



# Big Data et Graphes : Quelques pistes de recherche

**Hamamache Kheddouci**

Laboratoire d'InfoRmatique en Image et Systèmes d'information

LIRIS UMR 5205 CNRS/INSA de Lyon/Université Claude Bernard Lyon 1/Université Lumière Lyon 2/Ecole Centrale de Lyon

<http://liris.cnrs.fr>



# Big Data : Grandes Masses de Données



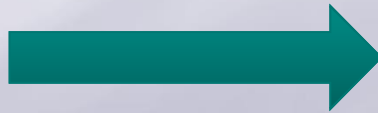
Age du Big Data !



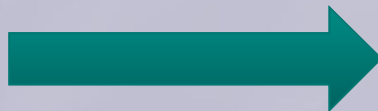
# Nouveau modèle de données

## Le Modèle de Génération/Consommation de la donnée a changé

**Ancien modèle :** Quelques compagnies génèrent des données, les autres sont des consommateurs de données



**Nouveau Modèle :** nous sommes tous des générateurs de données, et nous sommes tous des consommateurs de données



# Générateurs des Big Data



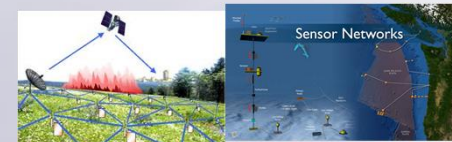
**Instruments scientifiques**  
(collecter toute sorte de données)



**Média et réseaux sociaux**  
(tous des générateurs de données)



**Mobiles**  
(tracer tous les objets tout le temps)



**Réseaux de capteurs**  
(mesurer tout type de données)



# Générateurs des Big Data



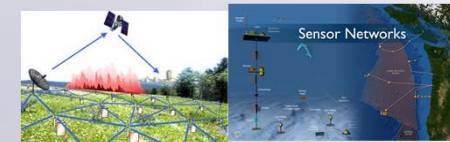
**Instruments scientifiques**  
(collecter toute sorte de données)



**Média et réseaux sociaux**  
(tous des générateurs de données)



**Mobiles**  
(tracer tous les objets tout le temps)



**Réseaux de capteurs**  
(mesurer tout type de données)



# Age du Big Data

**Í Data is a new class of economic asset, like currency and gold.Î**

*Source: World Economic Forum 2012*





# Big Data

☰ Un enjeu **scientifique** important :



# Big Data

## ☰ Définitions Å

Í **Big Data** is a massive volume of both structured and unstructured data that is so large that it's difficult to process with traditional database and software techniques.Î

Í **Big Data** is data whose scale, diversity, and complexity require new architectures, models, techniques, algorithms, and analytics to manage it and extract value and hidden knowledge from itÅ

**Avec quels modèles ?**





# Big Data & Graphs ?

Naturellement, les graphes et les données sont liés :

- **Linked open Data (graphe d'interaction entre données)**
- **Des objets du Web sont des graphes (XML, RDF, Å )**
- **Graphes des amis de Facebook**
- **Graphe de connaissances de Google**
- **Graphes extraits de grandes base de données**

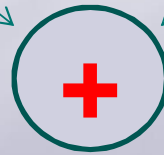
Base de données

| Données | temps | Emetteur | Récepteur | Type de mess. | .... | Attribut n. |
|---------|-------|----------|-----------|---------------|------|-------------|
| D1      | 1     | S1       | S2        | A             | .... | 3444        |
| D2      | 2     | S1       | S3        | C             | .... | 2112        |
| D3      | 3     | S2       | S4        | B             | .... | 5858        |
| D4      | 4     | S4       | S2        | A             | .... | 600         |
| D5      | 5     | S3       | S5        | C             | .... | 2333        |
| ....    | ....  | ....     | ....      | ....          | .... | ....        |

# Big Data & Big Graphs

**Big Data**

**Big Graphs**



**Big Data Graphs**



# Big Graphs for Big Data

| Verrous Big Data           | Solutions à base de graphes |
|----------------------------|-----------------------------|
| Indexation et stockage     | Partitionnement de graphes  |
| Flux de données (Vélocité) | Analyse de flux de graphes  |
| Visualisation des données  | Visualisation de graphes    |



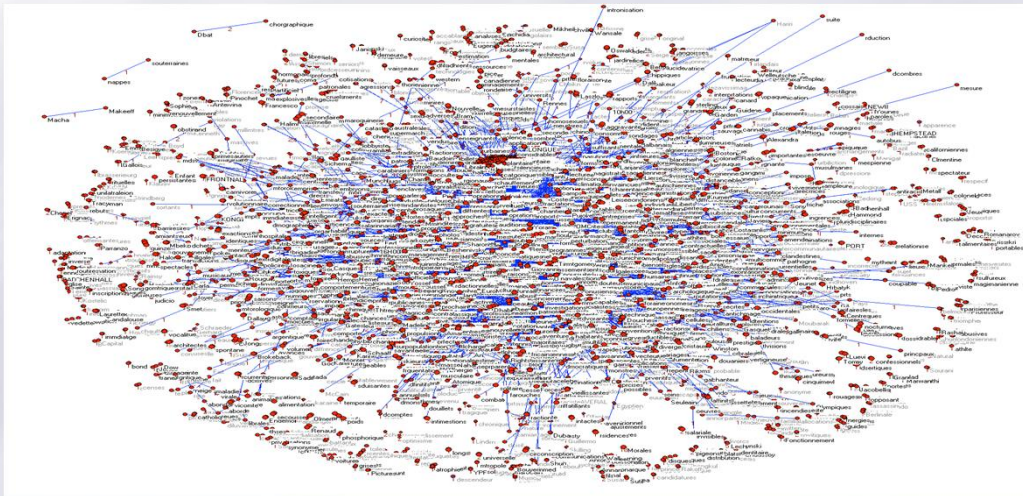
# Big Graphs for Big Data

| Verrous Big Data              | Solution à base de graphes        |
|-------------------------------|-----------------------------------|
| <b>Indexation et stockage</b> | <b>Partitionnement de graphes</b> |
| Analyse de flux de données    | Analyse de flux de graphes        |
| Visualisation des données     | Visualisation de graphes          |



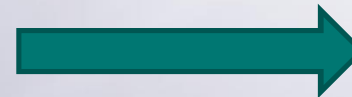
# Big Graphs for Big Data

## Partitionnement de graphes de données



**Big Graphs**  
(Milliards de nœuds et arêtes)

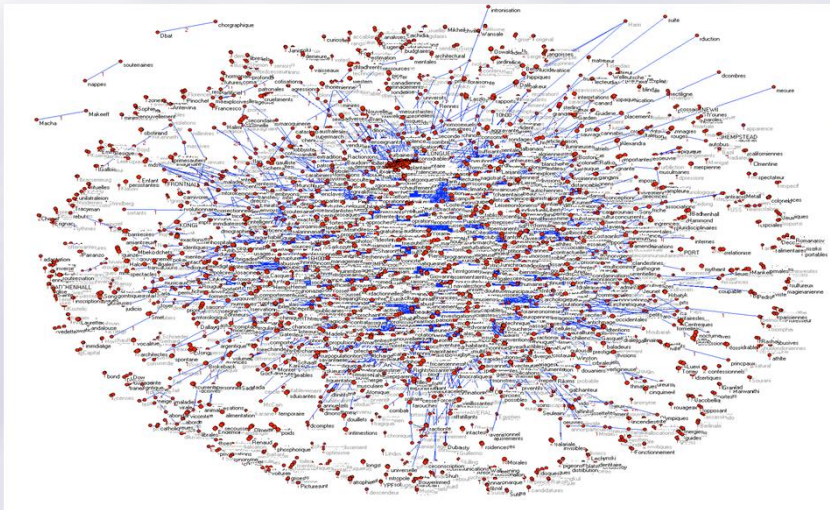
possible ?



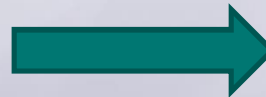
1 machine ?  
Combien de :  
- RAM ?  
- Disque ?

# Big Graphs for Big Data

## Partitionnement de graphes de données



possible ?



k machines

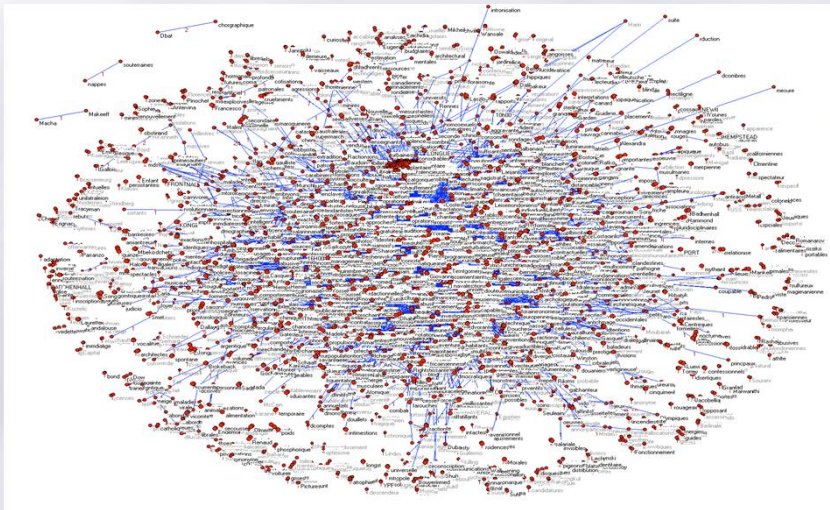
**Big Graphs**  
(Milliards de nœuds et arêtes)



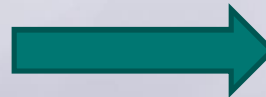


# Big Graphs for Big Data

## Partitionnement de graphes de données



possible ?



OUI



k machines

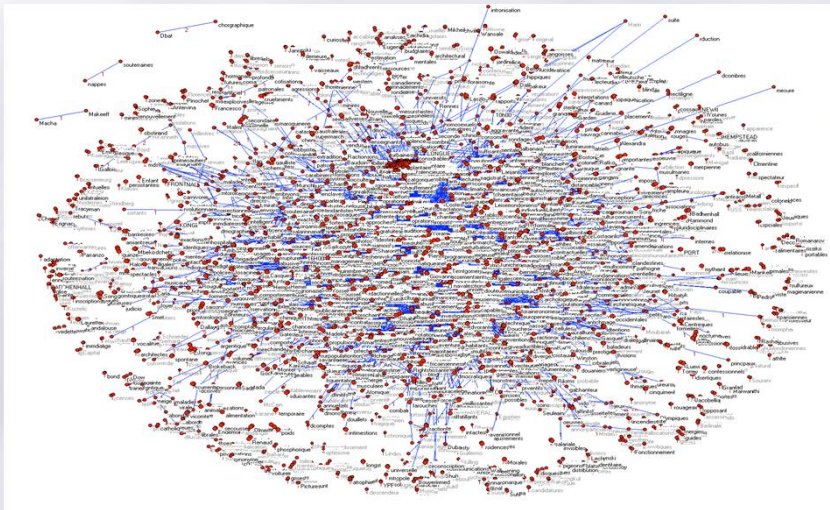
**Big Graphs**  
(Milliards de nœuds et arêtes)



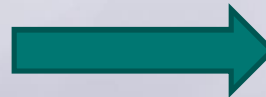


# Big Graphs for Big Data

## Partitionnement de graphes de données



possible ?



OUI



k clusters

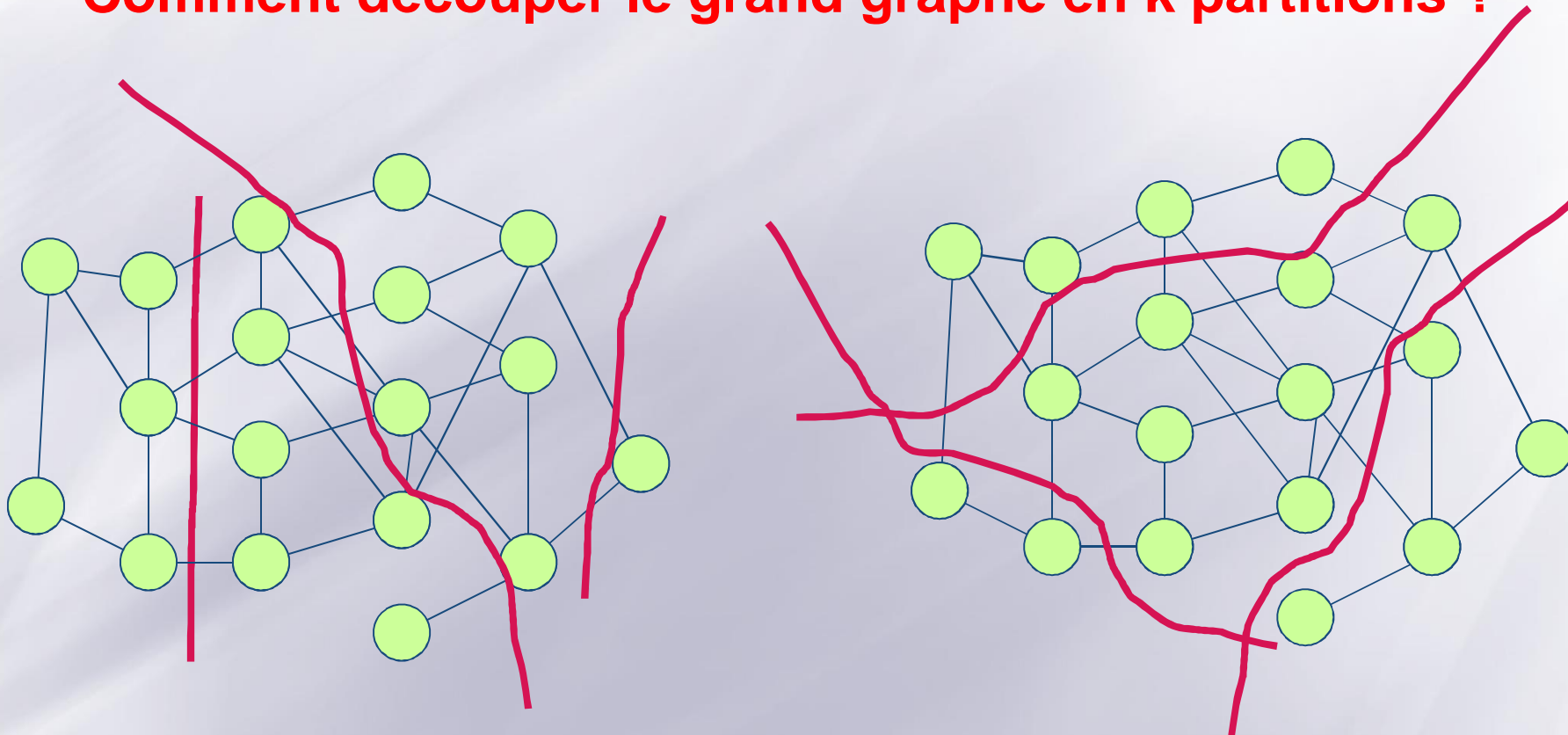
**Big Graphs**  
(Milliards de nœuds et arêtes)

**Comment découper le grand graphe ?**

# Big Graphs for Big Data

## ☰ Partitionnement de graphes de données

**Comment découper le grand graphe en k partitions ?**



# Big Graphs for Big Data

## Partitionnement de graphes de données

Etant donné un graphe  $G = (N, E, W_N, W_E)$

- $N$  = sommets,
- $W_N$  = poids sur les sommets
- $E$  = arêtes
- $W_E$  = poids des arêtes

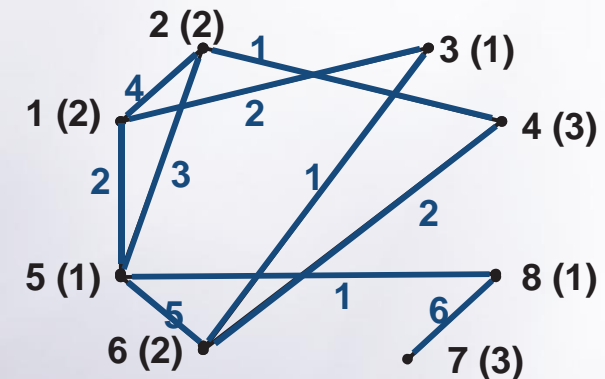
Ex:  $N = \{\text{données}\}$ ,  $W_N = \{\text{vecteurs de attributs de données}\}$ ,  
arête  $(j,k)$  dans  $E$  :  $j$  envoie  $W_E(j,k)$  mots au  $k$

Choisir une partition  $N = N_1 \cup N_2 \cup \dots \cup N_p$  telle que

- La somme des poids des nœuds dans chaque  $N_j$  est à presque le même
- La somme des poids des arêtes connectant toutes les différentes paires  $N_j$  et  $N_k$  est minimisée

Ex: équilibrage des chargements de données, en minimisant la communication entre les machines

Cas particulier,  $N = N_1 \cup N_2$



# Big Graphs for Big Data

## Partitionnement de graphes de données

Etant donné un graphe  $G = (N, E, W_N, W_E)$

- $N$  = sommets,
- $W_N$  = poids sur les sommets
- $E$  = arêtes
- $W_E$  = poids des arêtes

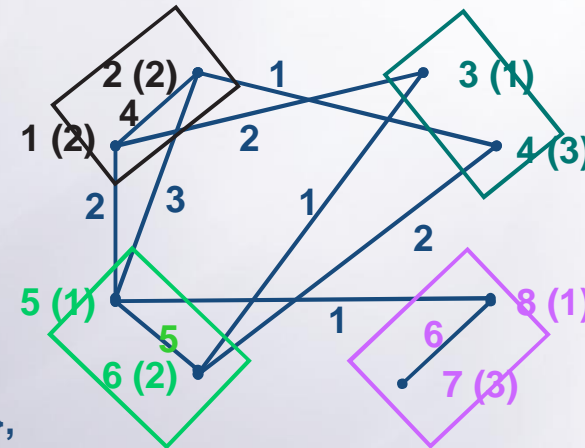
Ex:  $N = \{\text{données}\}$ ,  $W_N = \{\text{vecteurs d'attributs de données}\}$ ,  
arête  $(j,k)$  dans  $E$  :  $j$  envoie  $W_E(j,k)$  mots au  $k$

Choisir une partition  $N = N_1 \cup N_2 \cup \dots \cup N_p$  telle que

- La somme des poids des nœuds dans chaque  $N_j$  est à presque le même
- La somme des poids des arêtes connectant toutes les différentes paires  $N_j$  et  $N_k$  est minimisée

Ex: équilibrage des chargements de données, en minimisant la communication entre les machines

Cas particulier,  $N = N_1 \cup N_2$



# Big Graphs for Big Data

## Partitionnement de graphes de données

 NP-complet

 Plusieurs algorithmes existent :

- partitionnement spectral
- partitionnement géométrique
- partitionnement en graphes Multi-niveaux

# Big Graphs for Big Data

## Partitionnement de graphes de données

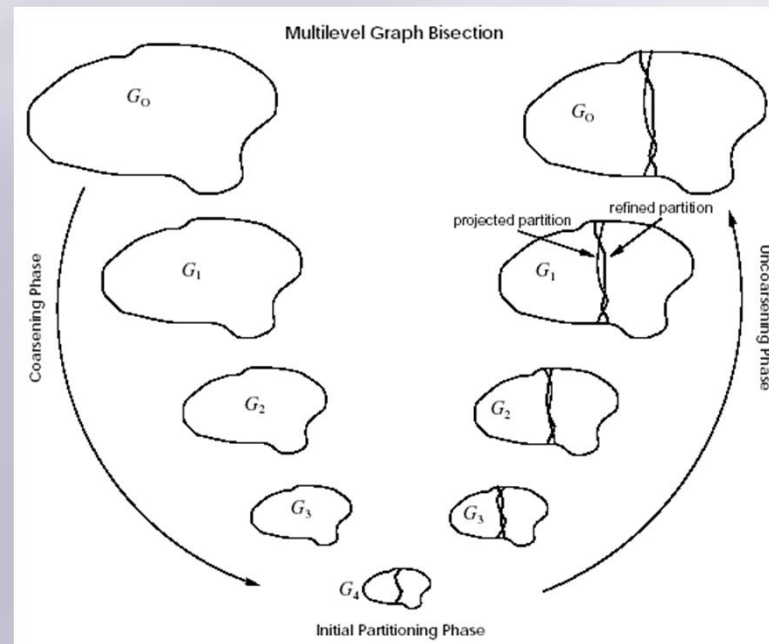
NP-complet

Plusieurs algorithmes existent :

- partitionnement spectral
- partitionnement géométrique
- partitionnement en graphes Multi-niveaux

3 Phases

- compresser
- Partitionner
- Décompresser

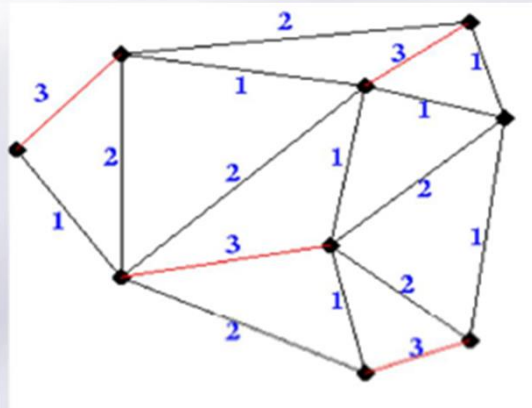




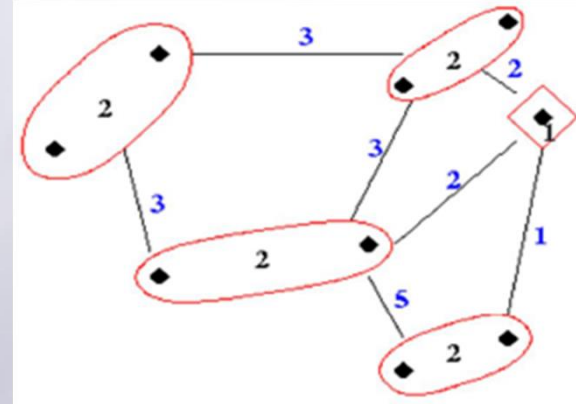
# Big Graphs for Big Data

## Partitionnement de graphes de données

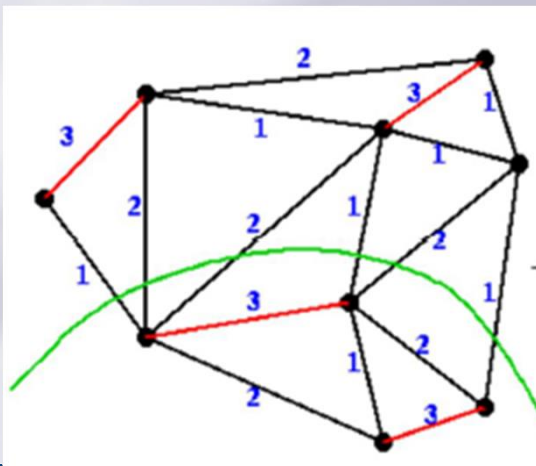
### Partitionnement en graphes Multi-niveaux



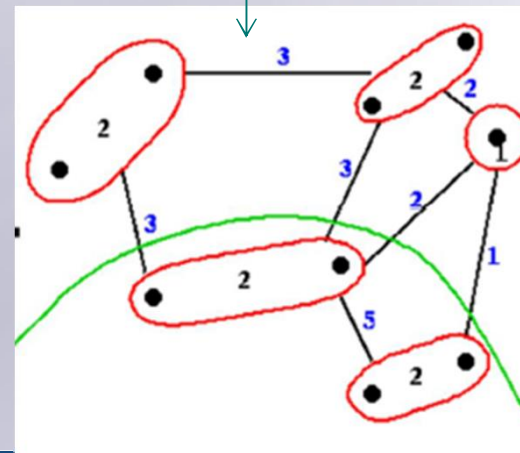
Compresser  
→  
(couplage maximum)



partitionnement  
équilibré



décompresser





# Big Graphs for Big Data

## Partitionnement de graphes de données

### Quelques paramètres de graphes liés au partitionnement

- Ensemble d'arêtes Séparateur :  $E_s$  (sous-ensemble de  $E$ ) sépare  $G$  si en retirant  $E_s$  de  $E$  donne 2 composantes connexes de tailles égales,  $N: N_1$  and  $N_2$
- Ensemble de sommets Séparateur :  $N_s$  (sous-ensemble de  $N$ ) sépare  $G$  si en retirant  $N_s$  et toutes leurs arêtes incidentes donne 2 composantes connexes de tailles égales,  $N: N_1$  and  $N_2$

$G = (N, E)$ , sommets  $N$  et arêtes  $E$   
 $N_s =$  sommets verts



# Big Graphs for Big Data

| Verrous Big Data           | Solution à base de graphes |
|----------------------------|----------------------------|
| Indexation et stockage     | Partitionnement de graphes |
| Analyse de flux de données | Analyse de flux de graphes |
| Visualisation des données  | Visualisation de graphes   |



# Big Graphs for Big Data

## ☰ Analyse de graphes de données en flux

### Flux de données :

- ☰ Un **flux de données** est une séquence de données :  $a_1, a_2, \dots, a_n$ .
  - Flux de prix
  - Flux de paquets IP
- ☰ Les données ont différentes formes dans différentes applications.
  - Valeur scalaire
  - Tuple
  - $\mathbb{A}$
- ☰ La sémantique des données est également différente dans différentes applications.

# Big Graphs for Big Data

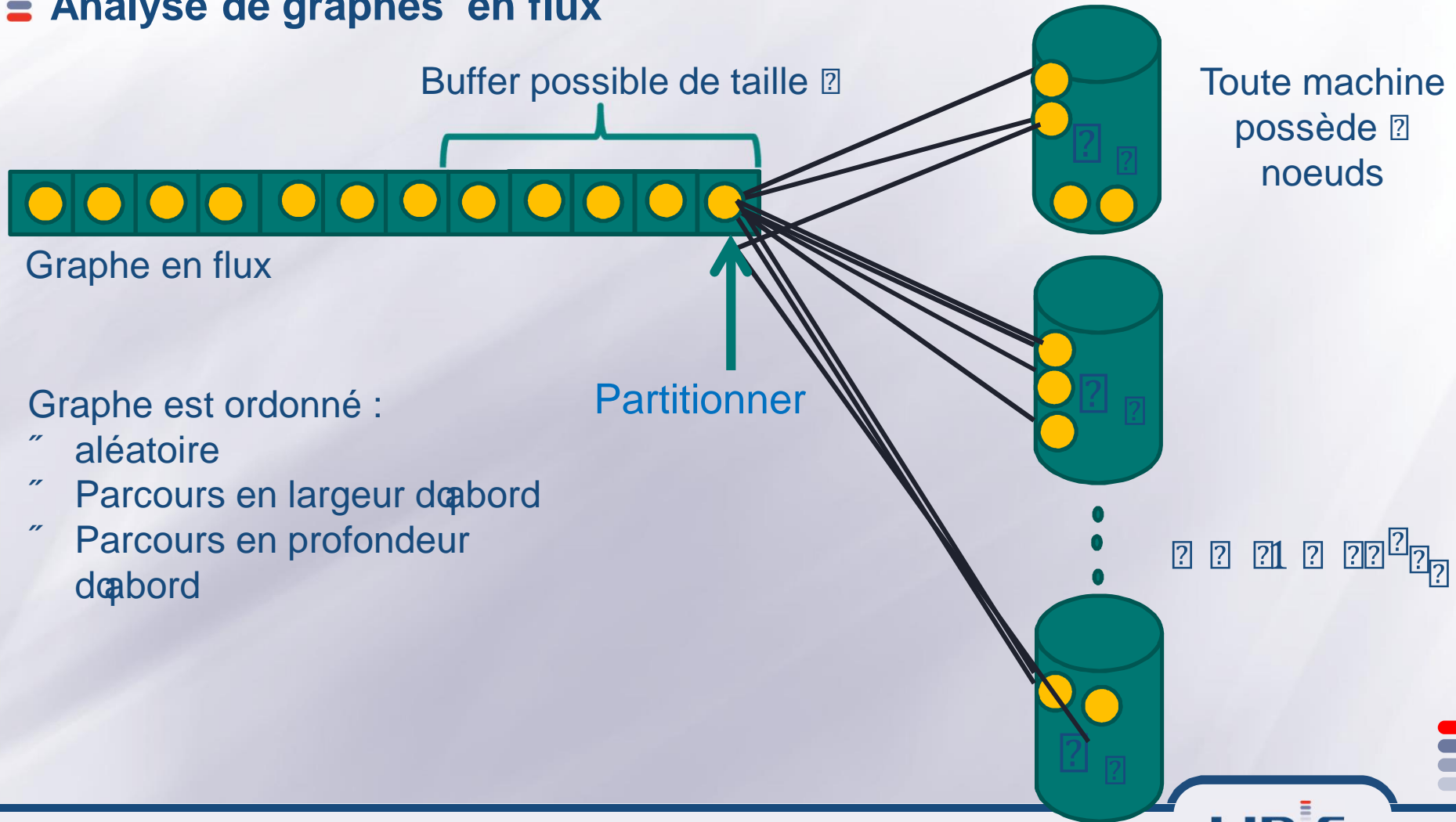
## ☰ Analyse de graphes en flux

### Modèle de traitement de flux :

- ☰ **Accès séquentiel** au flux de données
- ☰ **Ordre des données** dans le flux n'est pas contrôlé par l'algorithme et peut être artificiel.
- ☰ **Petit espace de travail** comparé à la longueur du flux  $n$  :
  - Polylog  $n$
  - $n^\epsilon$
- ☰ **Petit nombre de passes** sur le flux :
  - Une passe
  - Un nombre constant de passes
- ☰ **Temps de traitement d'une donnée est court**

# Big Graphs for Big Data

## ☰ Analyse de graphes en flux



- Graphes ordonnés :
- " aléatoire
  - " Parcours en largeur d'abord
  - " Parcours en profondeur d'abord

# Big Graphs for Big Data

## ☰ Analyse de graphes en flux

### ☰ Des travaux existent :

- Tester la connectivité,
- Tester la planarité,
- construction d'arbre couvrant,
- Å

### ☰ Re-penser les **problèmes de graphes classiques** dans le **modèle streaming**, notamment pour l'organisation des données : clustering, détection de propriétés structurelles, etc

# Big Graphs for Big Data

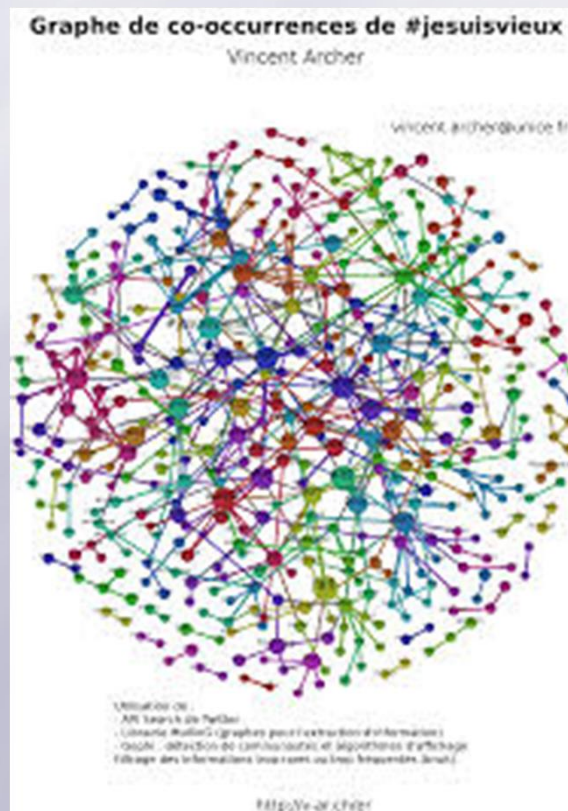
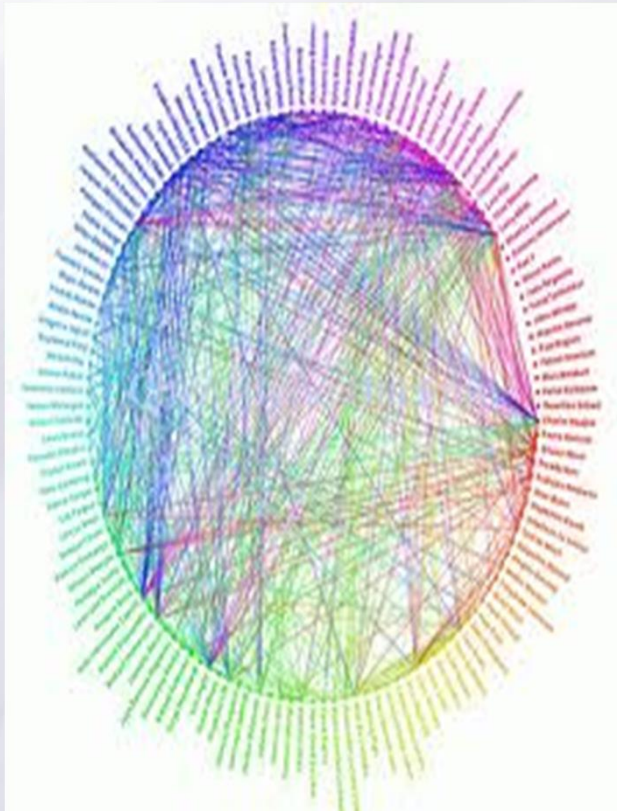
| Verrous Big Data           | Solution à base de graphes |
|----------------------------|----------------------------|
| Indexation et stockage     | Partitionnement de graphes |
| Analyse de flux de données | Analyse de flux de graphes |
| Visualisation des données  | Visualisation de graphes   |





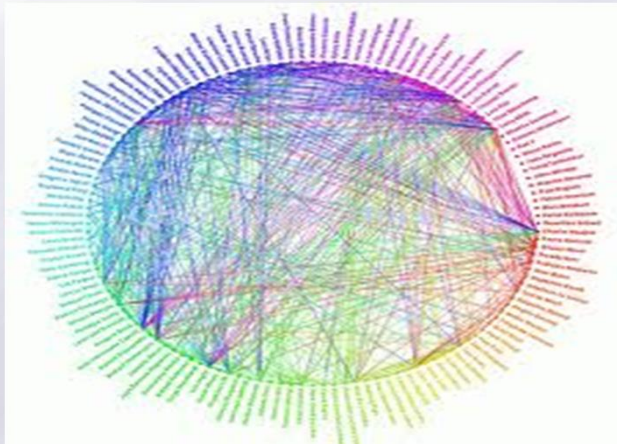
# Big Graphs for Big Data

## Visualisation des grands graphes



# Big Graphs for Big Data

## Visualisation de grands graphes



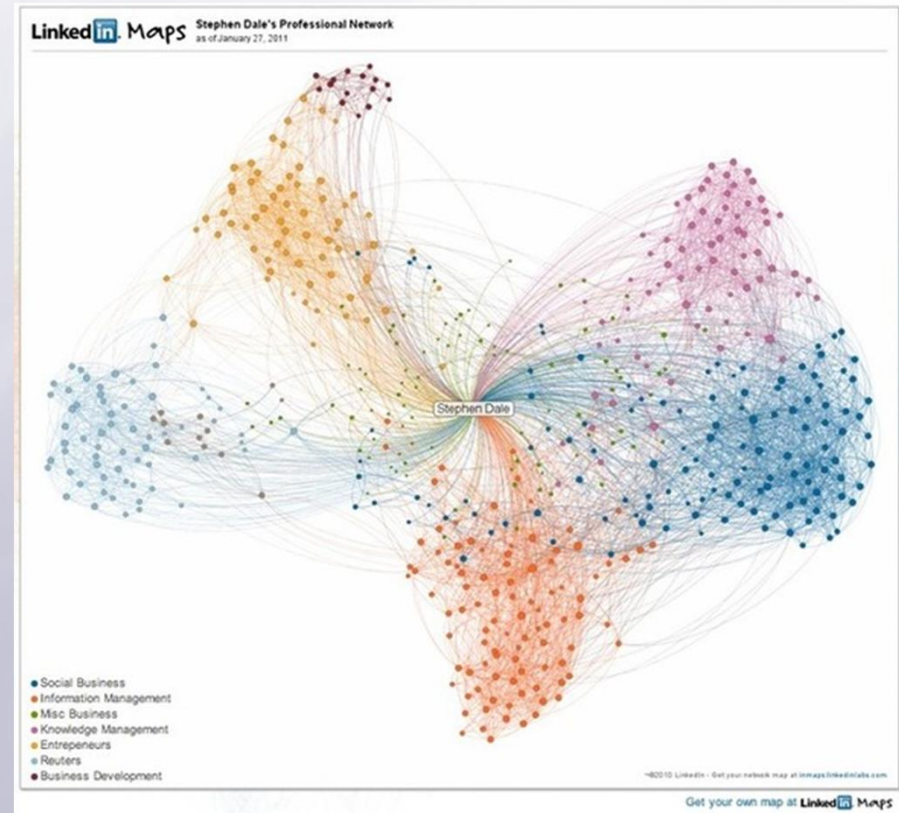
- “ Comment visualiser les grands graphes de données ?
- “ Quels algorithmes et techniques pour explorer visuellement le grand graphe ?



# Big Graphs for Big Data

## Visualisation de grands graphes

- “ Communautés/clustering/classification
- “ Recherche de motifs fréquents
- “ Visualisation de ~~de~~chantillons représentatifs, de sous-graphes, etc
- “ Visualisation 2D, 3D
- “ Combiner la fouille visuelle avec ~~de~~exploration algorithmique des grands graphes.



**Big Data et Graphes,**

**Ce n'est que le début de l'histoire . !!**