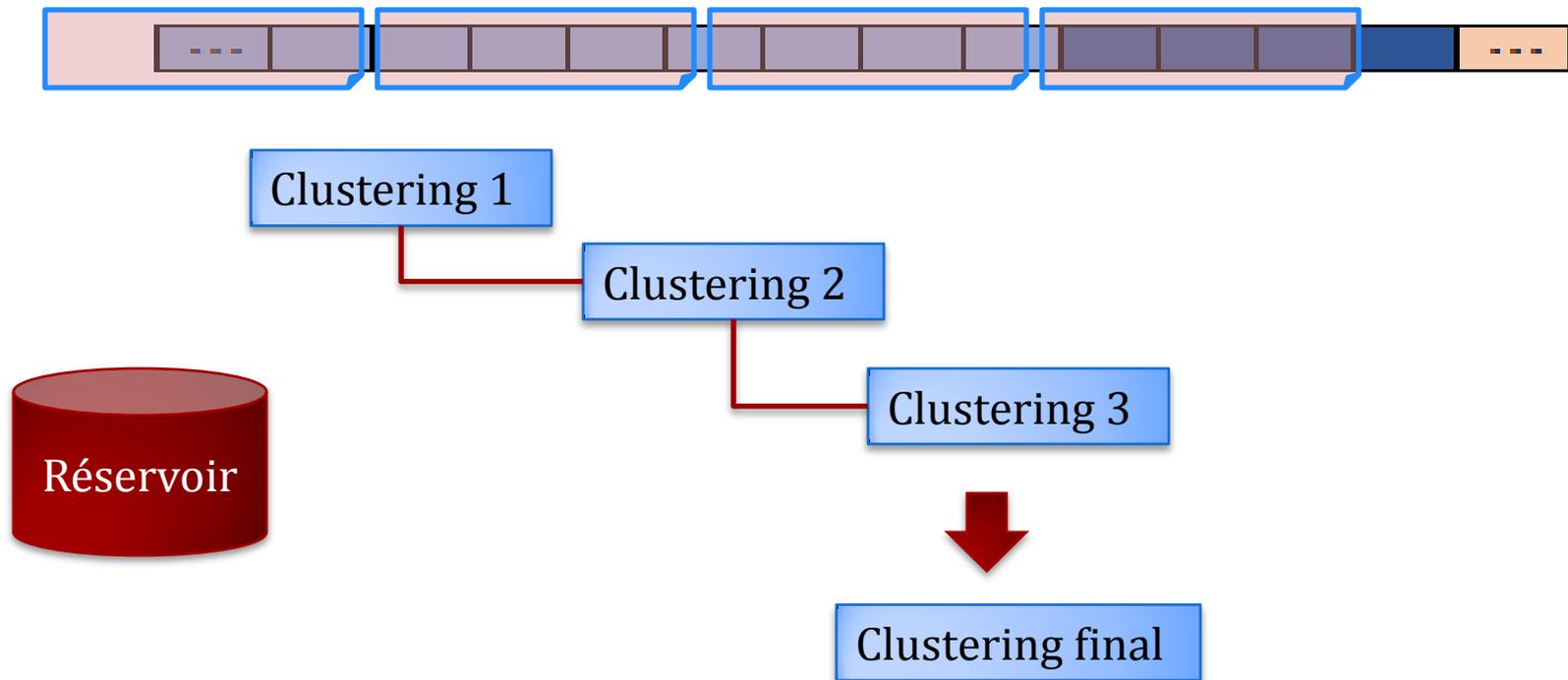
Décroissance exponentielle pour diminuer l'influence des données affectées aux micro-clusters (traiter la dérive de concept)

CLUSTERING DE FLUX DE DONNÉES

PRINCIPE DE BASE

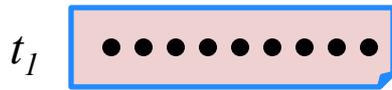
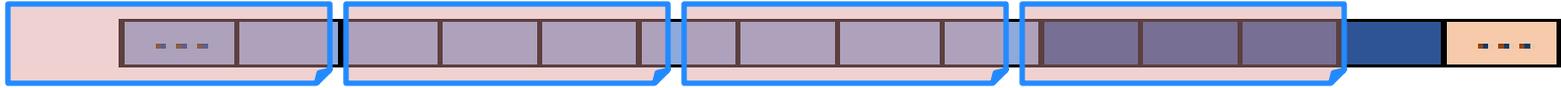
o Stratégie

- **Fenêtre de temps**: diviser le flux en fenêtres temporelles

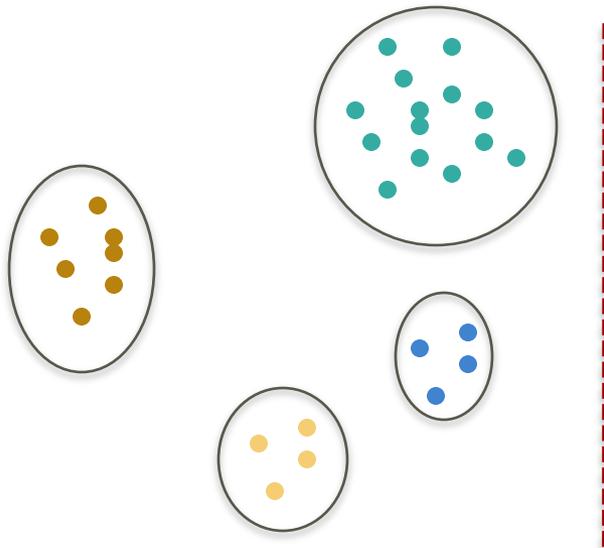


Partitionnement des micro-clusters: utilisation d'algorithme de clustering traditionnel (k-means, dbscan) en mode «off-line» de combiner.

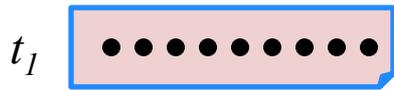
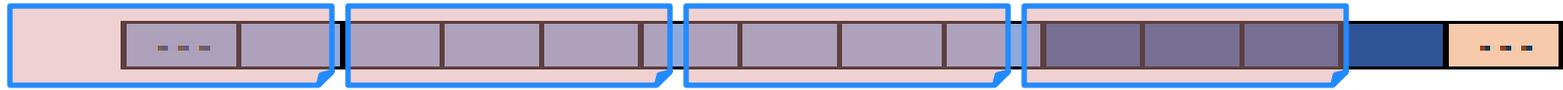
EVOLUTION D'UN CLUSTERING



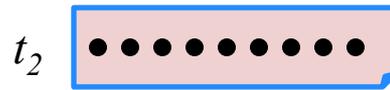
Clustering



EVOLUTION D'UN CLUSTERING

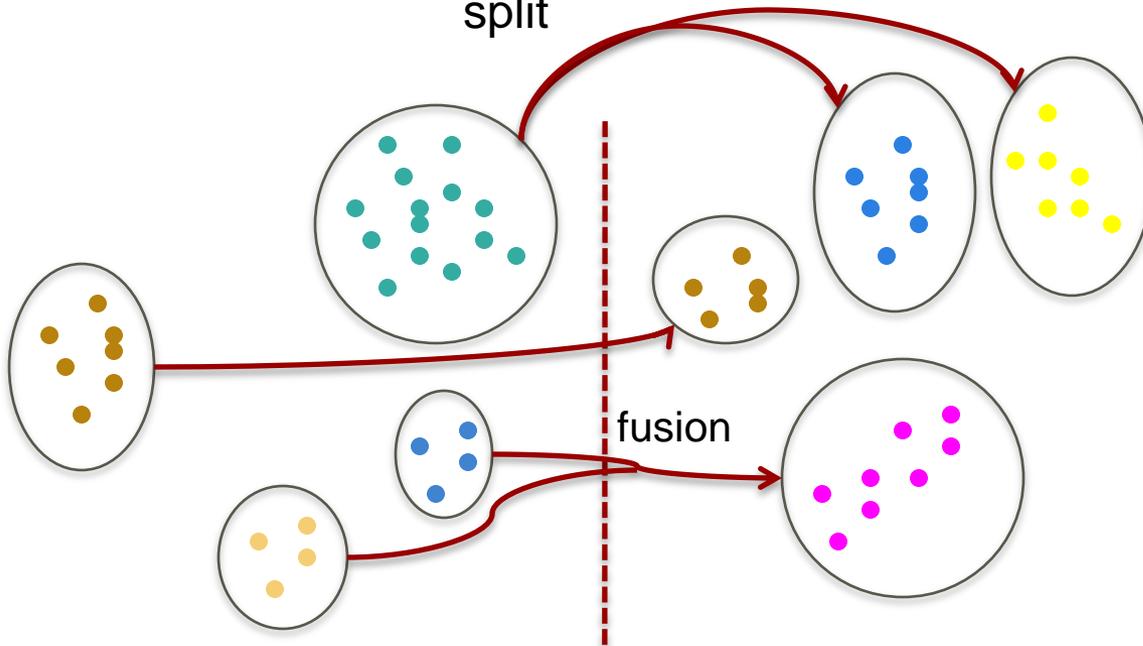


Clustering

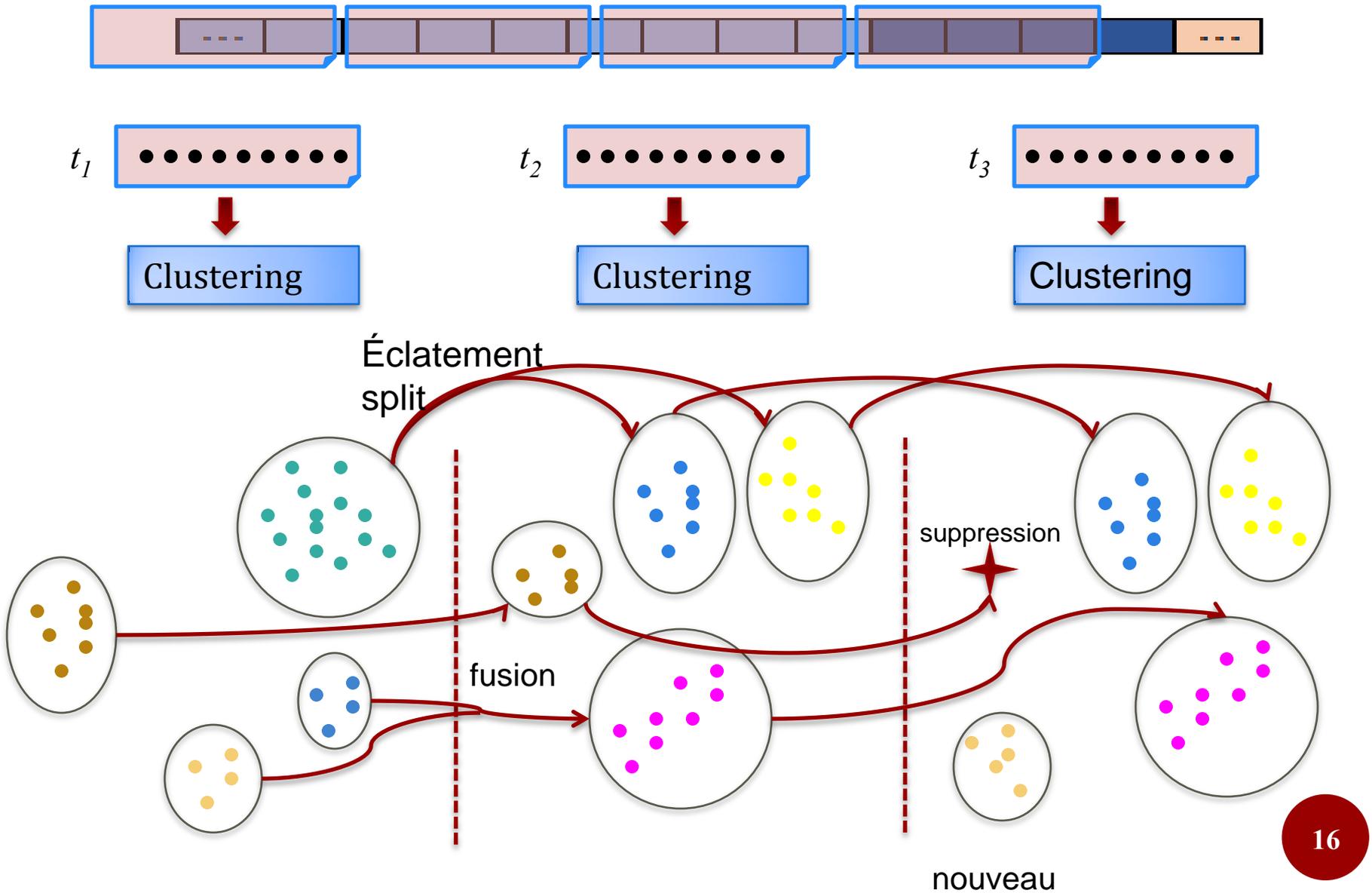


Clustering

Éclatement
split



EVOLUTION D'UN CLUSTERING



LES QUESTIONS

- Comment traiter les nouvelles données?
 - Affecter à un ancien micro-cluster
 - Création d'un nouveau micro-cluster
- Comment traiter les anciens micro-clusters
 - Supprimer
 - Fusionner des micro-clusters (les plus proches)

ALGORITHME SIMPLE

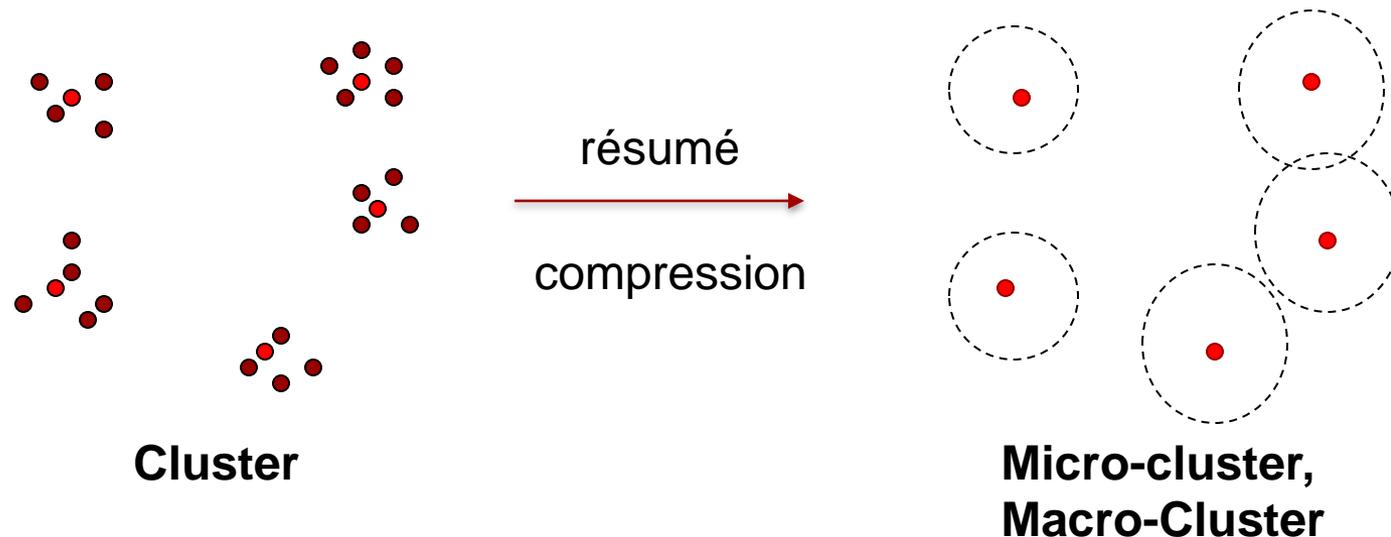
- Initialiser des **micro-cluster** vides
- Pour chaque nouvelle observation
 - Trouver le **micro-cluster** le plus proche
 - Affecter x au **micro-cluster** ou créer un autre
- Mise à jour de la structure sous-jacente des **micro-clusters** (une topologie, une hiérarchie)

CLUSTERING DE FLUX DE DONNÉES

PRINCIPE DE BASE

○ Stratégie : Micro-Clusters

- Un micro-cluster est une quantification « **résumé** » d'un ensemble de données qui sont proches les uns aux autres et seront traitées comme une seule unité.

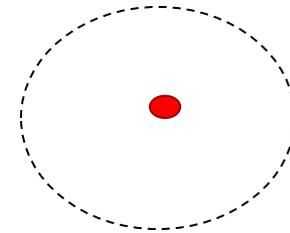


- Utiliser les prototypes : centre de gravité, vecteur caractéristique, variance...etc
- Utiliser le temps d'arrivée du cluster

VECTEUR CARACTÉRISTIQUE : BIRCH

- Utiliser la structure ou un vecteur caractéristique pour sauvegarder les données

- $CF = (N, \sum \mathbf{x}_i, \sum \mathbf{x}_i^2)$



$$CF = (N, LS, SS)$$

- Permet de déduire : le centre de gravité, rayon, diamètre
- **Toutes les informations nécessaires pour le clustering**

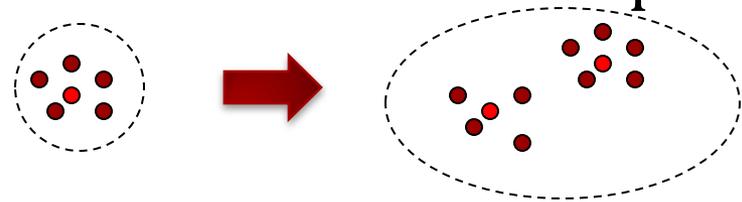
Birch: Balanced Iterative Reducing and Clustering using Hierarchies, by Zhang, Ramakrishnan, Livny 1996

VECTEUR CARACTÉRISTIQUE : BIRCH

- Permet de déduire : le centre de gravité, rayon, diamètre

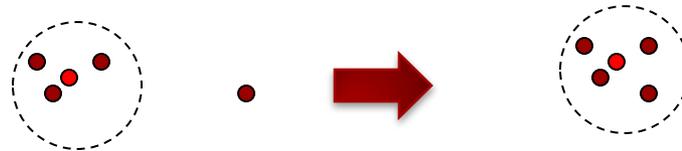
- **Incrémentalité** : à la présentation d'un nouveau point

- $LS=LS+X$
- $SS=SS+ (X)^2$
- $N=N+1$

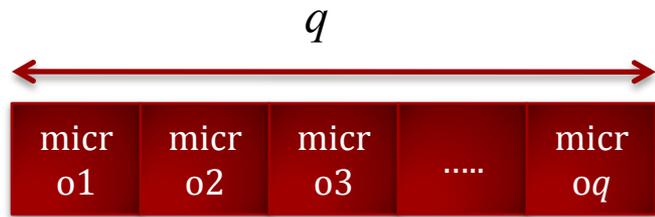


- **Additivité** : fusion de micro-cluster ($CF1, CF2$)

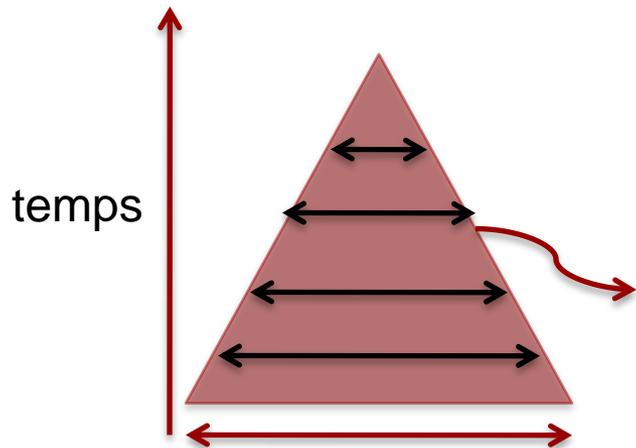
- $N3=N1+N2$
- $LS3=LS1+LS2$
- $SS3=SS1+SS2$



MICROCLUSTERS



$$CF = (N, \sum \mathbf{x}_i, \sum \mathbf{x}_i^2, LST, SST)$$



Taille de la fenêtre

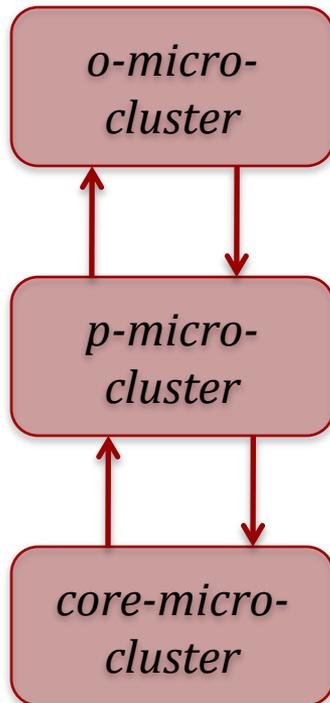
CluStream

- Sauvegarder la cardinalité, la somme linéaire et la somme au carré de des données du micro-cluster
- Fusion/suppression des clusters proches
- Pour chaque période on génère les macro-clusters
- le snapshot contient q micro-clusters, q dépend de la disponibilité de la mémoire

MICROCLUSTERS



$$CF = (N, \sum \mathbf{x}_i, \sum \mathbf{x}_i^2, LST, SST)$$



DenStream

• Introduction des clusters “outliers”
Les Micro-Clusters sont classés par rapport à leur poids w :

Si $w \geq \mu$, MC est un *core-micro-cluster*

Durant la partie online, on distingue:

1. **“Potential core-micro-clusters”** (p-micro-clusters), avec $w \geq \beta \cdot \mu$
2. **“Outlier micro-clusters”** (o-micro-clusters), avec $w < \beta \cdot \mu$

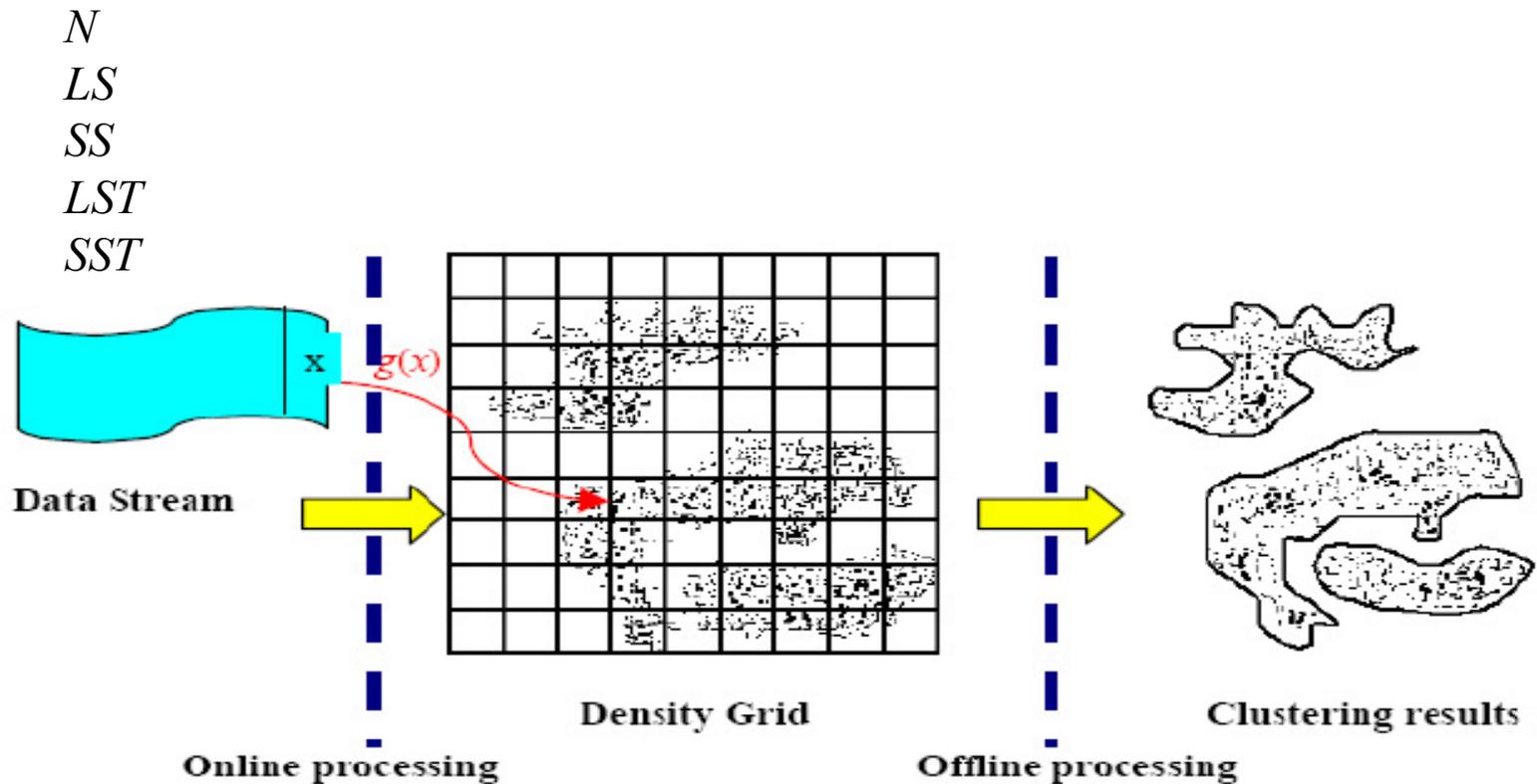
Fenêtre temporelle: poids des données diminue avec le temps: $f(t) = 2^{-\alpha t}$, $\alpha > 0$

MICROCLUSTERS



D-Stream

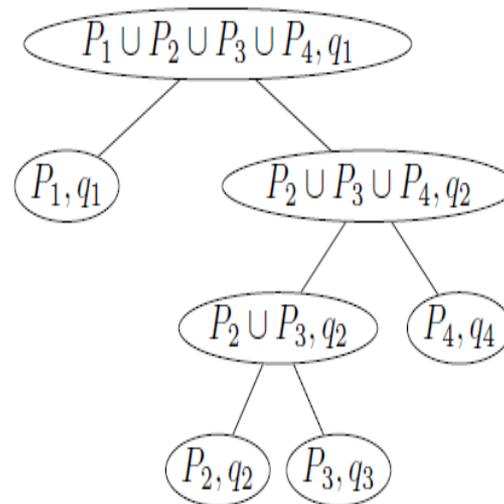
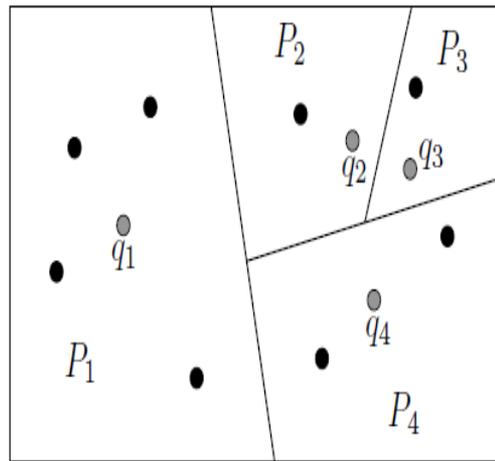
- Affecte la donnée à la grille
- Grilles pondérés par les récents points



STREAMKM++ (CORESETS)

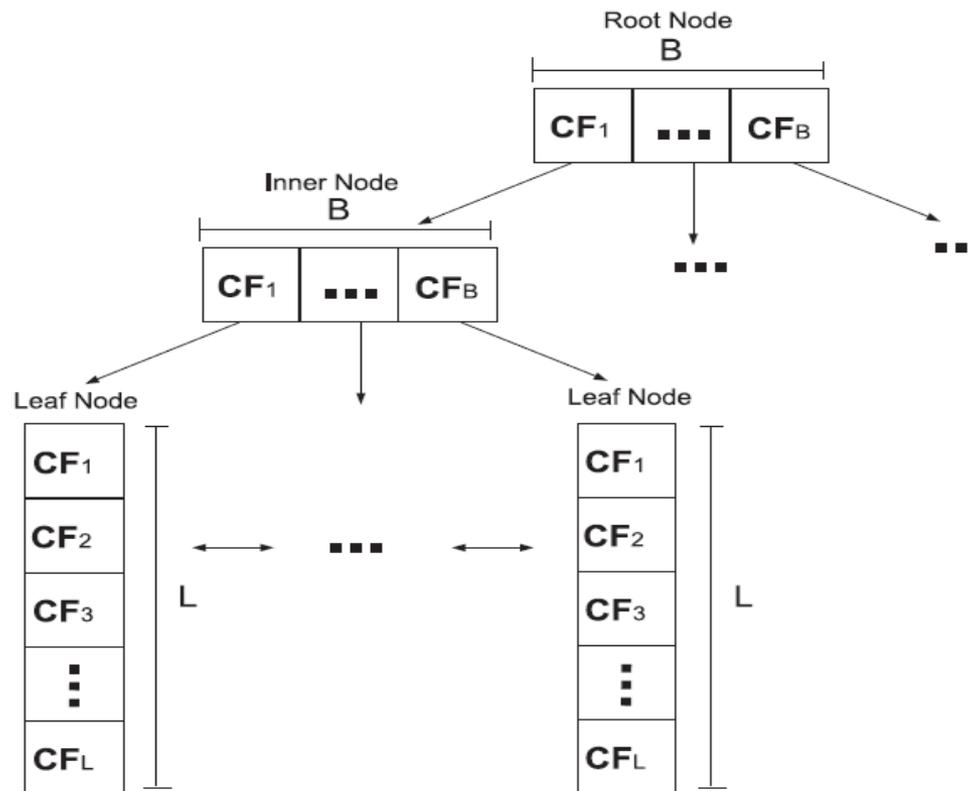
- Un ensemble S pondéré est un (k, ϵ) -corset pour un ensemble de données X si le partitionnement de S se rapproche du partitionnement de X avec une marge d'erreur de ϵ

$$(1 - \epsilon)dist(X, C) \leq dist_w(S, C) \leq (1 + \epsilon)dist(X, C)$$



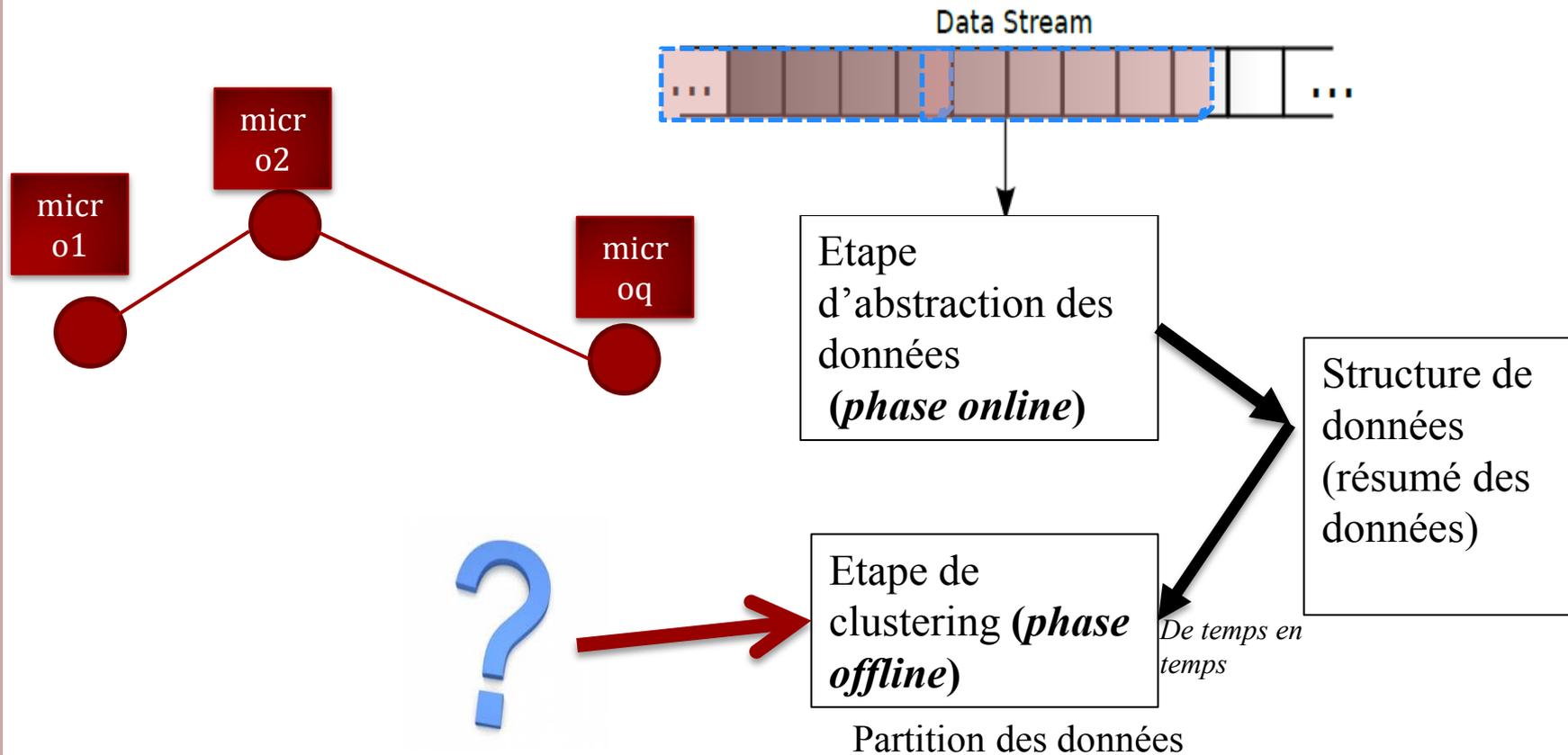
CONSTRUCTION D'UN ARBRE DE CF

- B-arbre, seuil de diamètre T et L le nombre de micro-clusters maximum par feuille



- Les “petits” micro-clusters sont considérés comme des “outliers”

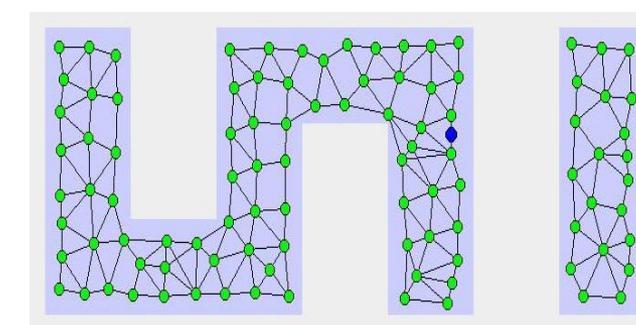
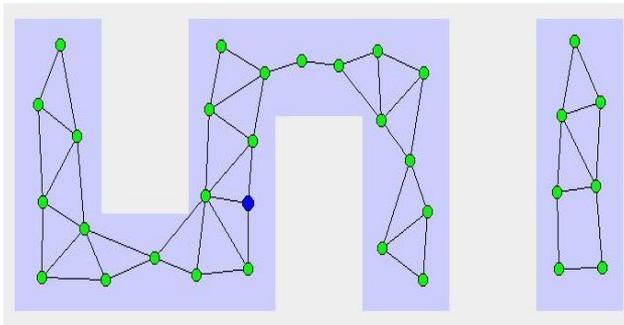
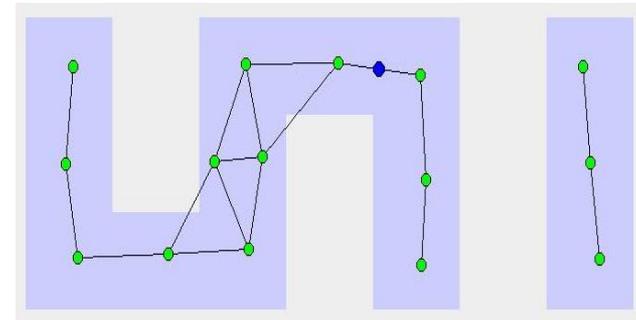
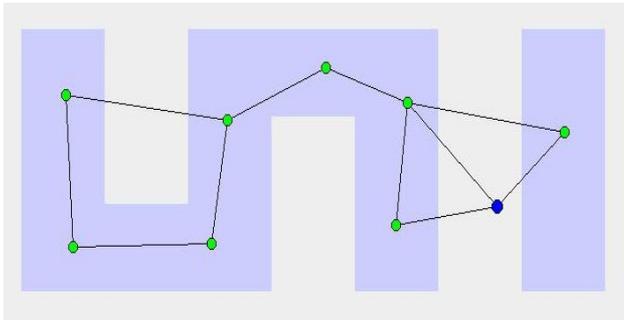
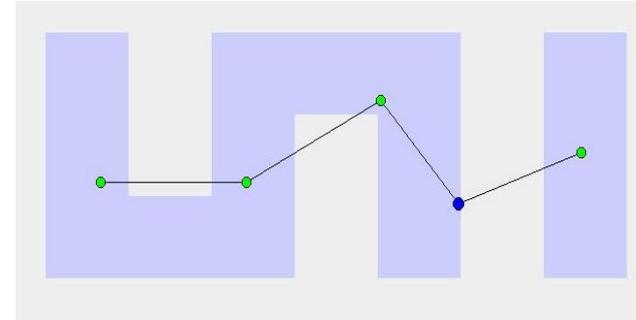
LE BIG DATA COMME LE DATA STREAM



Framework de clustering de flux de données online-offline

GNG

- Topologie évolutive
- Nombre de cellules non fixé

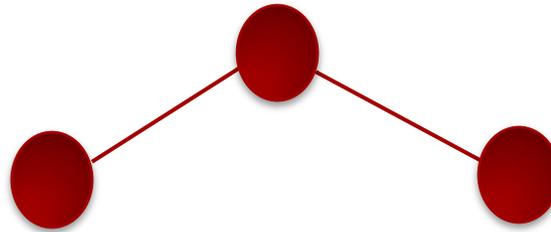


G-STREAM ET LE RESTE

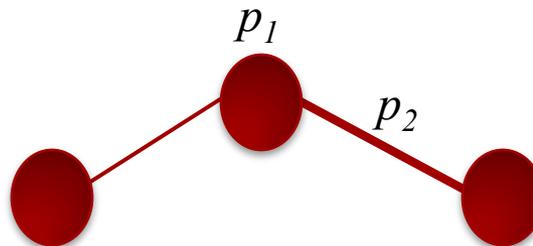
Algorithms	based on	topology	WL	phases	remove	merge	split	fade
G-Stream	NG	✓	✓	online	✓	✗	✗	✓
AING	NG	✓	✗	online	✗	✓	✗	✗
CluStream	<i>k</i> -means	✗	✗	2 phases	✓	offline	✗	✗
DenStream	DbScan	✗	✗	2 phases	✓	offline	✗	✓
SOSTream	DbScan, SOM	✗	✗	online	✓	✓	✗	✓
E-Stream	<i>k</i> -means	✗	✗	2 phases	✓	✓	✓	✓
StreamKM++	<i>k</i> -means	✗	✗	2 phases	✓	✓	✓	✓
StrAP	AP	✗	✗	2 phases	✓	✗	✗	✓
SVStream	SVC, SVDD	✗	✗	online	✓	✓	✓	✓

CARACTÉRISTIQUES

- Un graphe représentant la structure topologique,

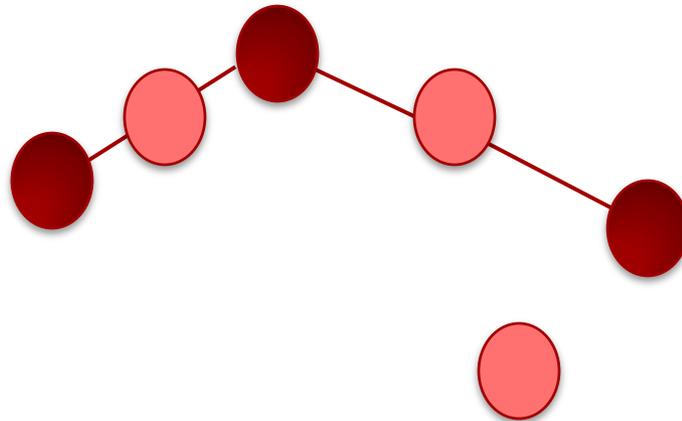


- Pondération des liens et des nœuds



G-STREAM: CARACTÉRISTIQUES ET AVANTAGES

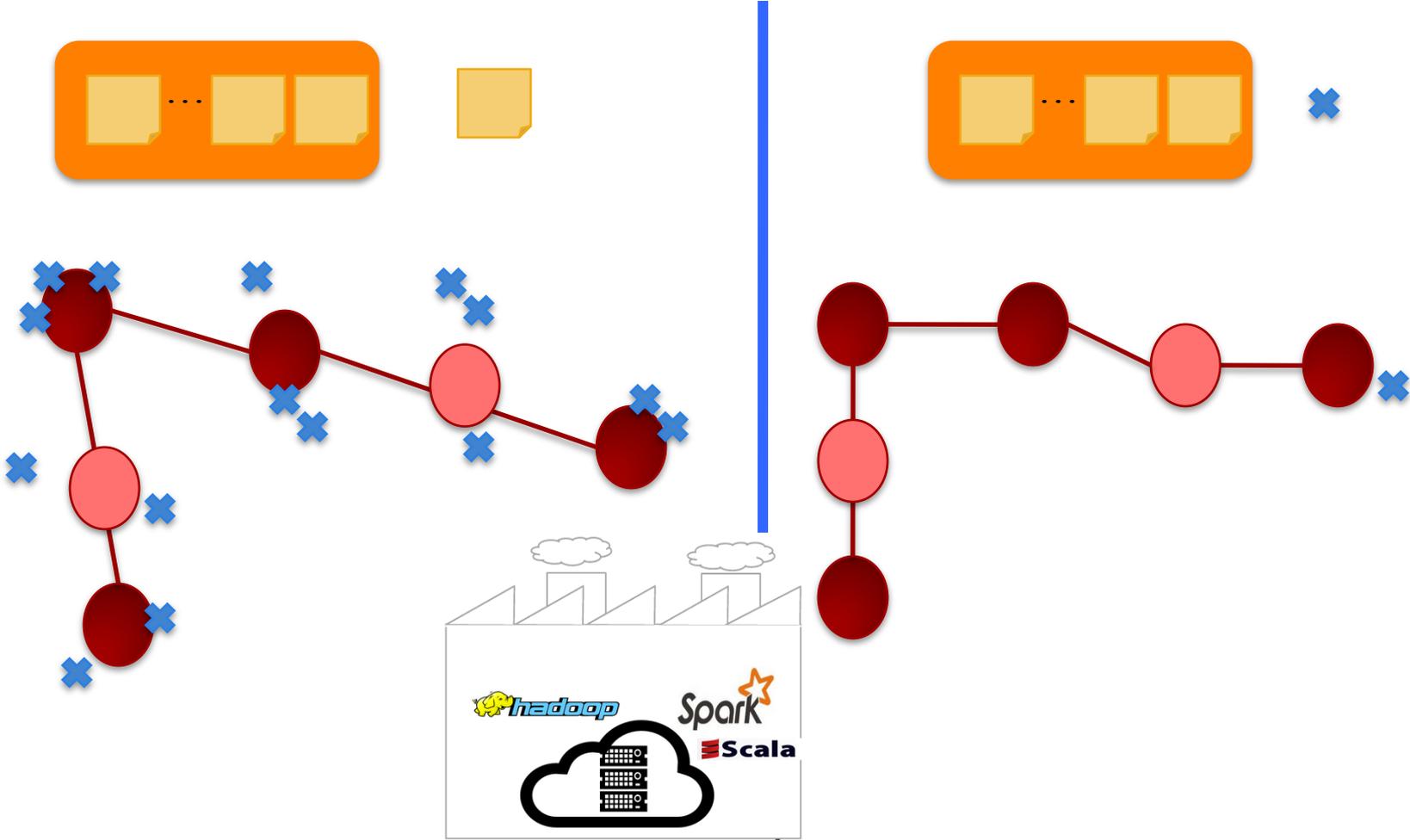
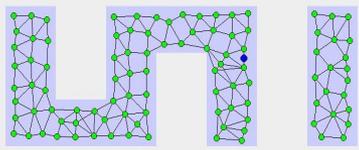
- Pas de phase d'initialisation du modèle,
- Création de plusieurs nœuds (potentiel-micro-cluster) à la fois,



- Utilisation d'un réservoir



G-STREAM

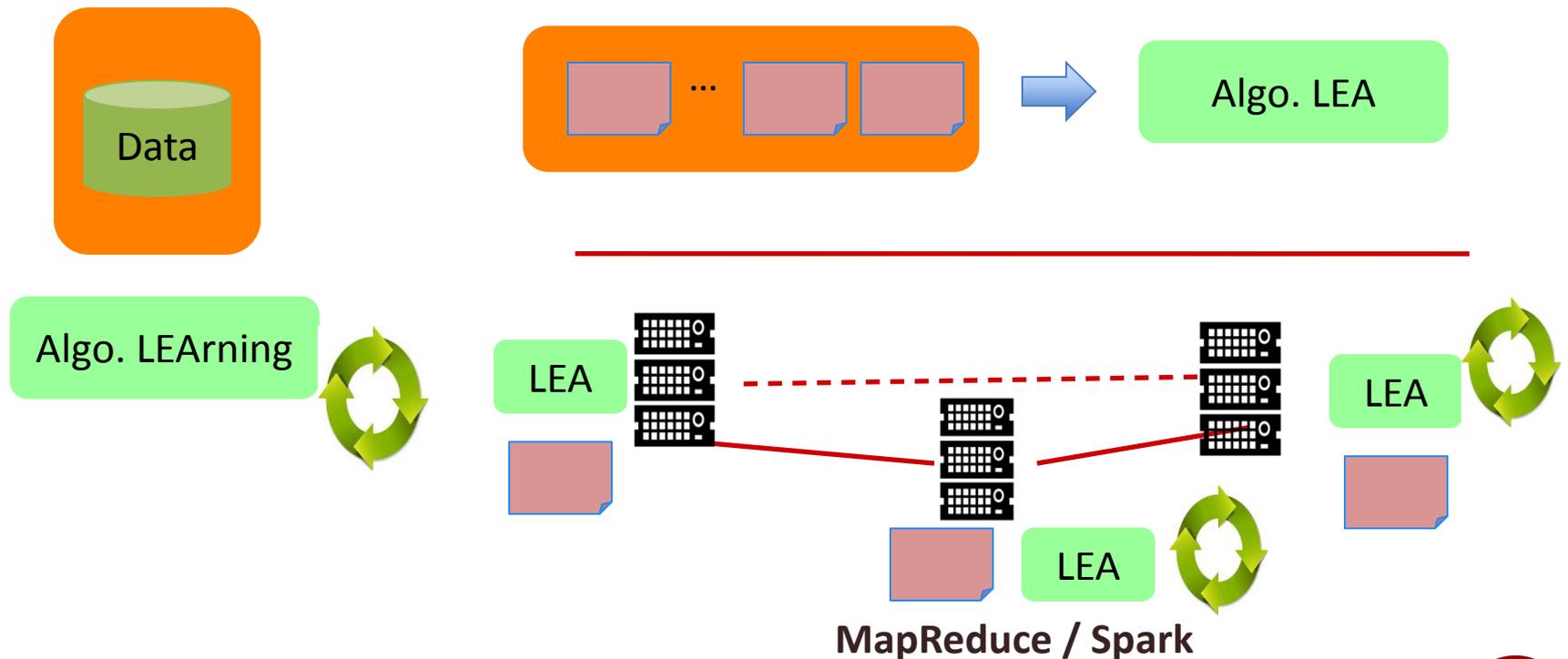


Mise à jour d'une manière «Batch»

Mise à jour d'une manière stochastique

Stream + Batch + MapReduce

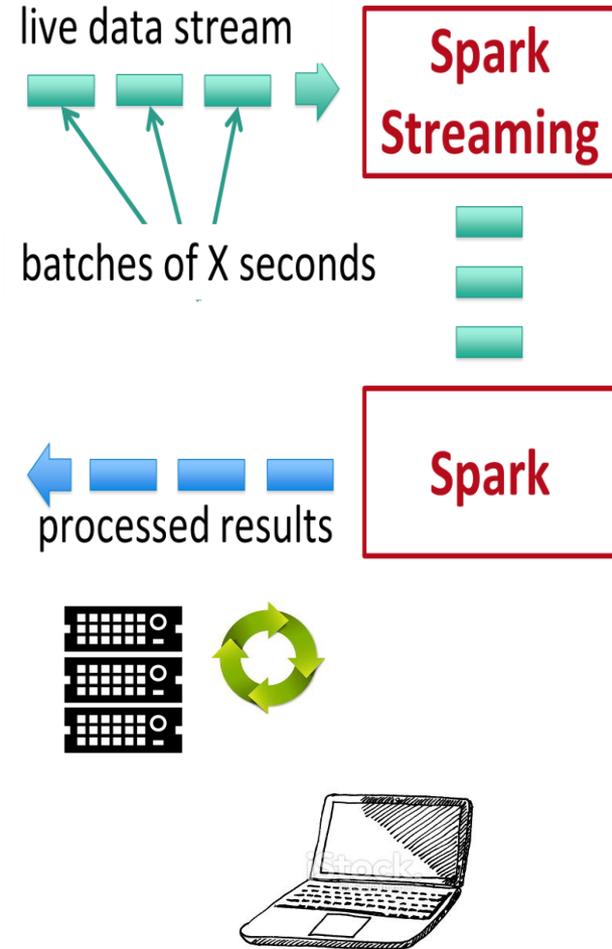
Aucune maîtrise de la distribution des données dans l'écosystème.



SPARK STREAMING

Combinaison du traitement du flux et le traitement en batch

```
for (i <- 1 until 10) {  
  val closest = data.map(p =>  
    (closestCentroid(p, centroids),  
    (p, 1))  
  )  
  val pointStats=closest.reduceByKey(  
    case ((p1, sum1), (p2, sum2)) =>  
    (p1 + p2, sum1 + sum2)  
  )  
  pointStats.foreach{case(id, value)  
=>  
    centroids(id) = value._1 / value._2  
  }  
}
```



G-STREAM : PRINCIPE DE LA VERSION BATCH

$$\mathcal{DS} = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_L\}$$

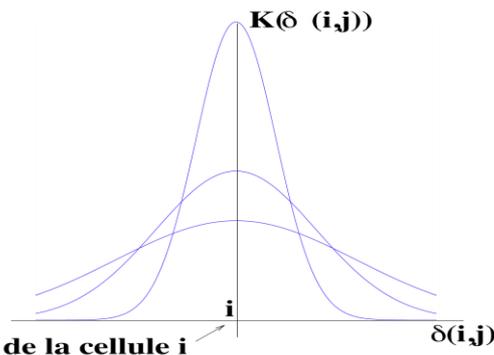
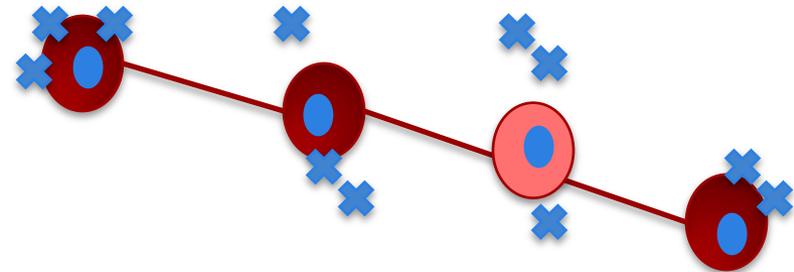
$$\mathbf{X}_i = \{\mathbf{x}_1, \dots, \mathbf{x}_p\}$$



Version batch de SOM

$\chi(\mathbf{x}_i)$ Affecter au plus proche

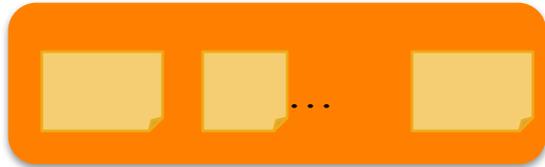
K Fonction de voisinage



G-STREAM : PRINCIPE DE LA VERSION BATCH

$$DS = \{X_1, X_2, \dots, X_L\}$$

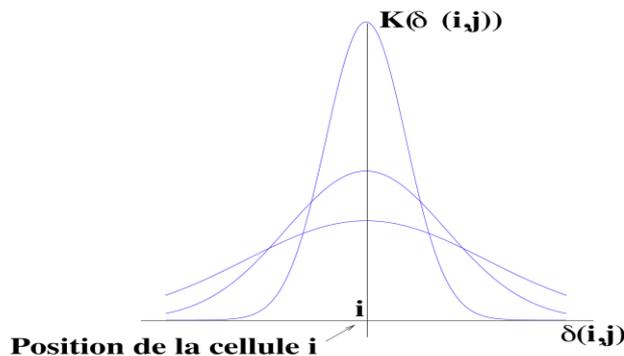
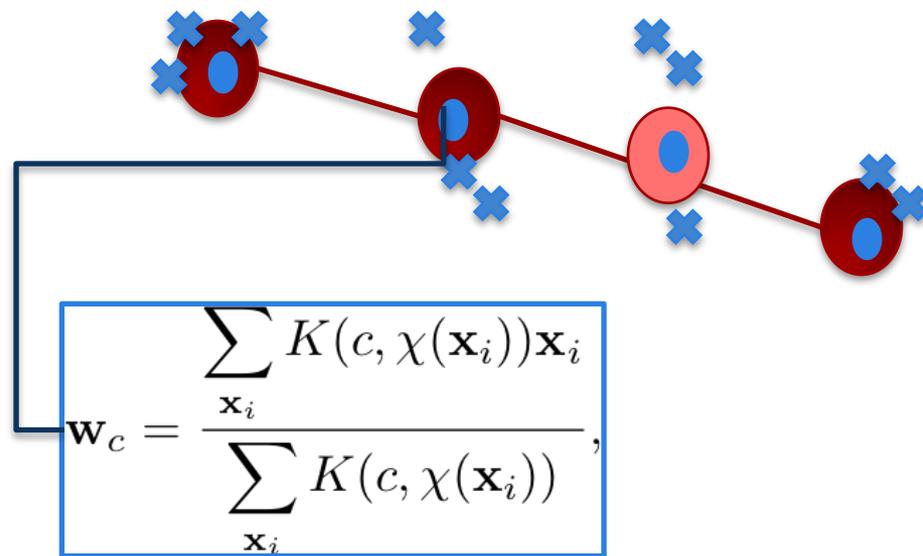
$$X_i = \{x_1, \dots, x_p\}$$



Version batch de SOM

$\chi(x_i)$ Affecter au plus proche

K Fonction de voisinage



Position de la cellule i

$$w_c = \frac{\sum_{x_i} K(c, \chi(x_i)) x_i}{\sum_{x_i} K(c, \chi(x_i))}$$

G-STREAM : LE PROTOTYPE



$DS^{(t)}$

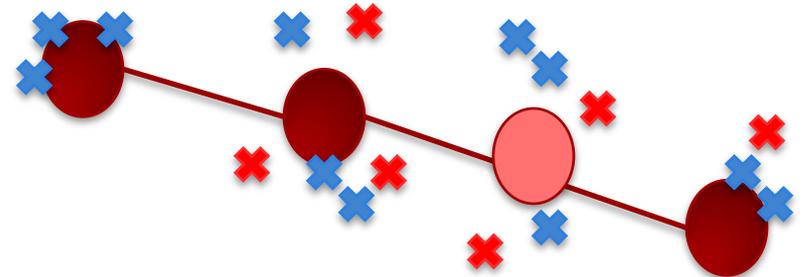


Nouvelles données

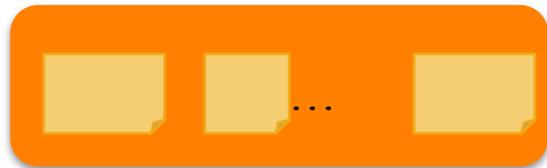


Anciennes données

$$\mathbf{w}_c^{(t+1)} = \frac{\sum_{\mathbf{x}_i \in DS^{(t)}} K(c, \chi(\mathbf{x}_i)) \mathbf{x}_i}{\sum_{\mathbf{x}_i \in DS^{(t)}} K(c, \chi(\mathbf{x}_i))},$$



G-STREAM : LE PROTOTYPE



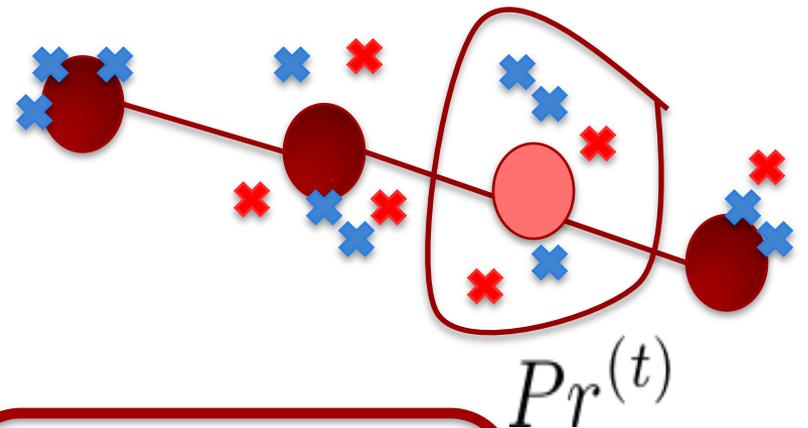
$DS^{(t)}$



Nouvelles données

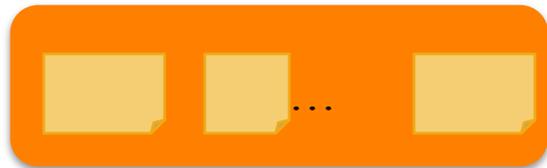
Anciennes données

$$\mathbf{w}_c^{(t+1)} = \frac{\sum_{\mathbf{x}_i \in DS^{(t)}} K(c, \chi(\mathbf{x}_i)) \mathbf{x}_i}{\sum_{\mathbf{x}_i \in DS^{(t)}} K(c, \chi(\mathbf{x}_i))},$$



$$\mathbf{w}_c^{(t+1)} = \frac{\sum_{r \in C} \sum_{\mathbf{x}_i \in P_r^{(t-1)}} K(c, r) \mathbf{x}_i + \sum_{r \in C} \sum_{\mathbf{x}_i \in P_r^{(t)}} K(c, r) \mathbf{x}_i}{\sum_{r \in C} \sum_{\mathbf{x}_i \in P_r^{(t-1)}} K(c, r) + \sum_{r \in C} \sum_{\mathbf{x}_i \in P_r^{(t)}} K(c, r)},$$

G-STREAM : LE PROTOTYPE



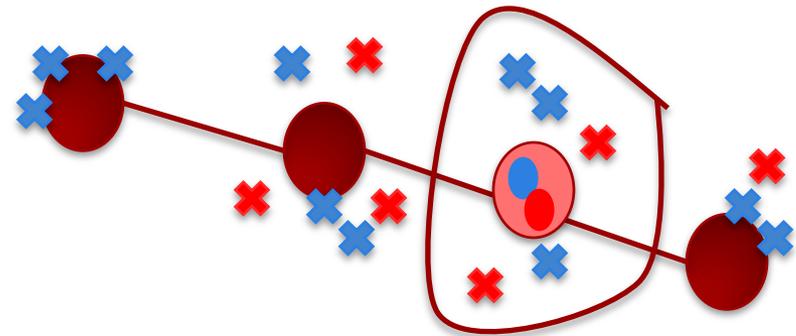
$DS^{(t)}$



Nouvelles données

Anciennes données

$$\mathbf{w}_c^{(t+1)} = \frac{\sum_{\mathbf{x}_i \in DS^{(t)}} K(c, \chi(\mathbf{x}_i)) \mathbf{x}_i}{\sum_{\mathbf{x}_i \in DS^{(t)}} K(c, \chi(\mathbf{x}_i))},$$

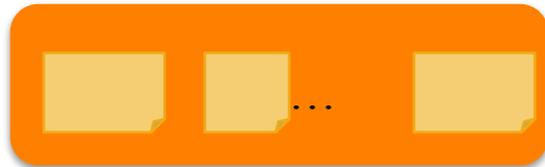


\mathbf{W}_c

\mathbf{Z}_r

$$\mathbf{w}_c^{(t+1)} = \frac{\sum_{r \in C} K(c, r) \mathbf{w}_c^{(t-1)} n_r^{(t-1)} + \sum_{r \in C} K(c, r) \mathbf{z}_r^{(t)} m_r^{(t)}}{\sum_{r \in C} K(c, r) n_r^{(t-1)} + \sum_{r \in C} K(c, r) m_r^{(t)}}$$

G-STREAM : LE PROTOTYPE



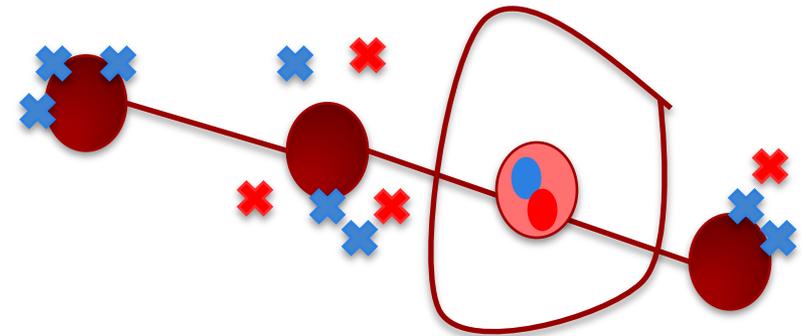
$DS^{(t)}$



Nouvelles données

Anciennes données

$$\mathbf{w}_c^{(t+1)} = \frac{\sum_{\mathbf{x}_i \in DS^{(t)}} K(c, \chi(\mathbf{x}_i)) \mathbf{x}_i}{\sum_{\mathbf{x}_i \in DS^{(t)}} K(c, \chi(\mathbf{x}_i))},$$



\mathbf{W}_c

\mathbf{Z}_r

$$\mathbf{w}_c^{(t+1)} = \frac{\sum_{r \in C} K(c, r) \mathbf{w}_c^{(t-1)} n_r^{(t-1)} + \sum_{r \in C} K(c, r) \mathbf{z}_r^{(t)} m_r^{(t)}}{\sum_{r \in C} K(c, r) n_r^{(t-1)} + \sum_{r \in C} K(c, r) m_r^{(t)}}$$

G-STREAM: AFFECTATION, MAP

$$DS = \{X_1, X_2, \dots, X_L\}$$



$$X_i = \{x_1, \dots, x_p\}$$



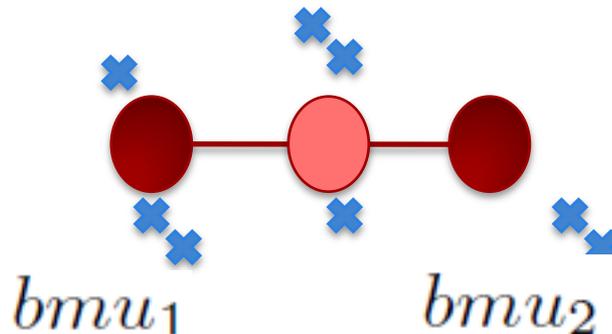
X_i



(bmu_1, bmu_2, X_i)



Numéro de cluster

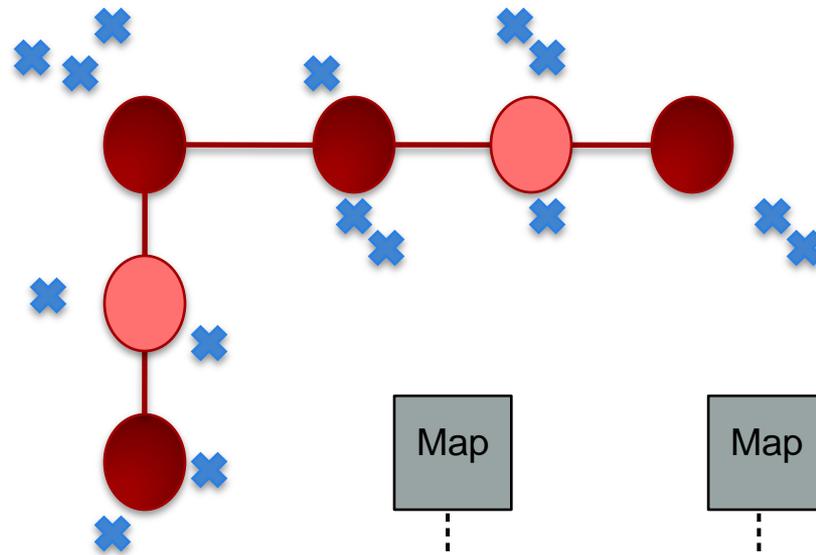


G-STREAM

$$\mathcal{DS} = \{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_L\}$$



$$\mathbf{X}_i = \{\mathbf{x}_1, \dots, \mathbf{x}_p\}$$

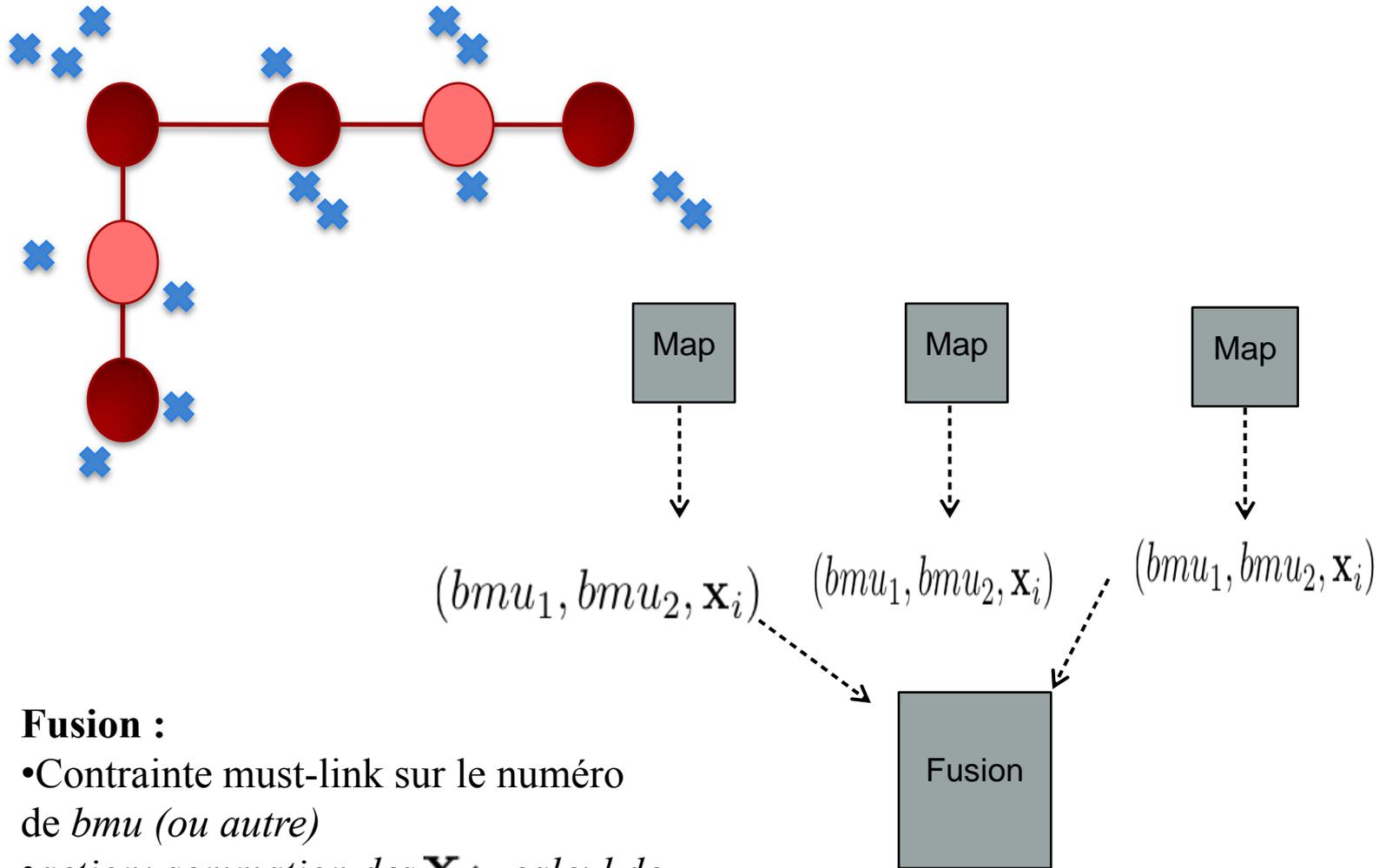


$(bmu_1, bmu_2, \mathbf{x}_i)$

$(bmu_1, bmu_2, \mathbf{x}_i)$

$(bmu_1, bmu_2, \mathbf{x}_i)$

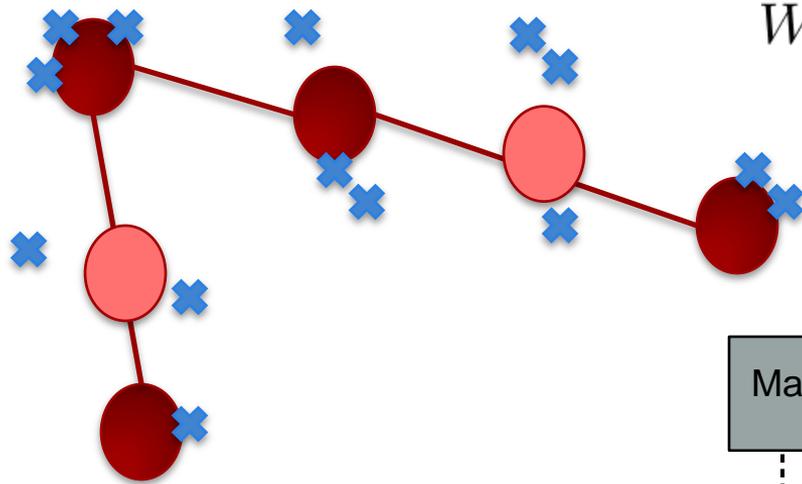
G-STREAM



Fusion :

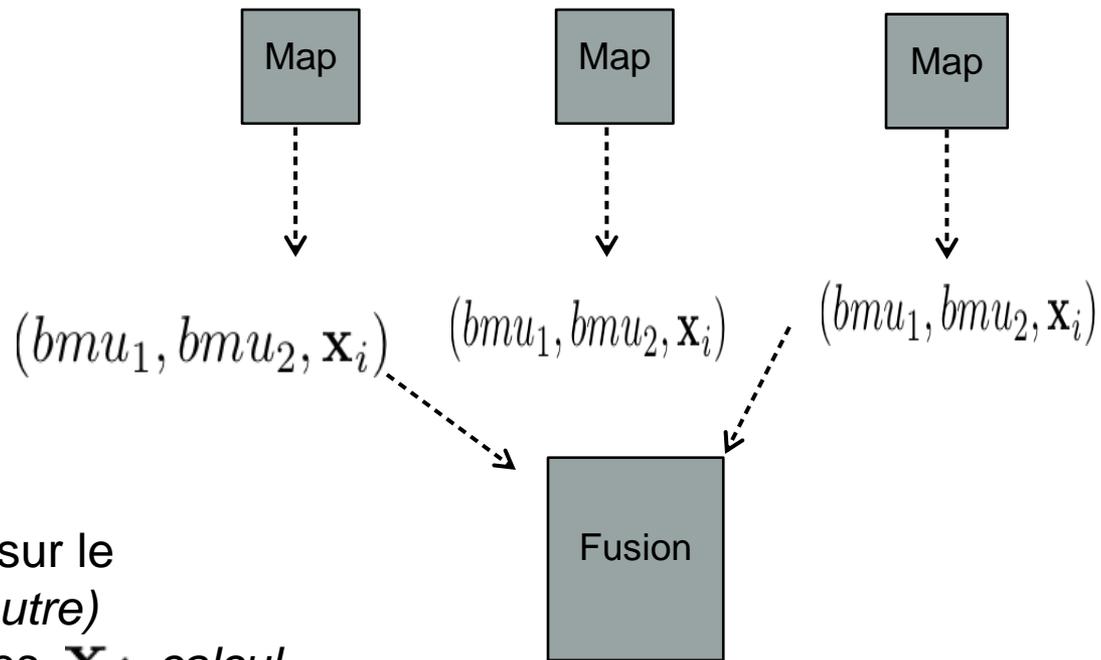
- Contrainte must-link sur le numéro de bmu (ou autre)
- action: sommation des \mathbf{X}_i , calcul de la fonction de voisinage, ...etc

G-STREAM



$$Weight(c) = \sum 2^{-\lambda(t-t_0)}$$

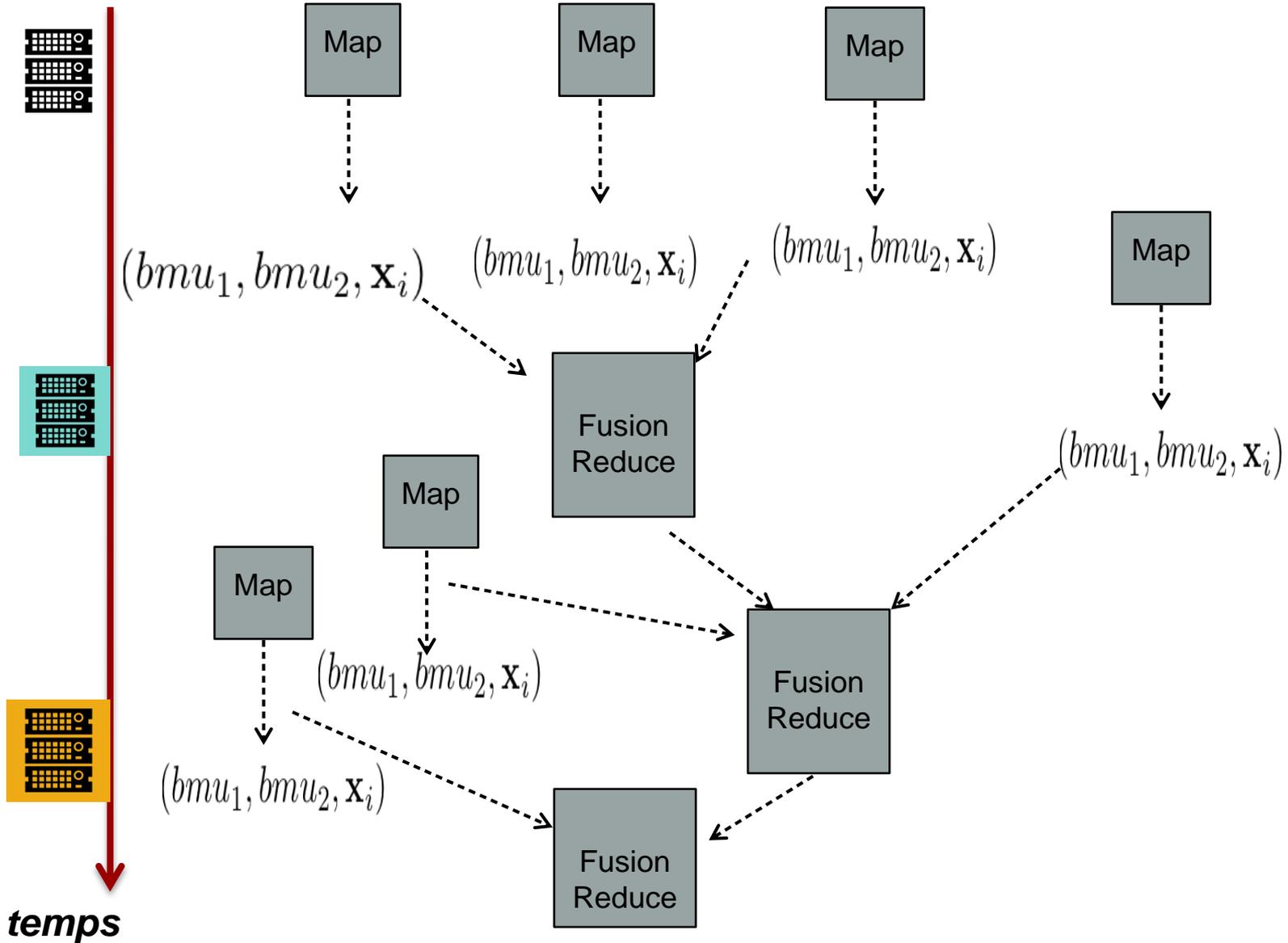
$$f(t) = 2^{-\lambda(t-t_0)}$$



Fusion :

- Contrainte must-link sur le numéro de bmu (ou autre)
- action: sommation des \mathbf{X}_i , calcul de la fonction de voisinage, ...etc

LA CHAINE MAP (AFFECTATION) ET RÉDUCE (FUSION)



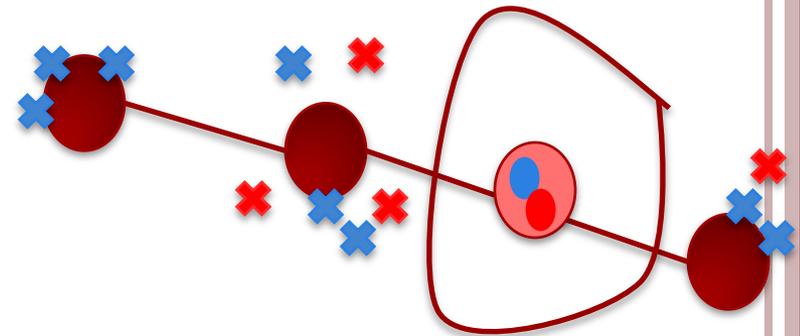
G-STREAM : LE PROTOTYPE



Nouvelles données

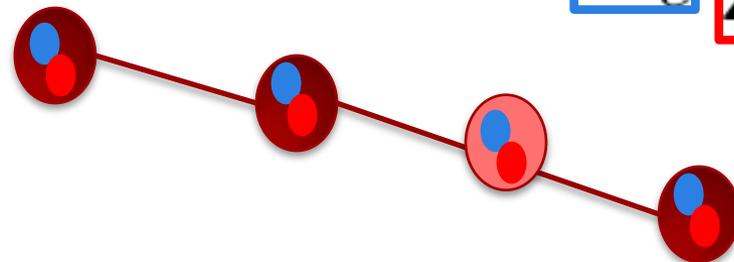
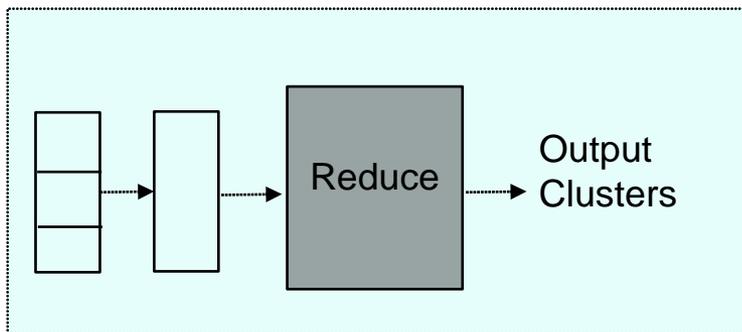
Anciennes données

$$\mathbf{w}_c^{(t+1)} = \frac{\mathbf{w}_c^{(t)} n_c^{(t)} \alpha + \sum_{r \in \mathcal{C}} \mathcal{K}(r, c) \mathbf{z}_r^{(t)} m_r^{(t)}}{n_c^{(t)} \alpha + \sum_{r \in \mathcal{C}} \mathcal{K}(r, c) m_r^{(t)}}$$



\mathbf{W}_c \mathbf{Z}_r

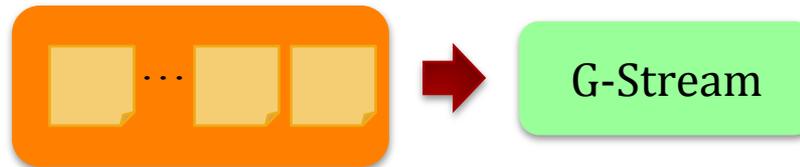
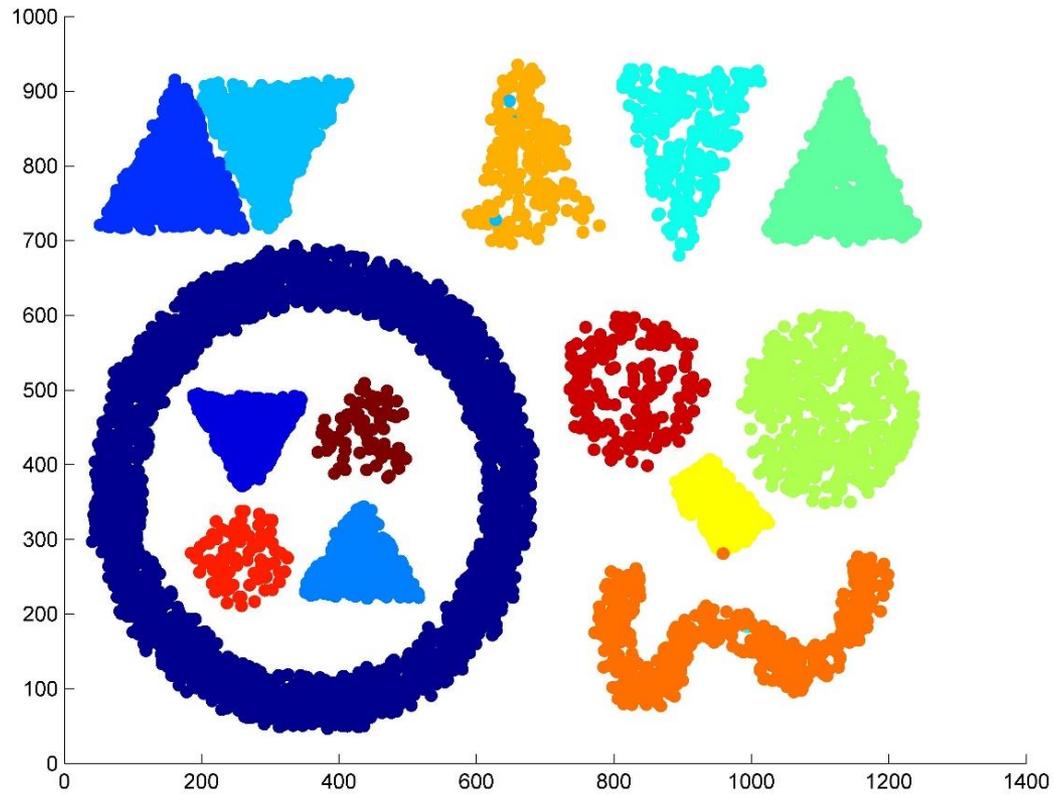
(key, [Value0, Value1, Value2,...])



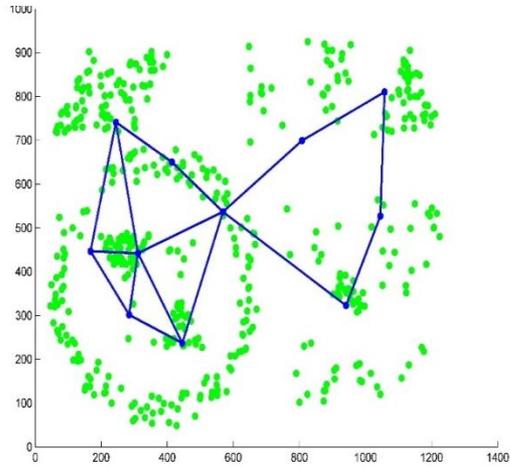
BASE DE DONNÉES

Datasets	#records	#features	#classes
DS1	9,153	2	14
DS2	5,458	2	13
letter4	9,344	2	7
Sea	60,000	3	2
HyperPlan	100,000	10	5
KddCup99	494,021	41	23
CoverType	581,012	54	7
Sensor	2,219,803	5	54

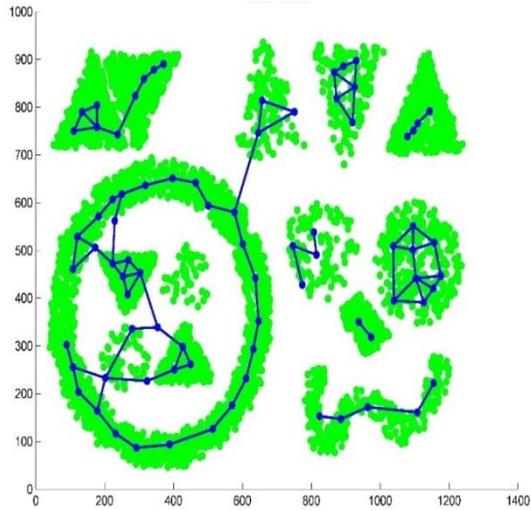
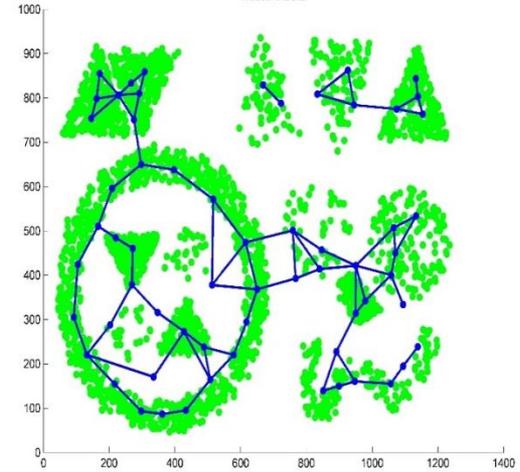
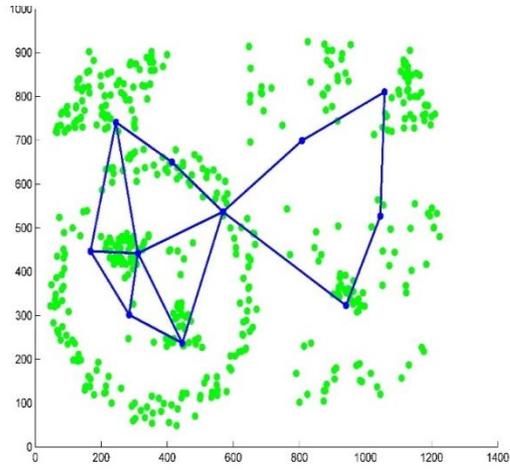
G-STREAM : EXEMPLE



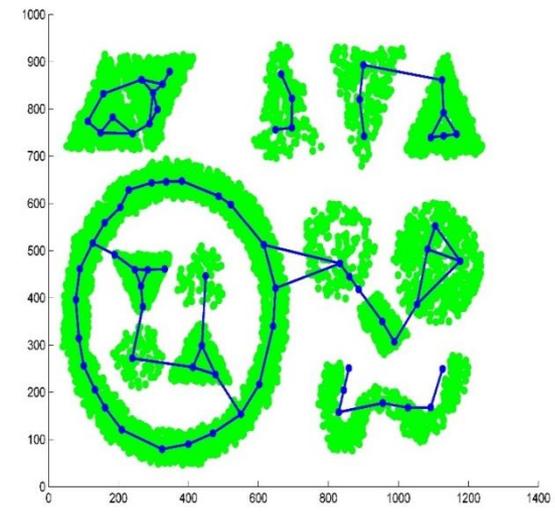
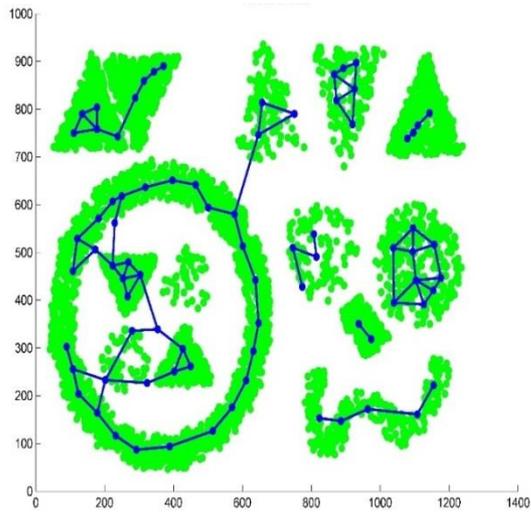
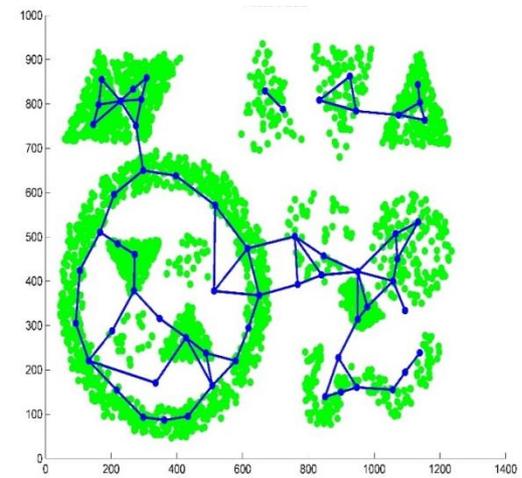
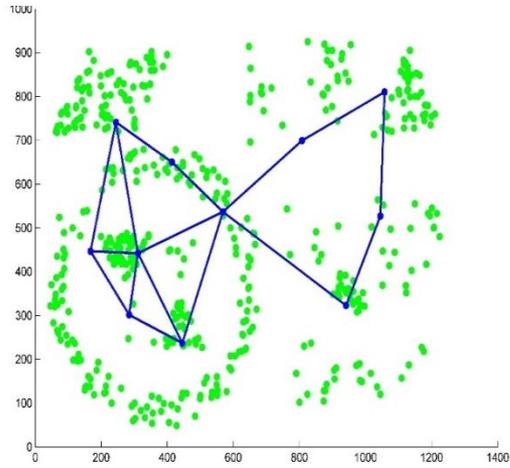
G-STREAM SUR DS1



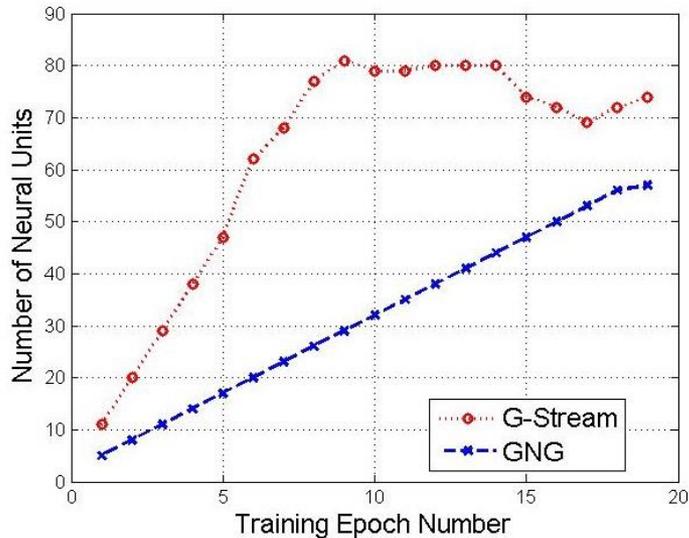
G-STREAM SUR DS1



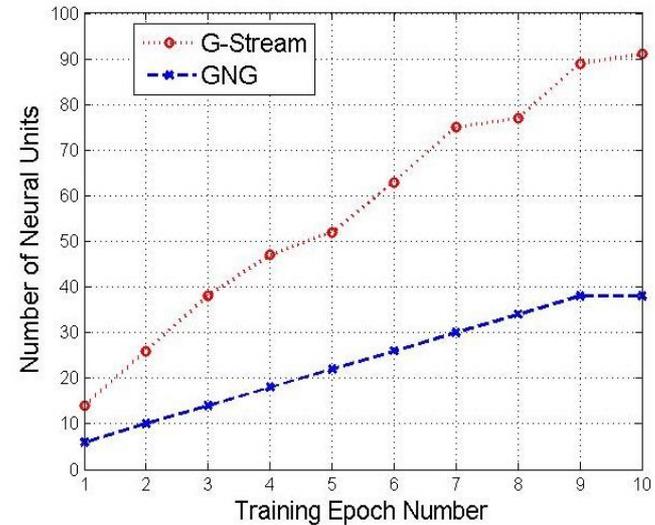
G-STREAM SUR DS1



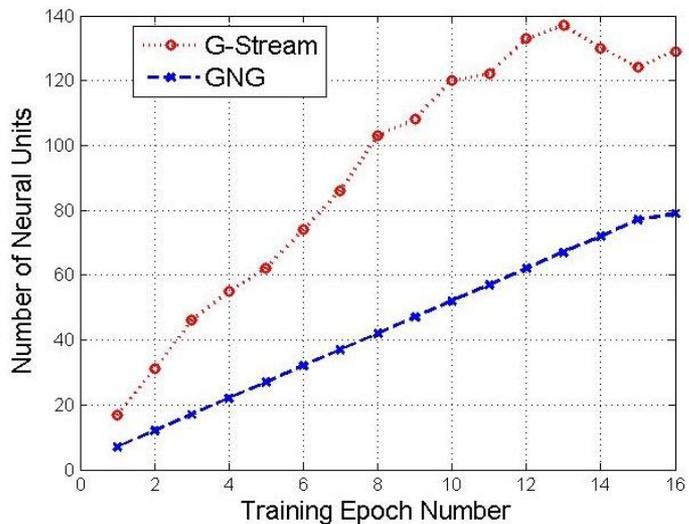
G-STREAM VS GNG-ONLINE: #NœUDS



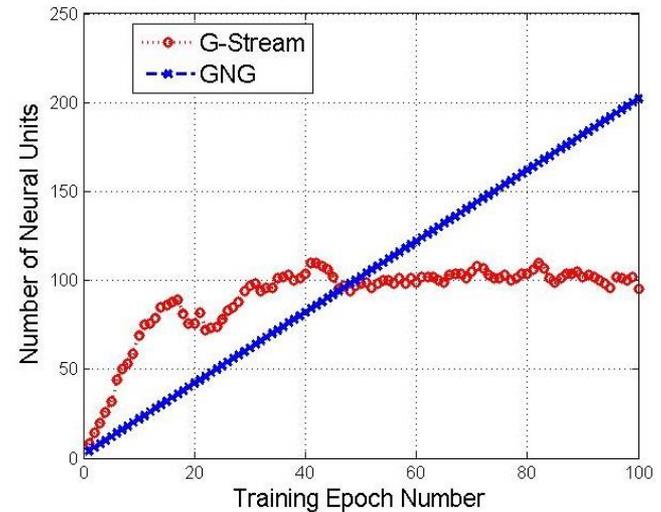
DS1



DS2

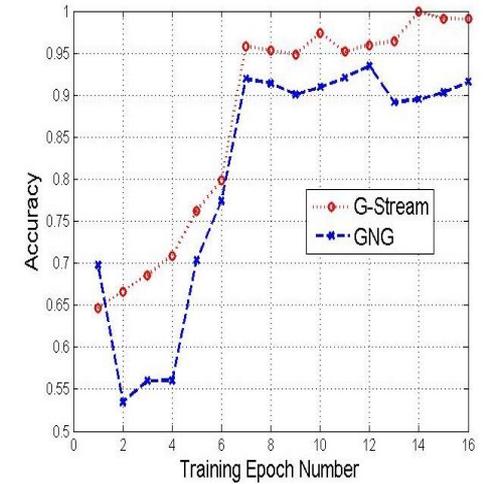
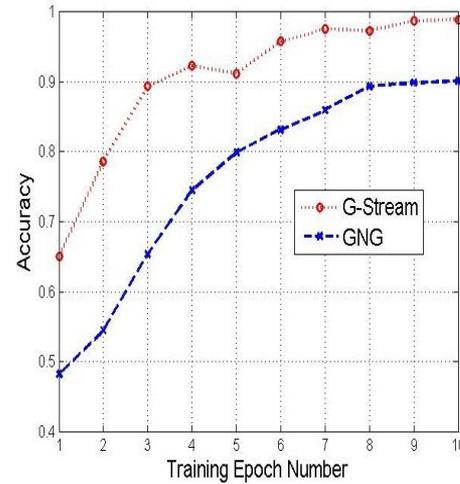
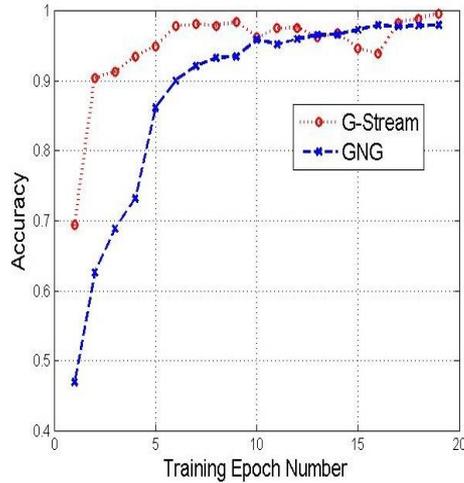


letter4



Sea

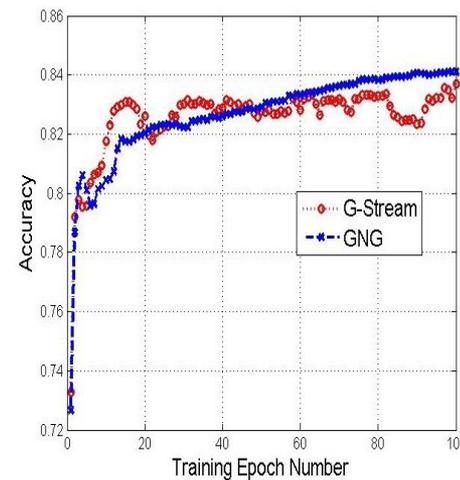
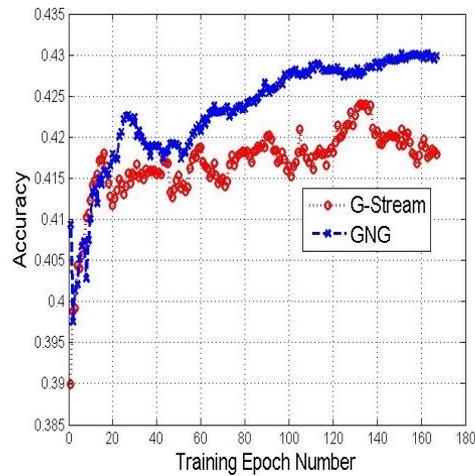
G-STREAM VS GNG ONLINE: ACCURACY



DS1

DS2

Letter4

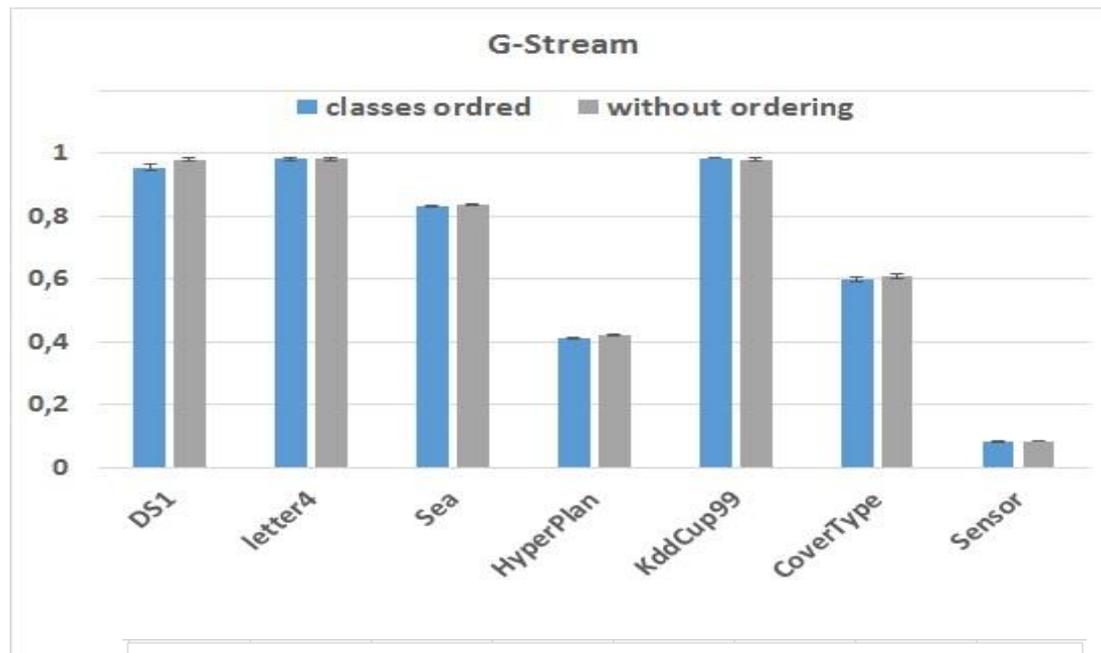


HyperPlan

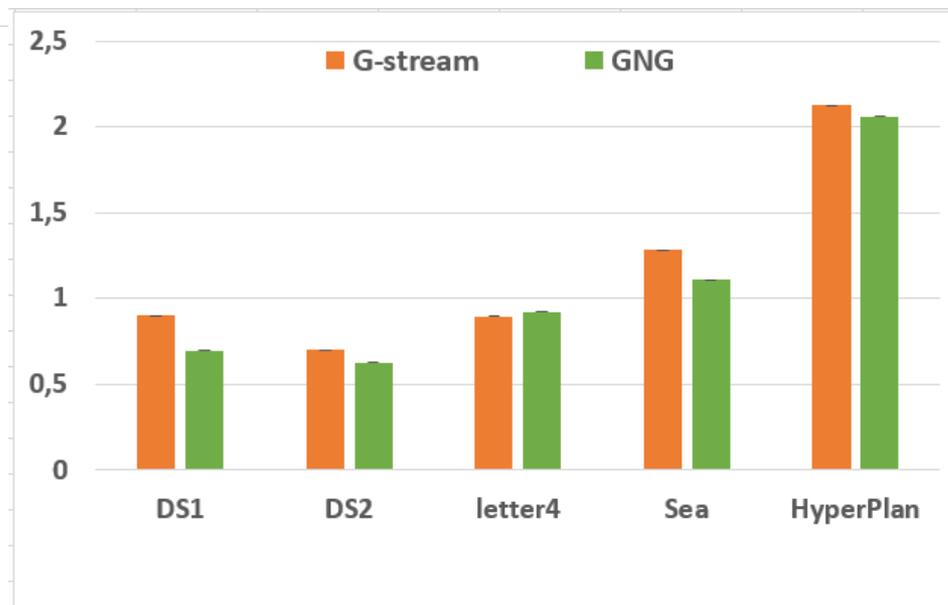
Sea

EVALUATIONS

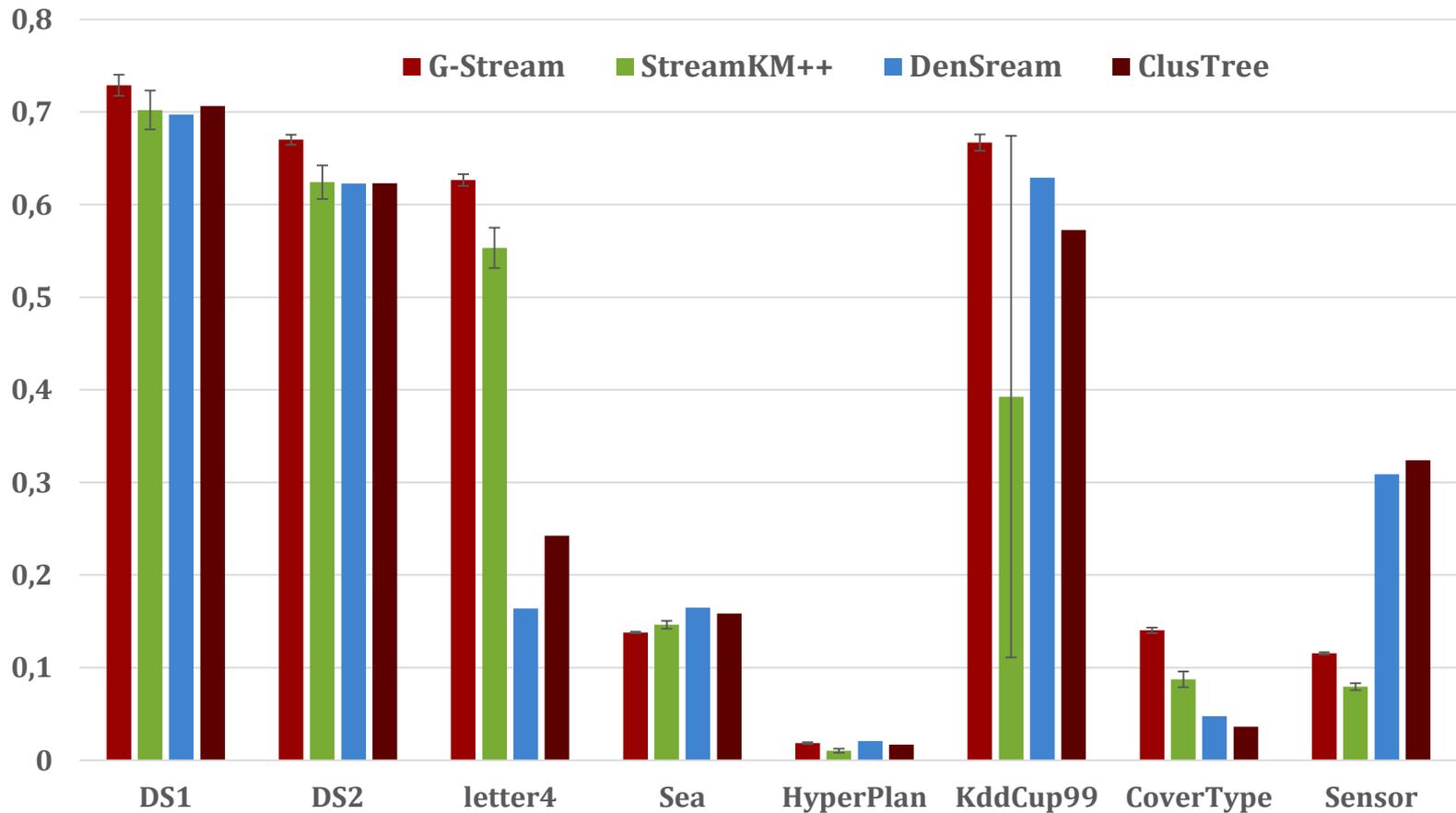
NMI



DB



NMI



CONCLUSION & PERSPECTIVES

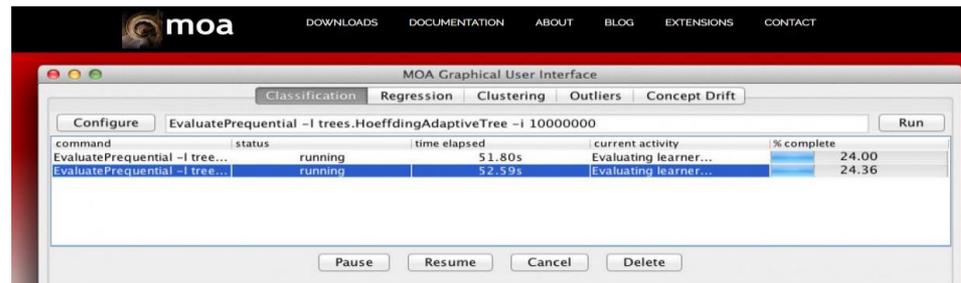
- Une diversité d'approches
 - Approches probabilistes: EM online, OVEM-DyMix (Online Variational Expectation Maximisation for Dynamic Mixture model)
- L'apprentissage en flux de données n'est pas généralement réalisé en un seul passage. C'est un processus évolutif ..!
- L'intérêt de combiner du «batch» et du «streaming»
- Différents outils (dont spark)

OUTILS DISPONIBLES

○ {M}assive {O}nline {A}nalysis MOA (Bifet et al. 2010)

- Connecté à WEKA

<http://moa.cms.waikato.ac.nz/details/stream-clustering/>



○ StreamDM

<http://streamdm.noahlab.com.hk/>



<http://spark.apache.org/>

OUTILS DISPONIBLES

- Apache Mahout



- SAMOA

- G. De Francisci Morales, et al . "SAMOA: Scalable Advanced Massive Online Analysis." Journal of Machine Learning Research, 16(Jan):149–153, 2015.

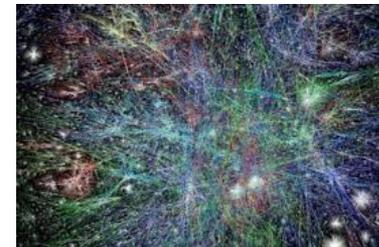
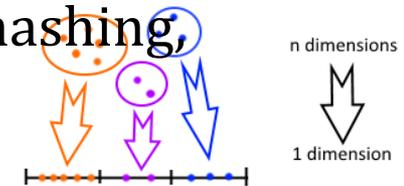
<http://samoa.incubator.apache.org/>



<https://flink.apache.org/>

QUELQUES CHALLENGES

- S'adapter constamment à l'évolution de l'environnement (auto-configurable)
- La prise en compte de l'évolution des données et du recouvrement
- La prise en compte des différents types de données
- Accélération des programmes: l'échantillonnage, hashing, projection aléatoire (LSH)
- La visualisation
- Facile à utiliser « **ease to use** »



RÉFÉRENCES

- G. De Francisci Morales, A. Bifet. "SAMOA: Scalable Advanced Massive Online Analysis." *Journal of Machine Learning Research*, 16(Jan):149–153, 2015.
- E. Sparks, A. Talwalkar, V. Smith, J. Kottalam, X. Pan, J. Gonzalez, M. Franklin, M. I. Jordan, T. Kraska. *MLI: An API for Distributed Machine Learning*. International Conference on Data Mining (ICDM), 2013
- Charu C. Aggarwal, Jiawei Han, Jianyong Wang, and Philip S. Yu. A framework for clustering evolving data streams. In *Proceedings of the International Conference on Very Large Data Bases (VLDB '03)*, pages 81–92, 2003.
- Aggarwal C C, Han J, Wang J, Yu P S. A framework for projected clustering of high dimensional data streams. In *Proc. the 30th International Conference on Very Large Data Bases, Volume 30, Aug. 29-Sept. 3, 2004*, pp.852-863.
- Amini A, Wah TY, Saboohi H. On density-based data streams clustering algorithms: A survey. *JOURNAL OF COMPUTER SCIENCE AND TECHNOLOGY* 29(1): 116–141 Jan. 2014. DOI 10.1007/s11390-013-1416-3
- Jonathan A. Silva, Elaine R. Faria, Rodrigo C. Barros, Eduardo R. Hruschka, André C. P. L. F. de Carvalho, and João Gama. 2013. Data stream clustering: A survey. *ACM Comput. Surv.* 46, 1, Article 13 (July 2013), 31 pages
- Gianmarco De Francisci Morales, Joao Gama, Albert Bifet, Wei Fan. Big Data Stream Mining Tutorial at IEEE BigData 201
- Albert Bifet, André C P L F de Carvalho, João Gama. BigData Stream Mining. ECML-PKDD Summer school 2015
- Hani El Assaad, Allou Samé, Gérard Govaert, Patrice Aknin: Model-Based Clustering of Temporal Data. ICANN 2013: 9-16
- Mohammed Ghesmoune, Mustapha Lebbah, Hanene Azzag. Micro-Batching Growing Neural Gas for Clustering Data Streams using Spark Streaming. *Procedia Computer Science journal* (2015) pp. 158-166. Doi 10.1016/j.procs.2015.07.290.
- Mohammed Ghesmoune, Mustapha Lebbah, and Hanene Azzag. Clustering over data streams based on growing neural gas. *PAKDD* (2) 2015: 134-145.