A complex network graph composed of numerous nodes (circles) of varying sizes and colors (ranging from light green to dark blue), connected by a dense web of lines representing relationships. The graph is set against a white background with a subtle radial gradient.

Graph & Network Analytics

The power of relations

Vandy BERTEN
Smals - Section Recherche
20 février 2019 - UMONS

Table des matières

Définir un réseau

Caractériser un réseau

Visualiser un réseau

Interroger un réseau





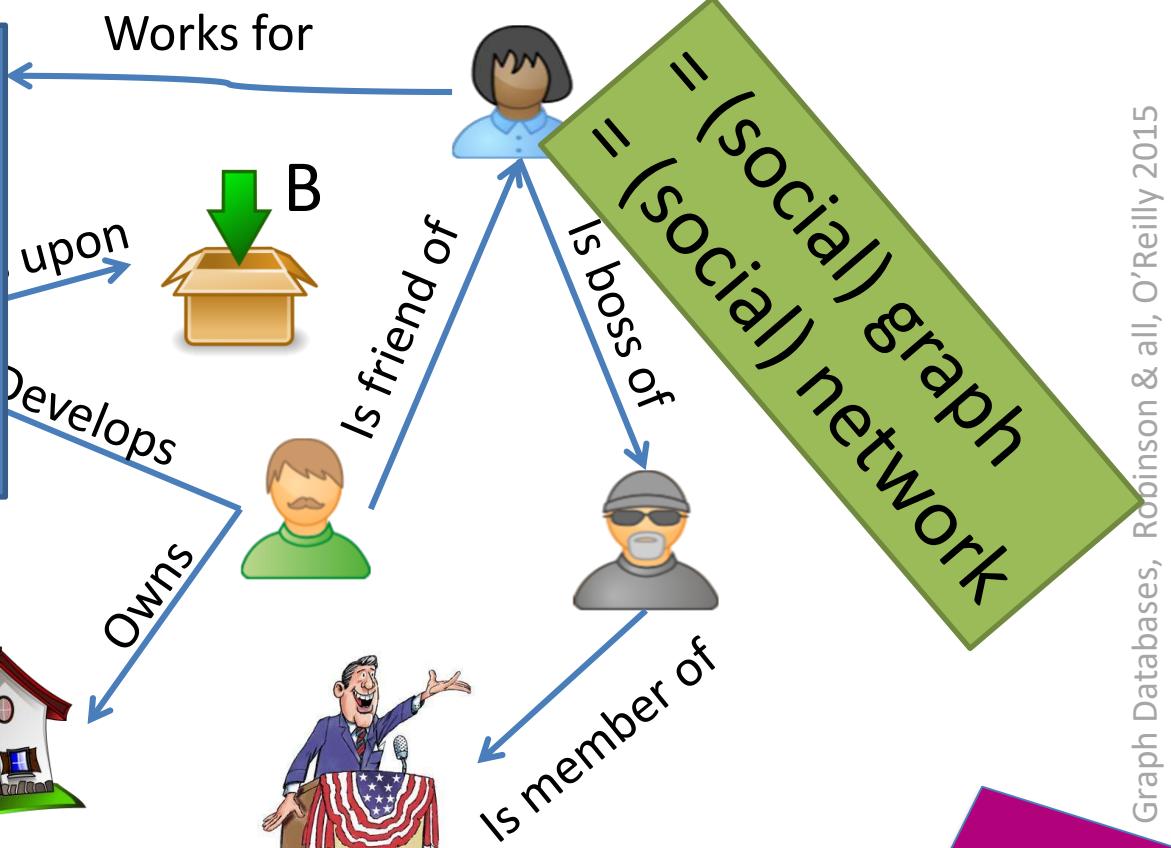
Définir un réseau

Définir un réseau

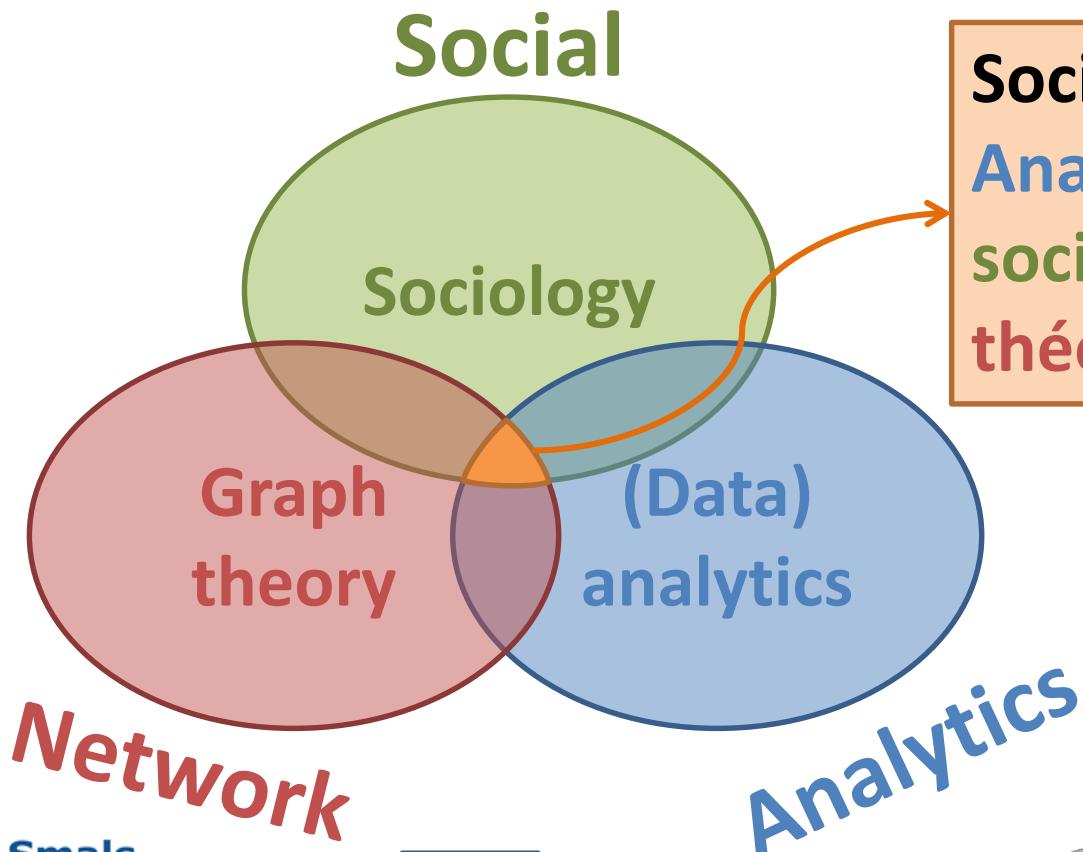
MOTIVATIONS

Motivations

"Facebook [...] was founded on the idea that while there's value in discrete information about people—their names, what they do, etc.—there's even more value in the relationships between them."

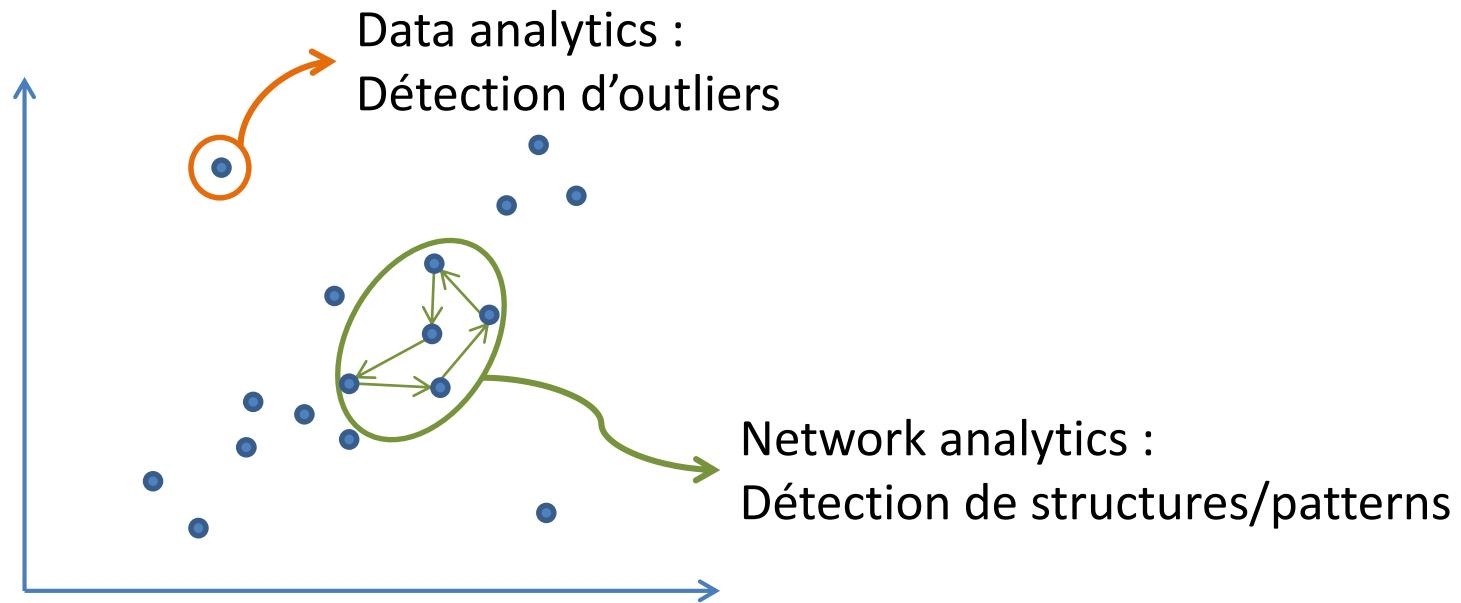


Social Network Analytics



Social Network Analytics :
Analyse des **structures sociales** au travers de la
théorie des graphes

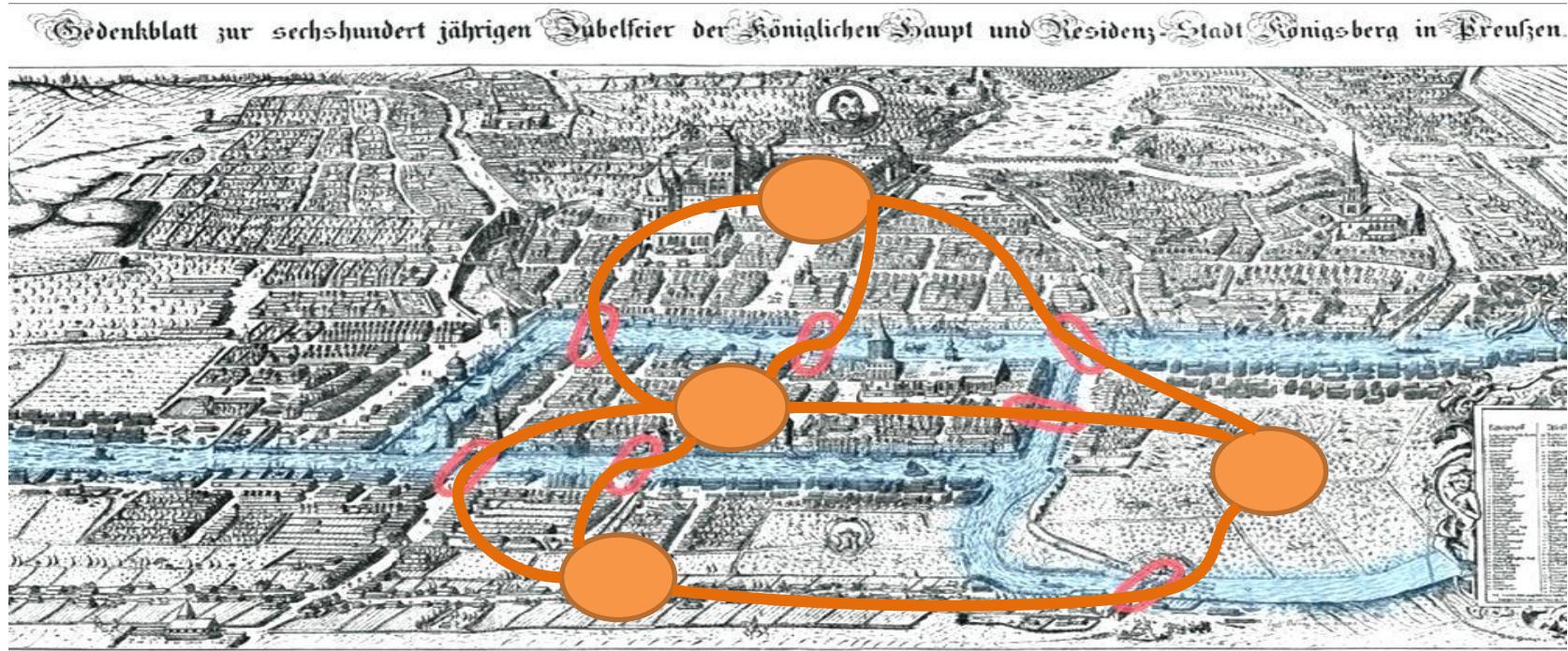
Avantages de Network Analytics



(Social) Graph/Network Analytics is not...



7 ponts de Königsberg (Euler, 1736)



Gartner : applications



Route optimization



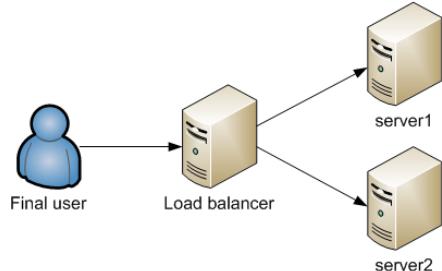
Epidemiology



Intelligence



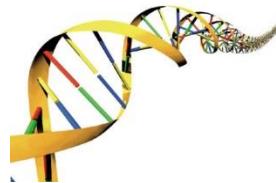
Fraud detection



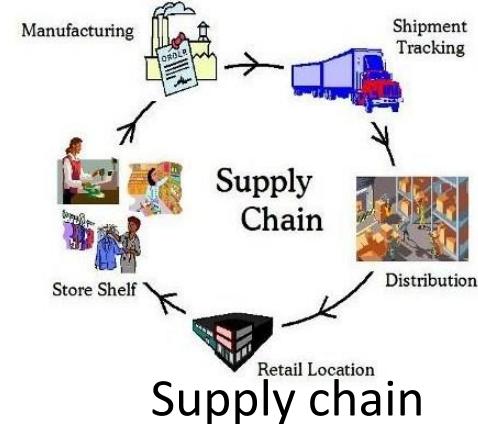
Load balancing



Market basket analysis



Genome

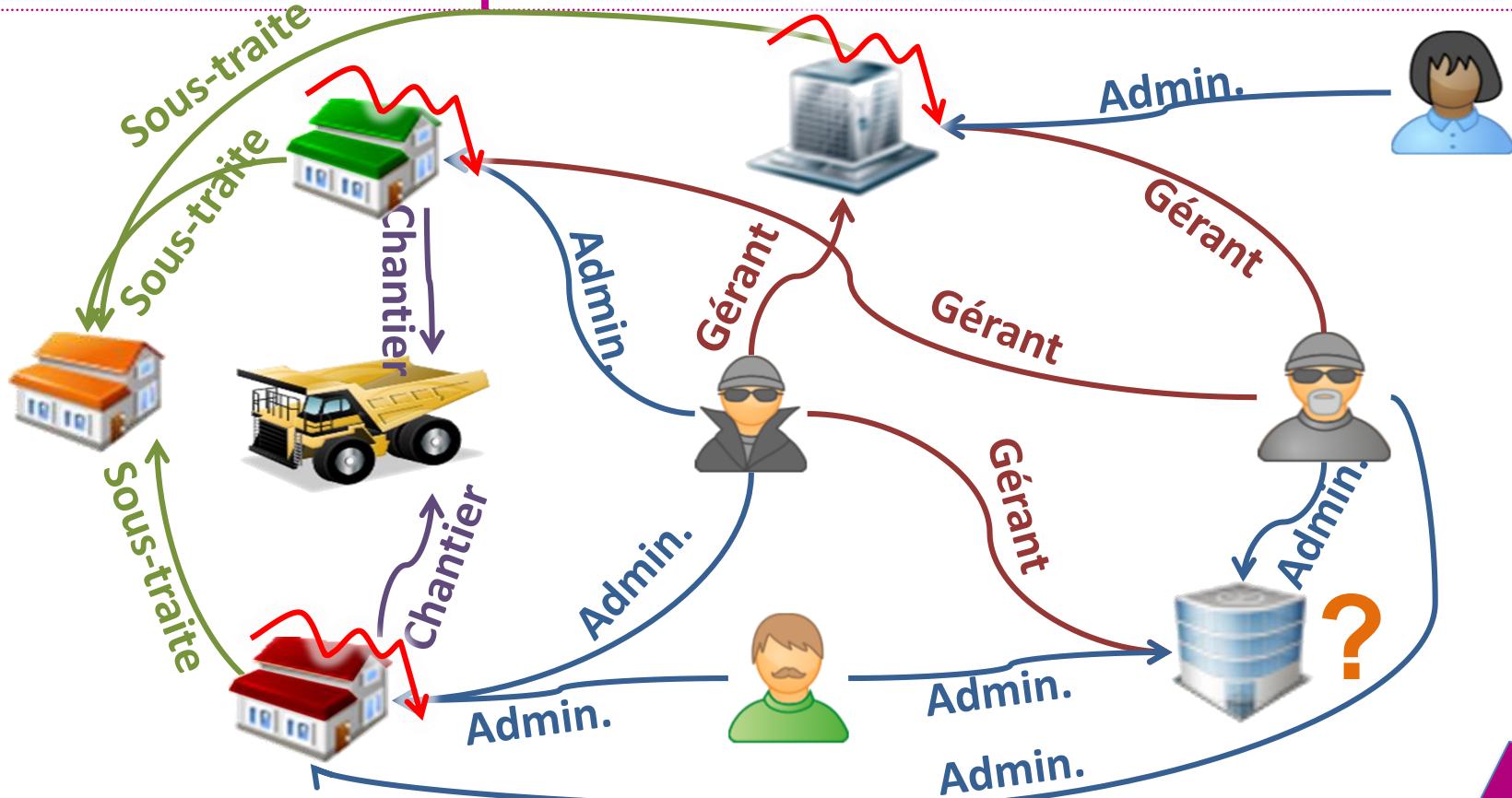


Money laundering

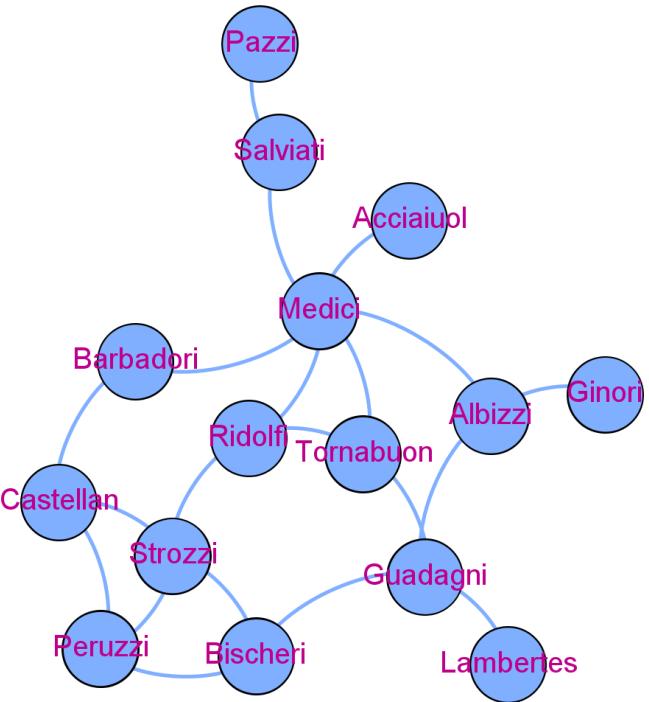
Définir un réseau

EXEMPLES

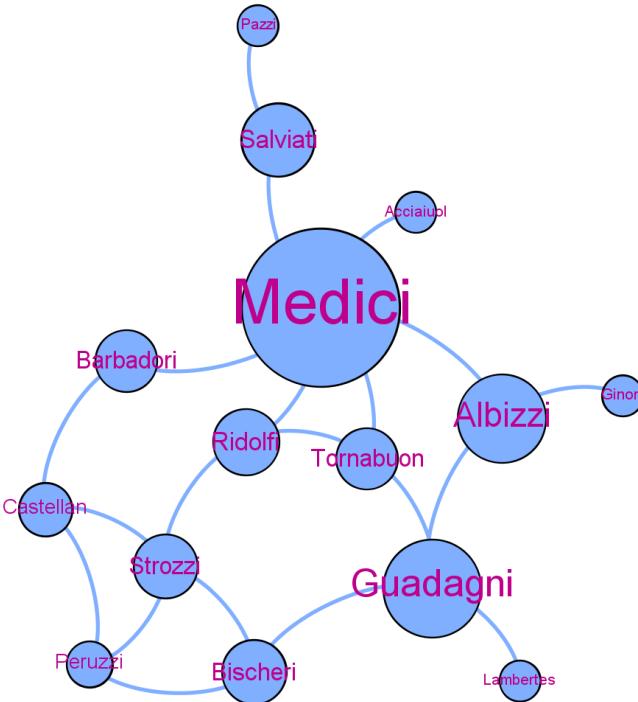
Spider construction



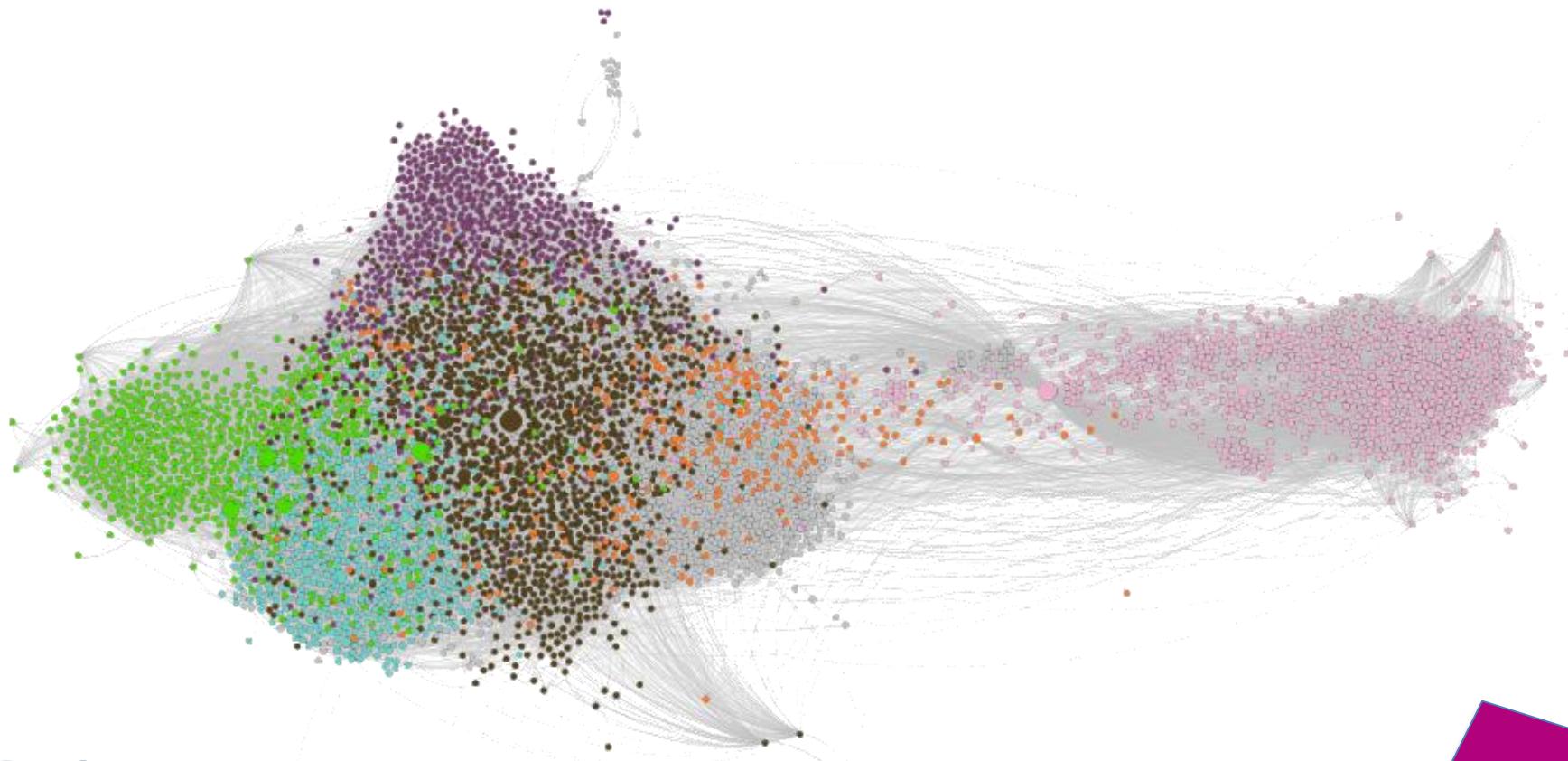
Mariages à Florence



Mariages à Florence



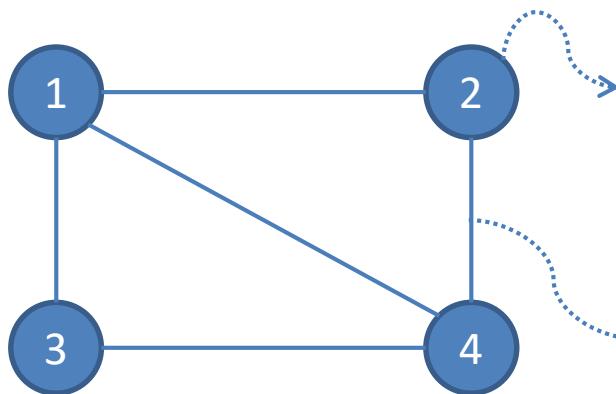
Dépendance de librairies Python



Définir un réseau

DÉFINITIONS

Nœuds, arcs

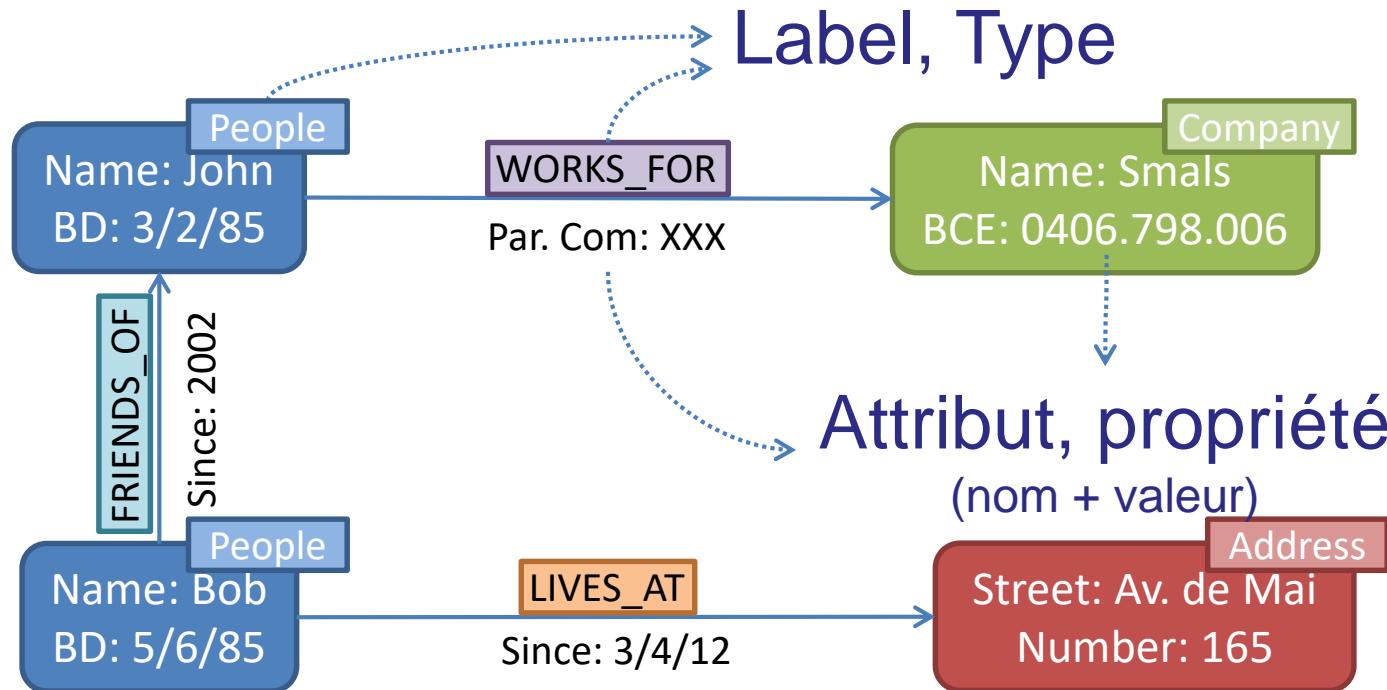


Nœud, sommet,
entité, *vertex*

Arc, relation,
lien, connexion,
edge

$$G = (V, E)$$

Labeled Property Graph



Réseau dirigé/non-dirigé

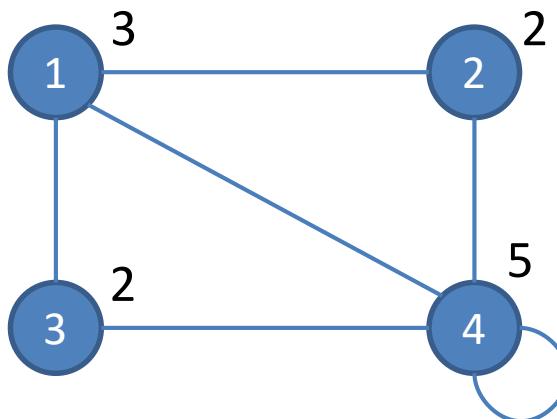


Non dirigé (*undirected*) : collègue, ami Facebook



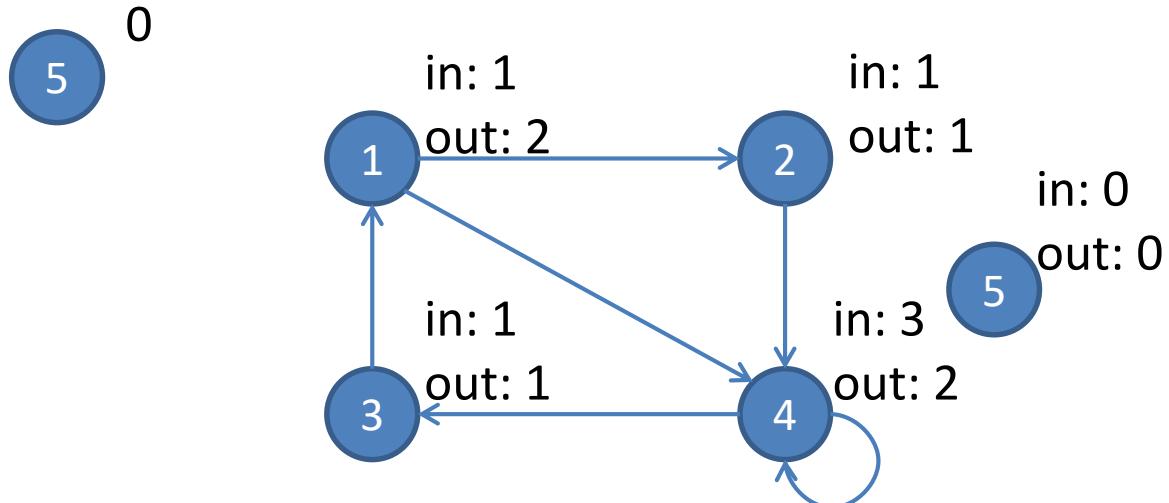
Dirigé (*directed*) : travaille pour, dépend de, follower Twitter

Degré



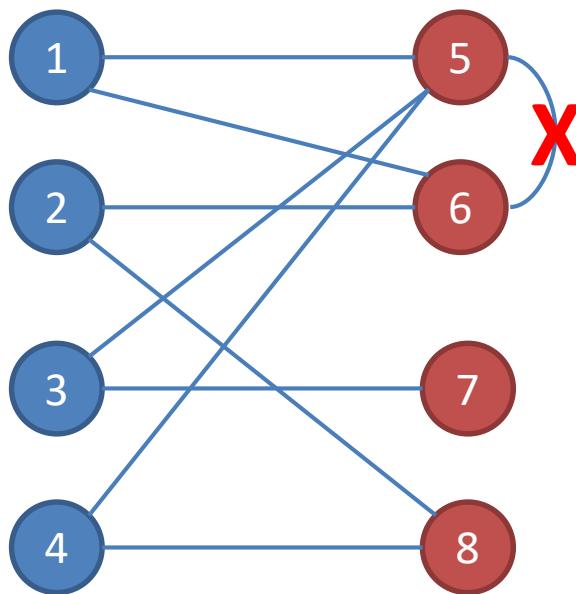
$$\sum_{v \in V} \deg(v) = 2 |E|$$

$$\forall v \in V \deg(v) = \deg^+(v) + \deg^-(v)$$



$$\sum_{v \in V} \deg^+(v) = \sum_{v \in V} \deg^-(v)$$

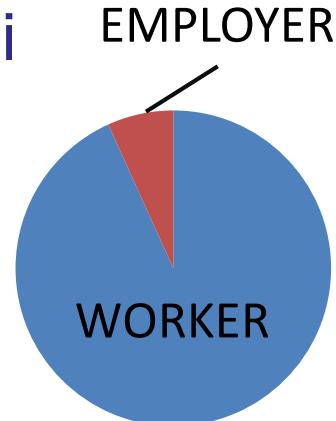
Graphe biparti (*bipartite graph*)



Exemple : Dimona



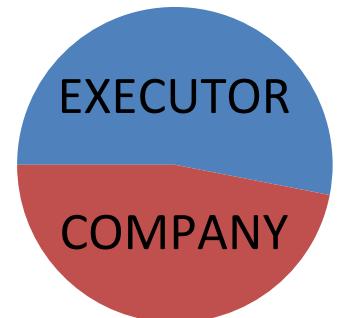
- Déclarations Dimona de 2007 à aujourd'hui
- 7 048 178 nœuds, 20 640 727 relations
- Graphe biparti
- Degré d'un nœud :
 - EMPLOYER : nombre de ses employés
 - WORKER : nombre de ses employeurs



Exemple : BCE/KBO



- Fonctions BCE de 2007 à aujourd'hui
- 3 572 876 nœuds, 3 244 227 relations
- Graphe biparti
- Degré d'un nœud :
 - EXECUTOR : nombre de compagnies où il officie
 - COMPANY : nombre de ses « executors » (=mandataires)



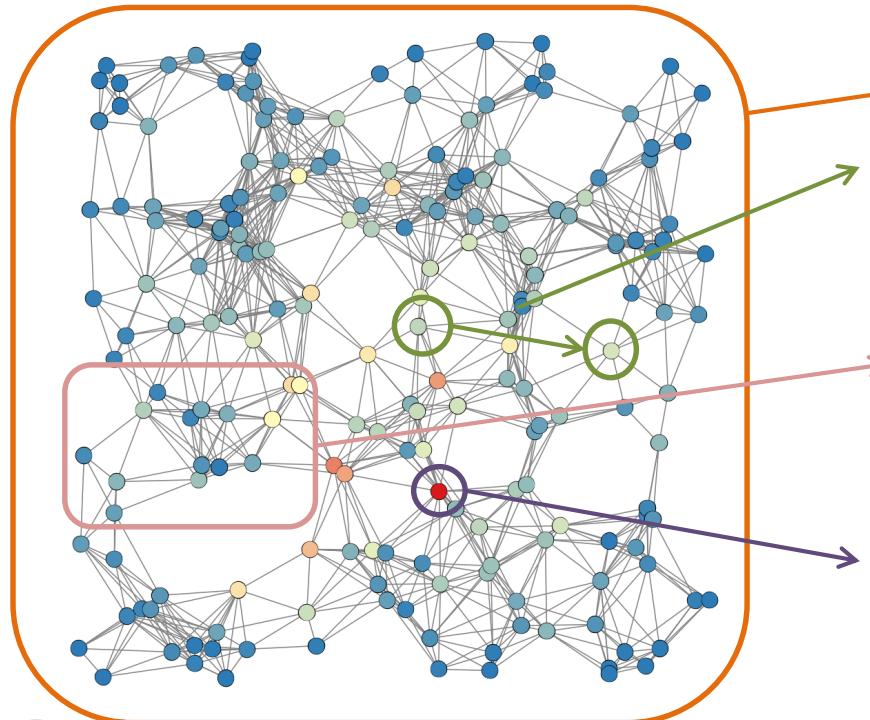


Caractériser un réseau

Caractériser un réseau

Le réseau dans son ensemble
(Métriques générales)

On peut caractériser :



La relation/distance/similarité
entre deux nœuds

Un groupe de nœuds (Clustering)

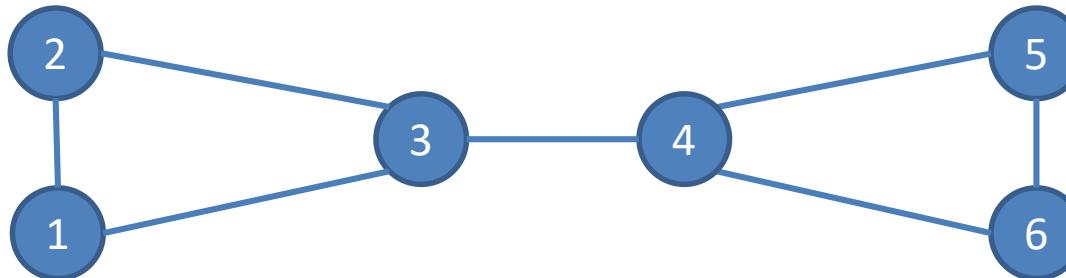
Un nœud en particulier (Centralité)

Caractériser un réseau

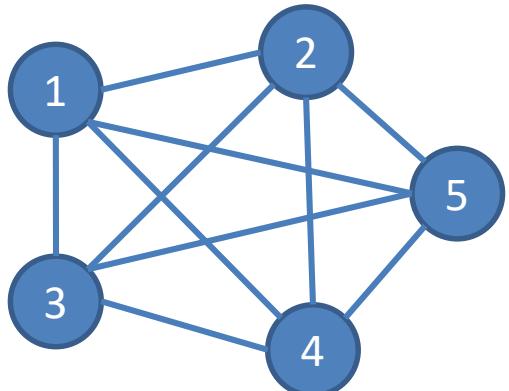
MÉTRIQUES GÉNÉRALES

Diamètre

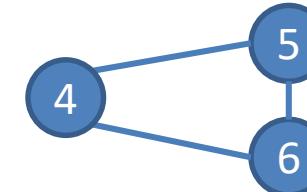
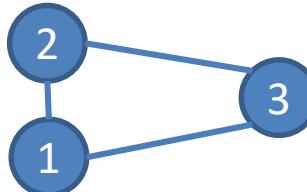
- Diamètre = plus long « plus court chemin »



→ 3

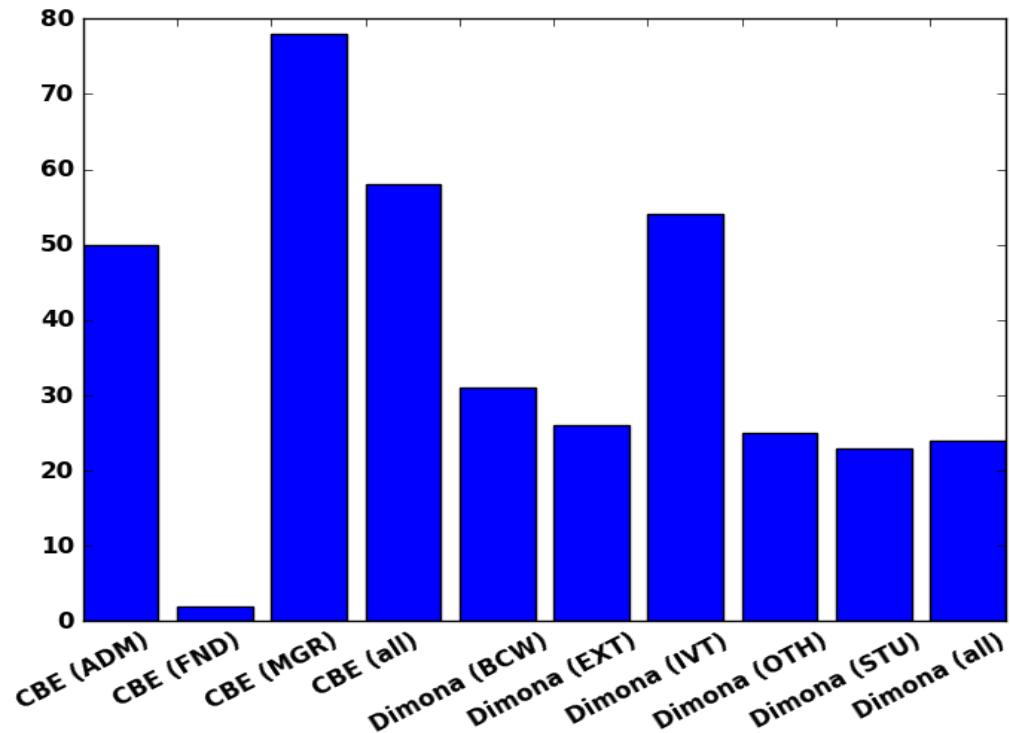


→ 1

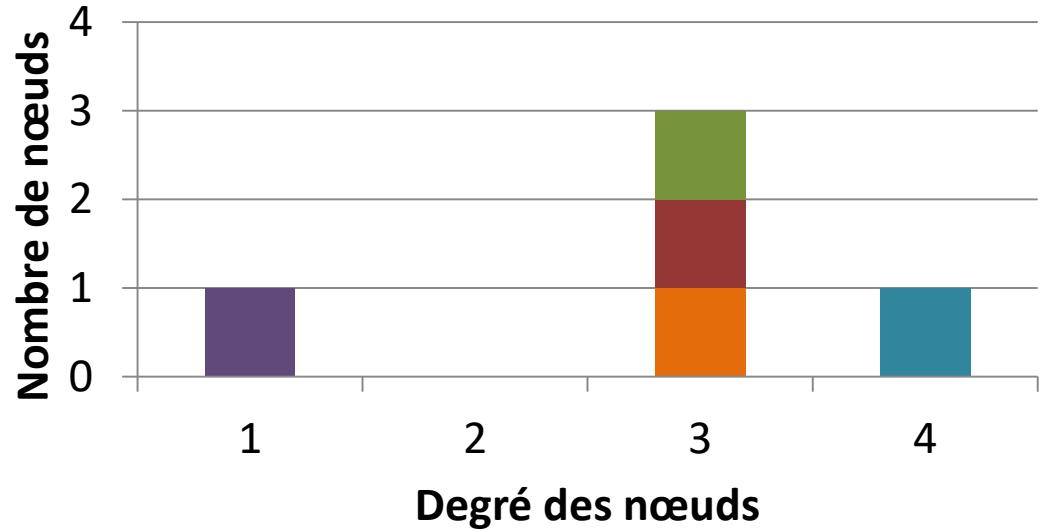
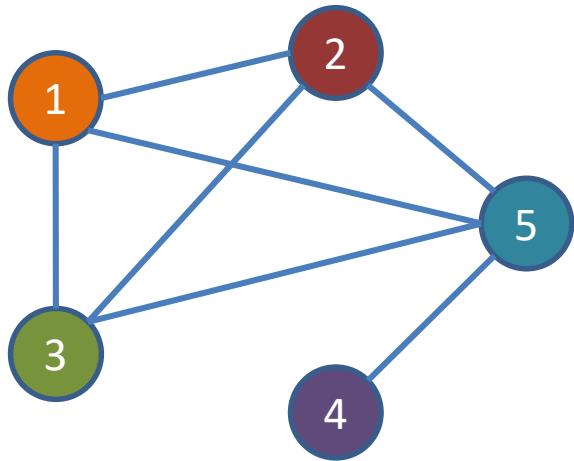


→ ∞

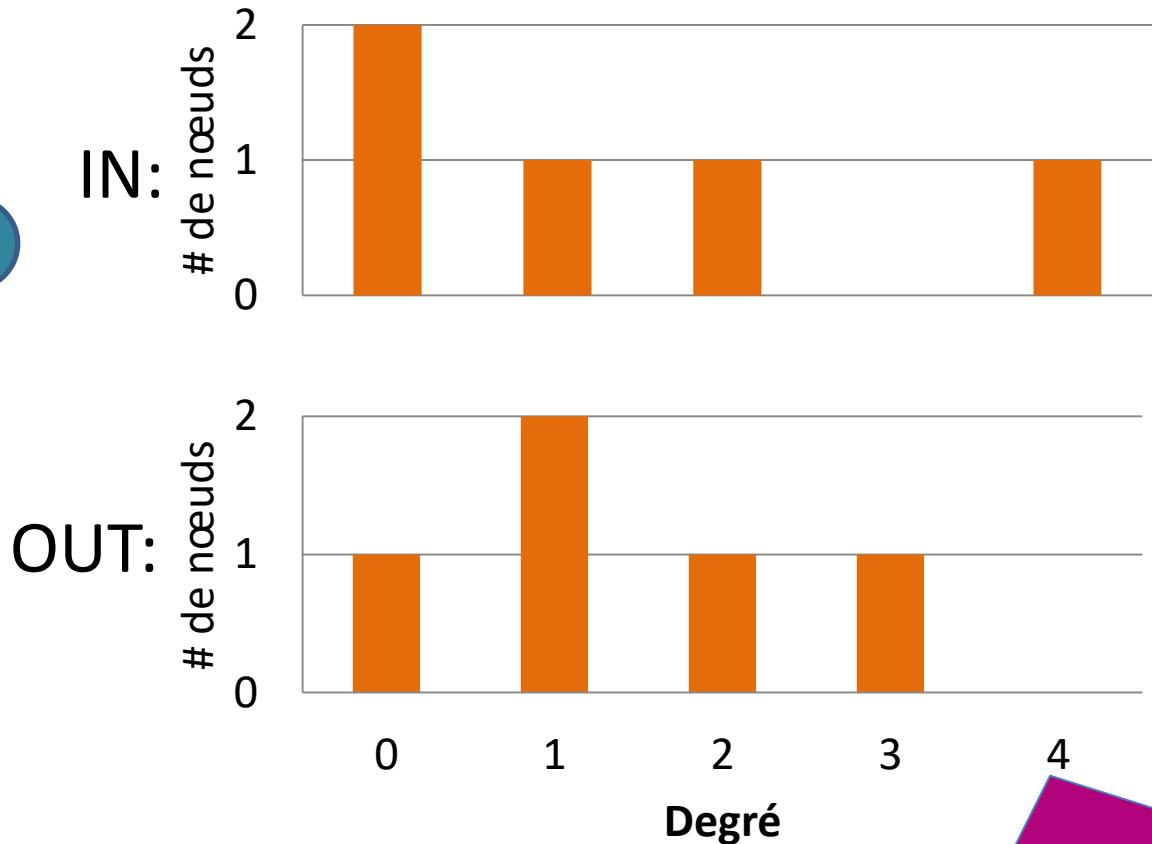
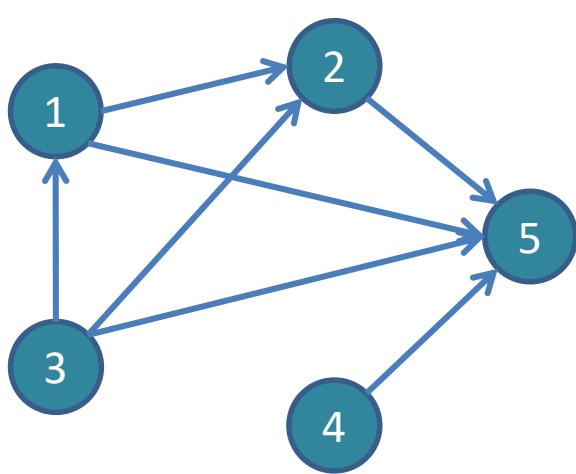
Diamètre



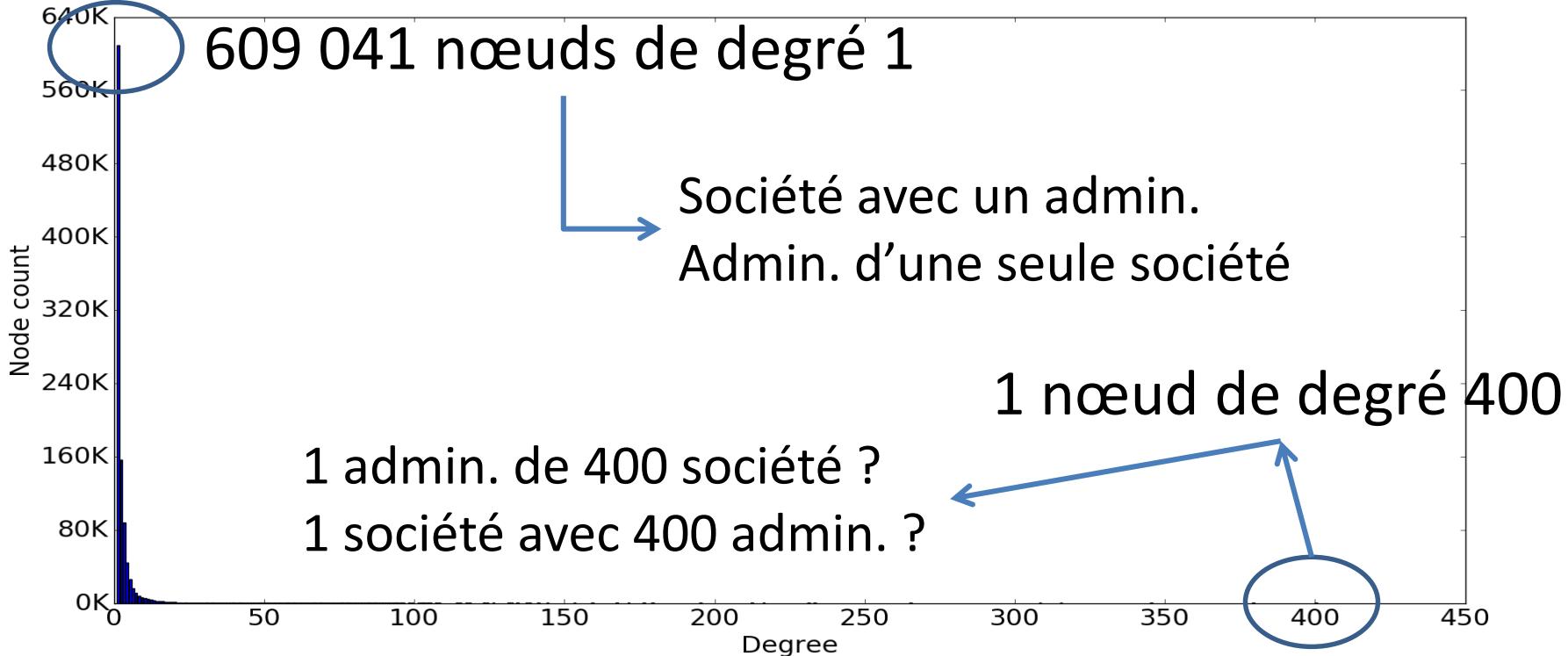
Distribution de degrés



Distribution de degrés (IN vs OUT)

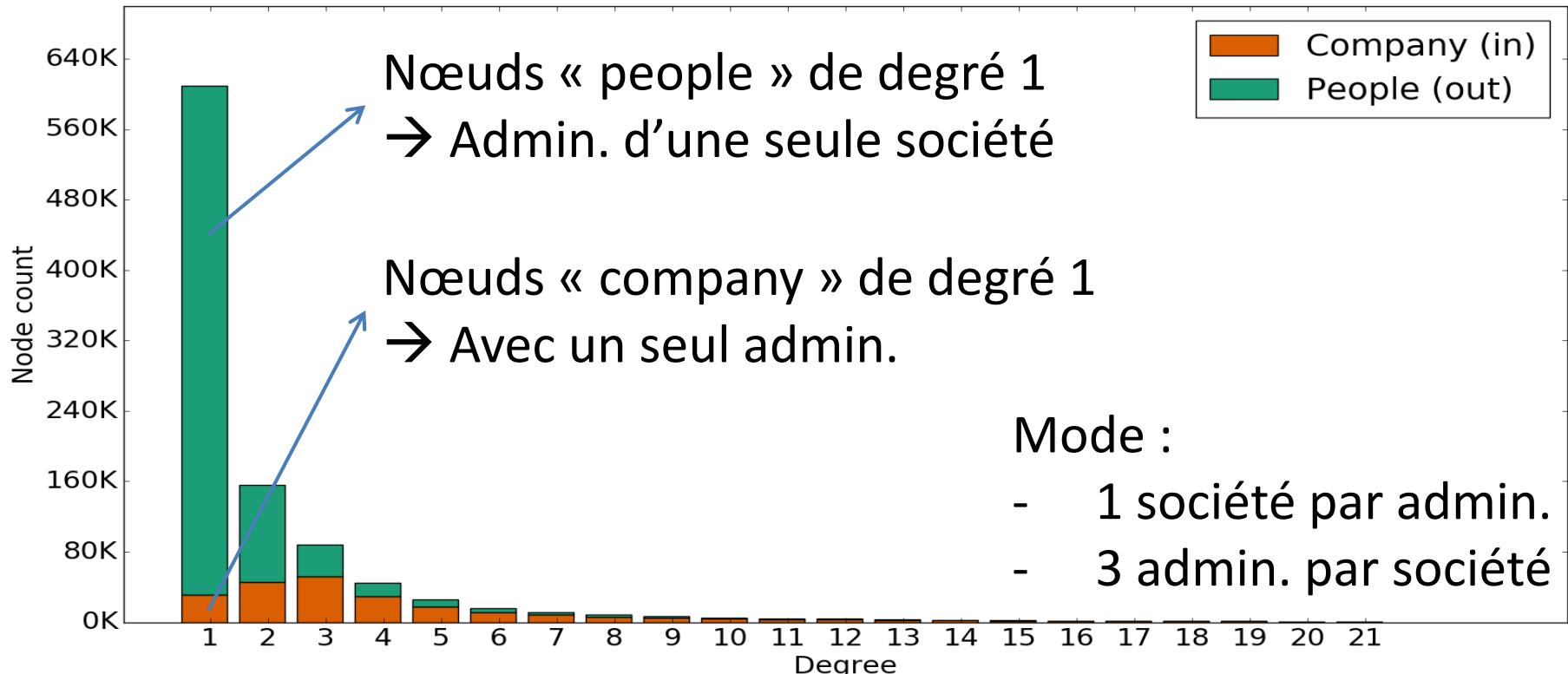


Distribution de degrés : BCE (Adm)



Sur une période de 10 ans

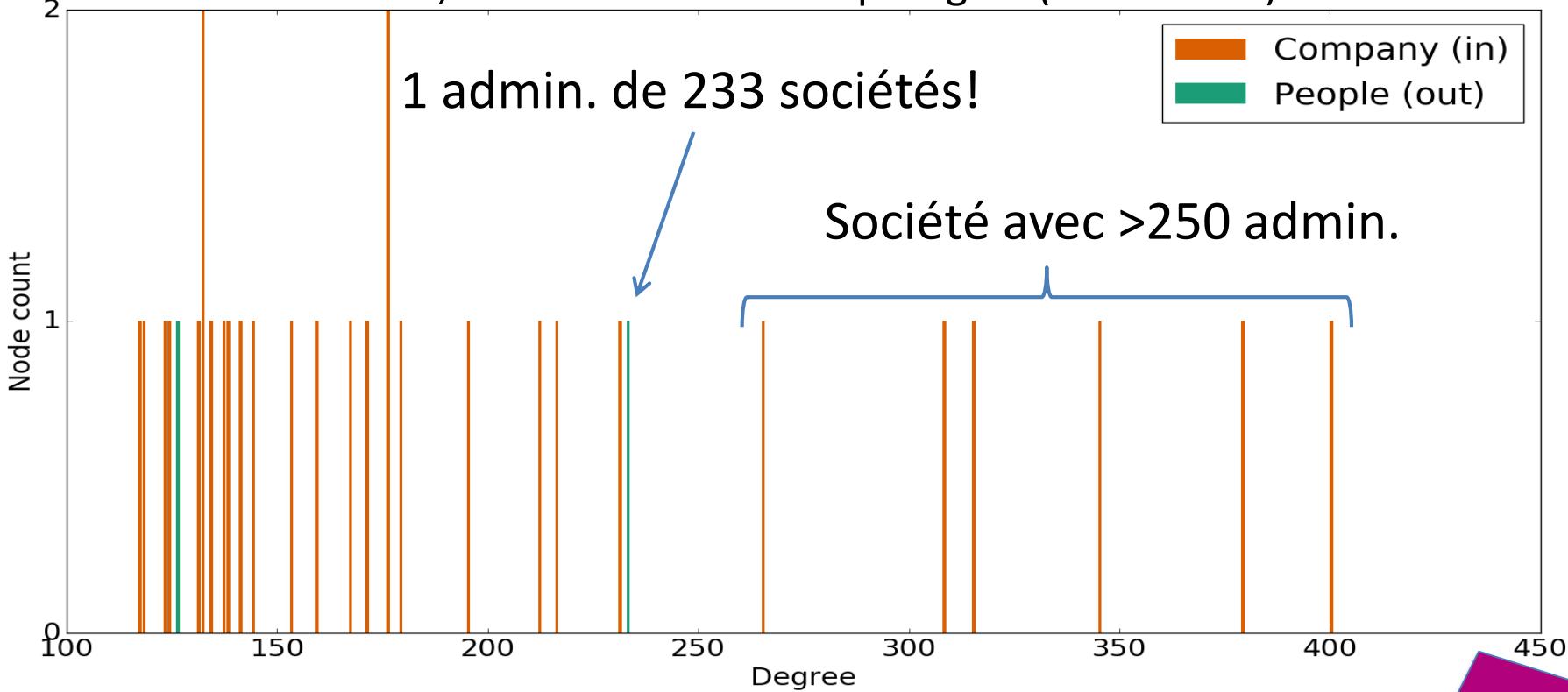
Distribution de degrés : BCE (Adm)



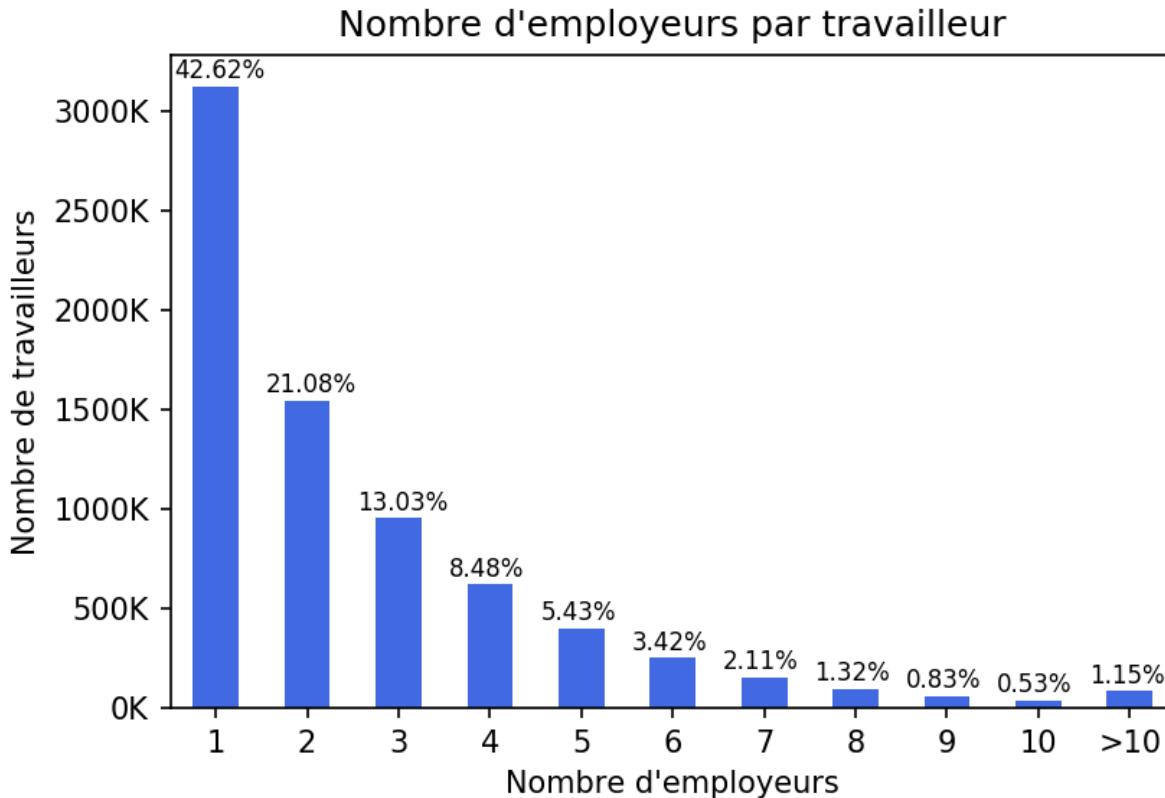
Suppression de 0,5% des nœuds les plus gros

Distribution de degrés : BCE (Adm)

0,003 % des nœuds les plus gros (« outliers »)



Distribution de degrés : Dimona



- 2003-2017
- Max: 110 employeurs !
- 63 : > 50 employeurs

Métriques générales

- Diamètre : étalement du réseau
- Distribution de degré : « forme » du réseau
- Hors présentation :
 - Densité : intensité des relations

Caractériser un réseau

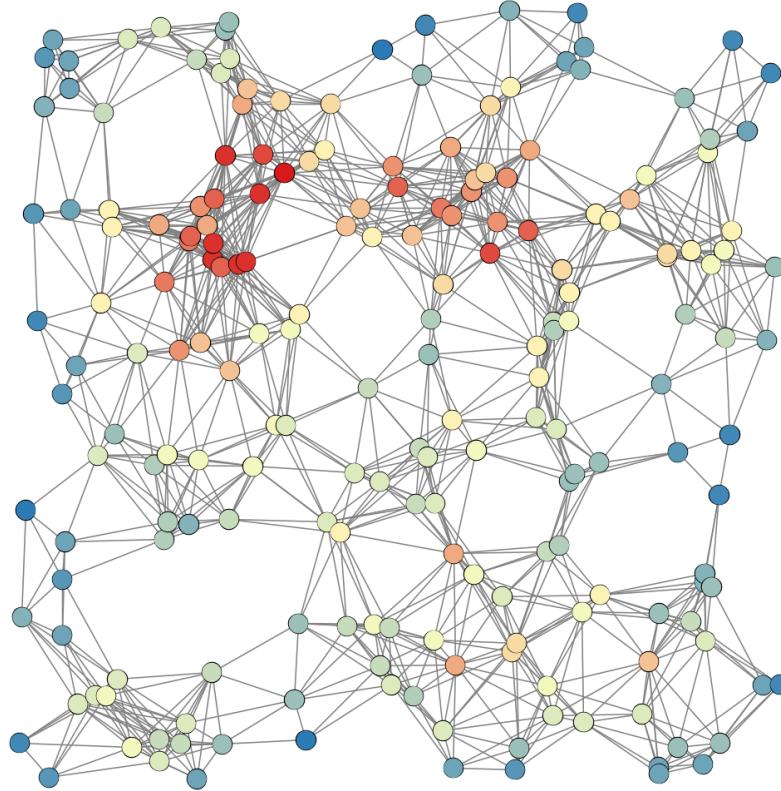
CENTRALITÉ

Centralité

- Centralité d'un nœud = importance d'un nœud au sein du réseau
- Peut se baser sur des caractéristiques :
 - Du nœud isolé (Degree centrality)
 - Du nœud et de son voisinage immédiat (K-core)
 - De l'ensemble du réseau (Betweenness, Page-rank...)

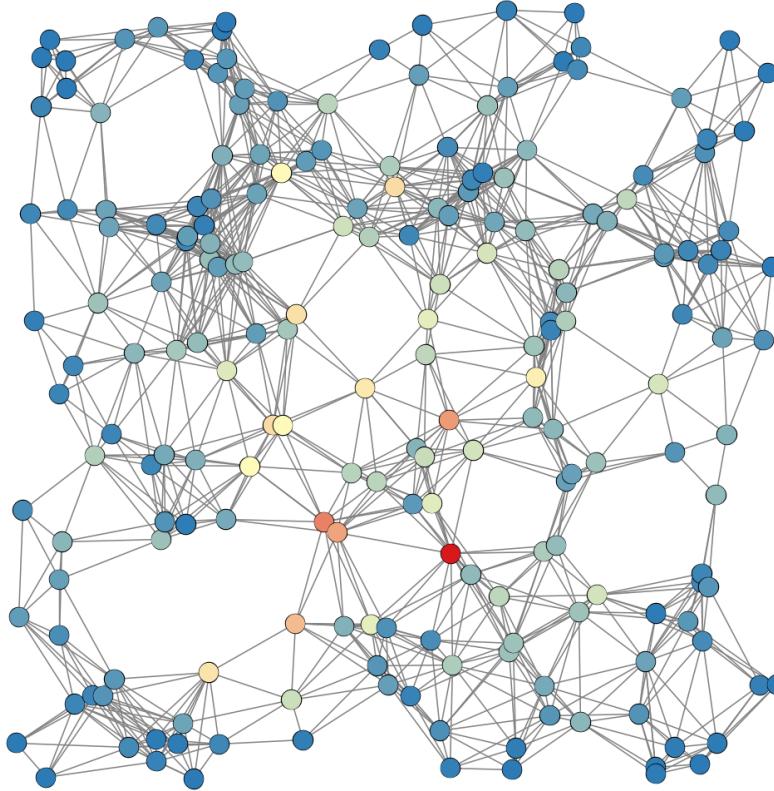
Centralité de degré

Centralité :
degré du
nœud



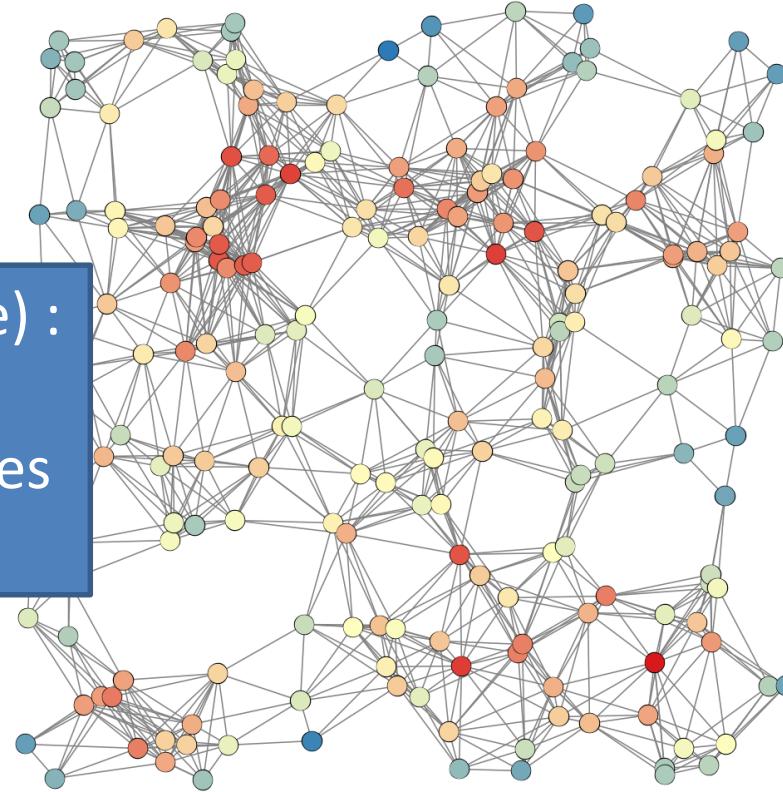
Centralité d'intermédiarité (betweenness)

Centralité :
Nombre de
« plus courts
chemins »
passant par le
nœud



Page-rank

Page-rank (Google) :
Diffusion de son
importance vers ses
voisins



Centralité

- Degré : nombre de voisins
- Betweenness : importance en tant que relai, connecteur, intermédiaire
- Page-rank : importance par diffusion
- Quoi utiliser ? Dépend très fort du business case
- Hors présentation :
 - K-coreness : nombre de voisins de même k-coreness
 - Closeness : centralité « topologique »

Caractériser un réseau

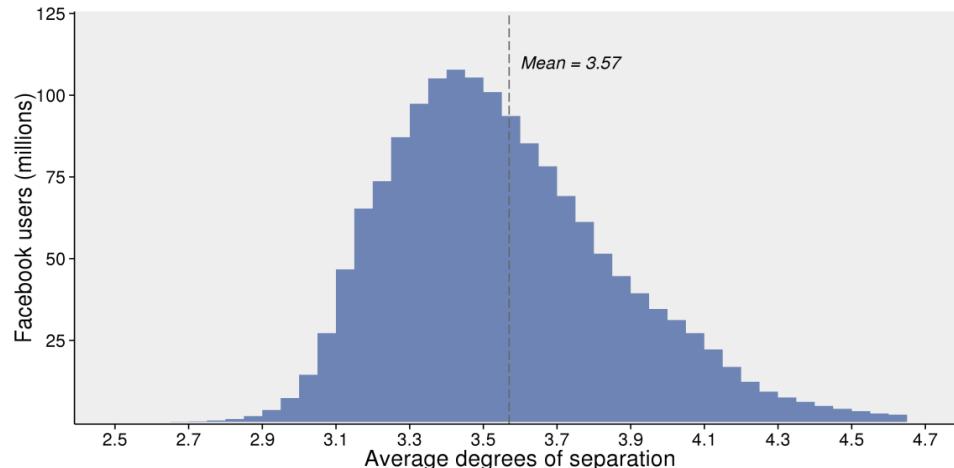
DISTANCE/SIMILARITÉ

Plus court chemin/Distance

- Dijkstra : Calcul du « plus court chemin » entre deux nœuds
- Chaque lien possède un « **poids** » (longueur, coût...)
- **Distance** = Σ poids (=nbr de liens si non pondéré)
- Variante « **A*** » : pour les réseaux **géographiques** (utilisé par les GPS)
- Peut ne pas être unique. Le nombre de « plus courts chemins » est aussi une métrique de proximité

Six degree of separation

- Théorie : distance entre chaque personne sur terre ≤ 6 (via 5 connaissances intermédiaires)
- « Degré de séparation » = plus petite distance
- Facebook : moyenne des moyennes : 4,57



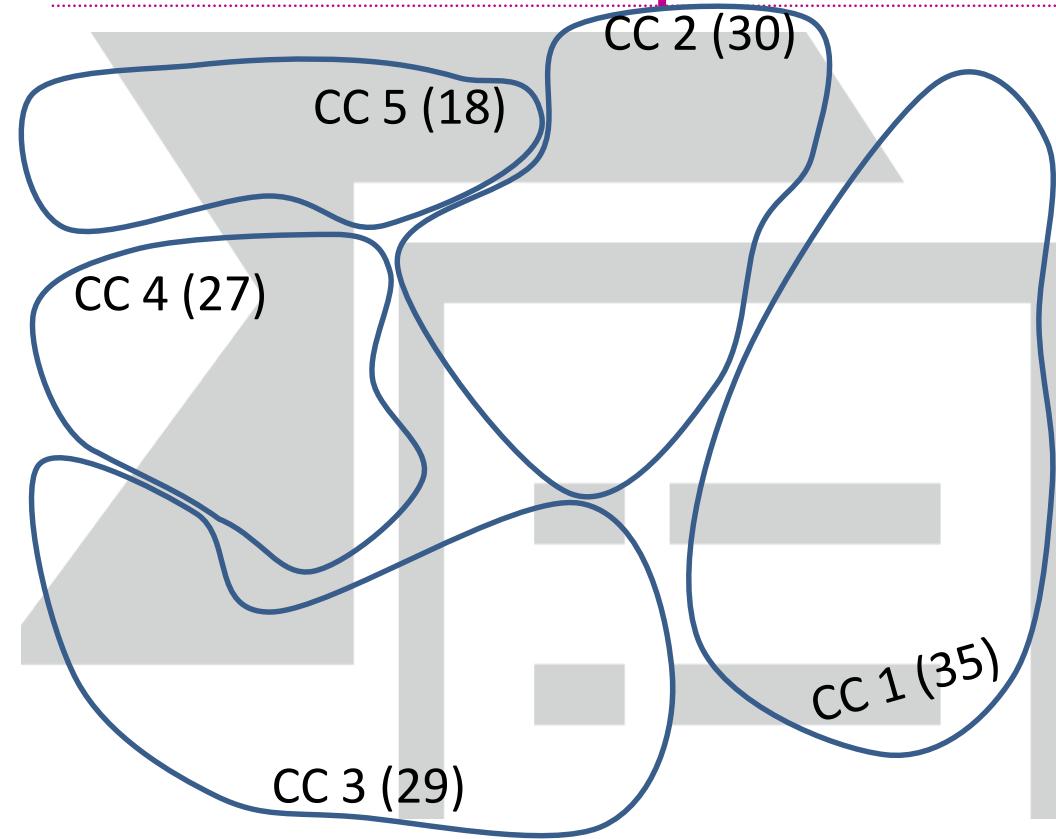
Distance

- Plus court chemin
- Hors présentation :
 - Similarité de Jaccard
 - Connectivité
 - Nombre de plus courts chemins
 - Dépendance sociale (social neighbors)

Caractériser un réseau

CLUSTERING

Composantes connexes



Composante connexe : ensemble (maximal) de nœuds avec un chemin entre chaque paire de nœuds

Composante géante (*giant component*) : la plus grande CC

Beaucoup d'analyses peuvent (doivent) se faire isolément sur chaque composante

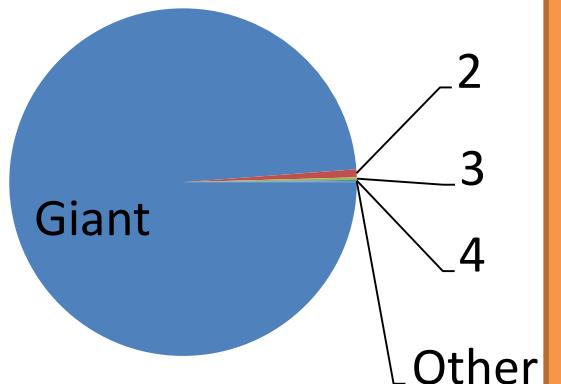
Composante connexe

Dimona (all)

	Taille de la CC	Nbre de CC
1	2	25 929
2	3	5 638
3	4	1 584
4	5	614
...
35	48	1
36	55	1
37	96	1
38	6 965 132	1

25 929 situations avec :

- Un employé sans aucun autre* employeur
- Dans une entreprise sans autre* employé



Giant component :
6,53 millions d'employés
(99,3%) sont « collègues
de (ex-)collègues de ... »
via 430K employeurs
(92%)

* Dans la période considérée

Composante géante

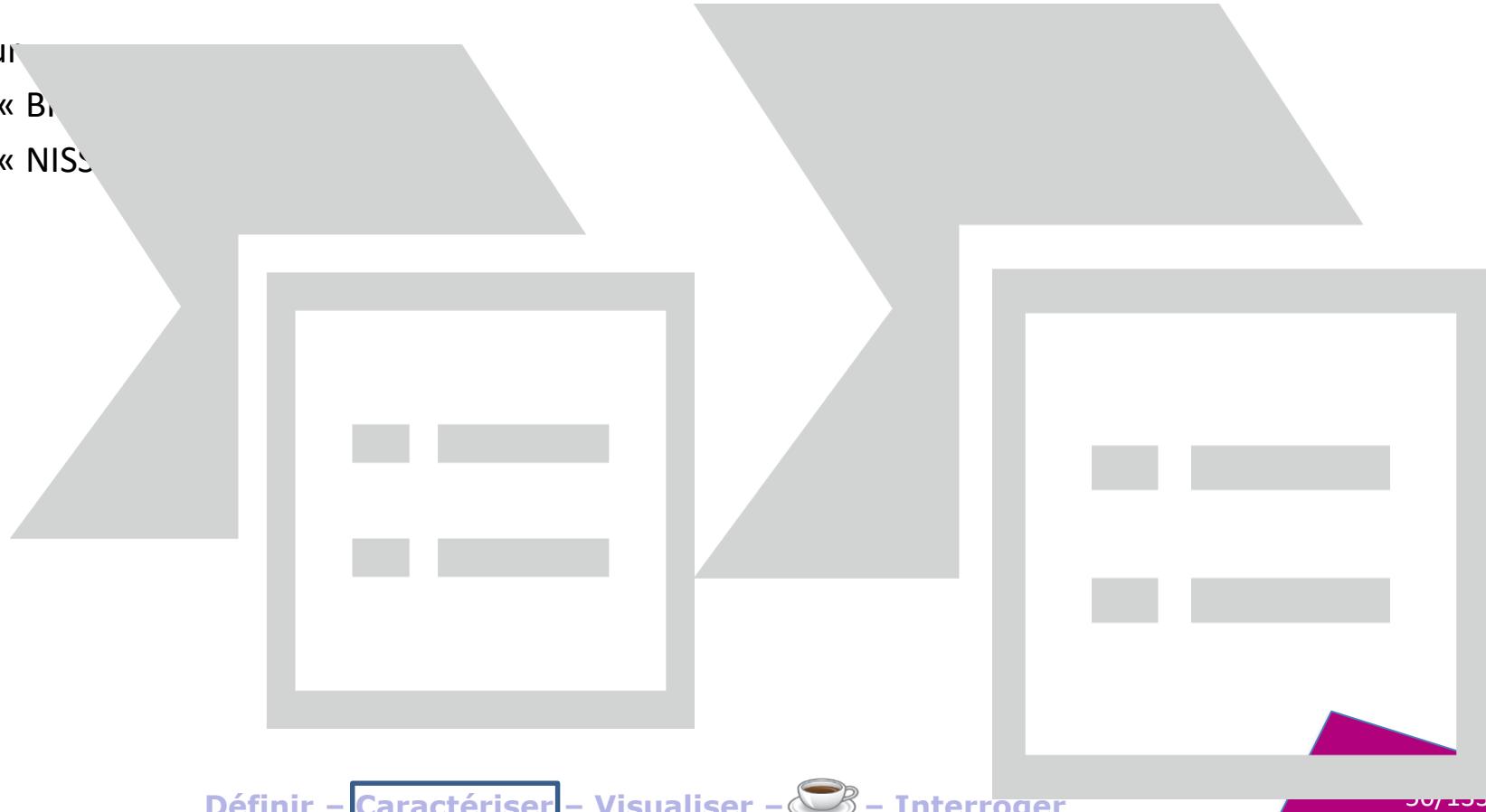
- Diamètre : 20
- Médiane : distance = 6



- → 50 % des « couples » sont « collègue de collègue de collègue »
- Percentile 95 : distance = 8 (+1 « de collègue »)
- Il y a des « super-connecteurs » : tous les enseignants sont collègues !

Composante connexe : sub-giants

- Employeur
- Employé « Br
- Employé « NISS



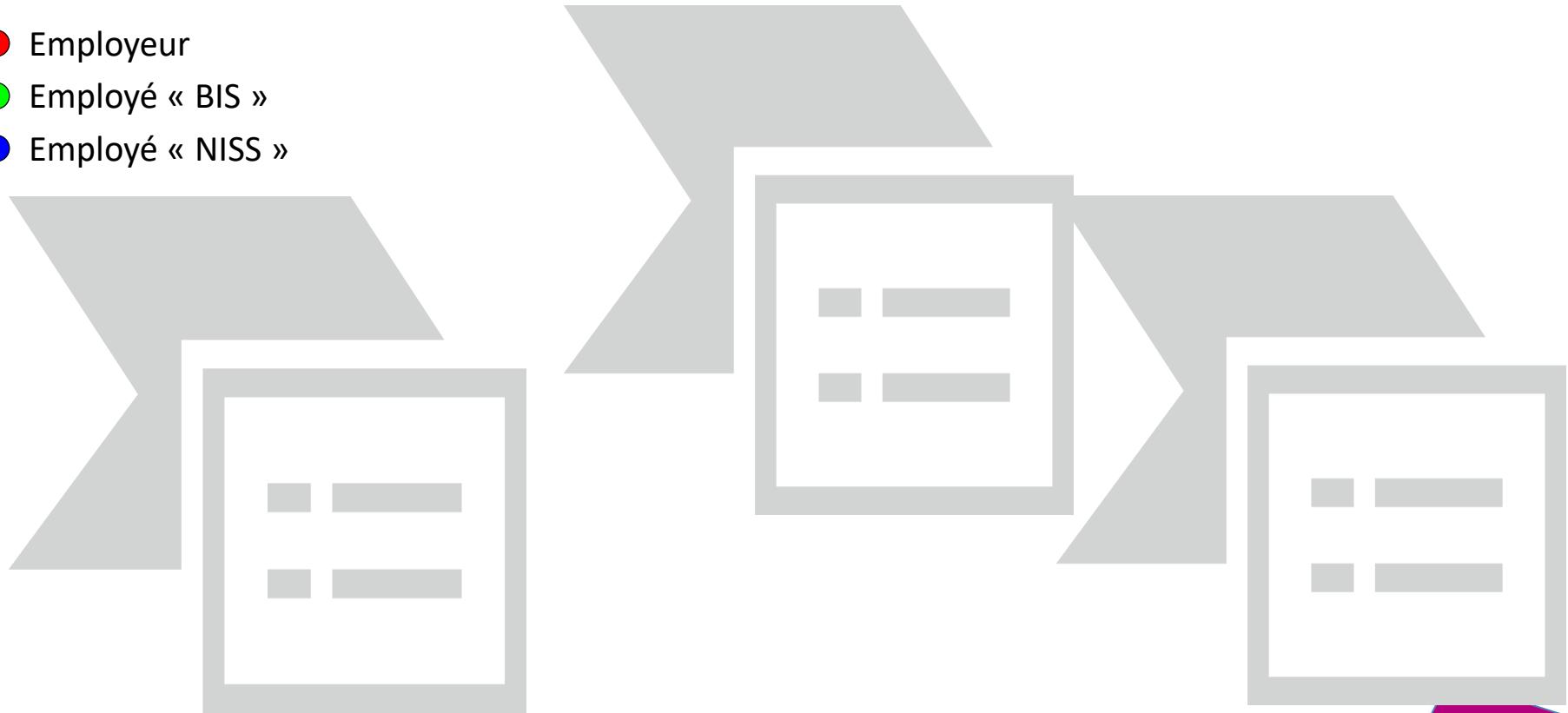
Composante connexe : sub-giants



Définir – Caractériser – Visualiser – – Interroger

Composante connexe : sub-giants

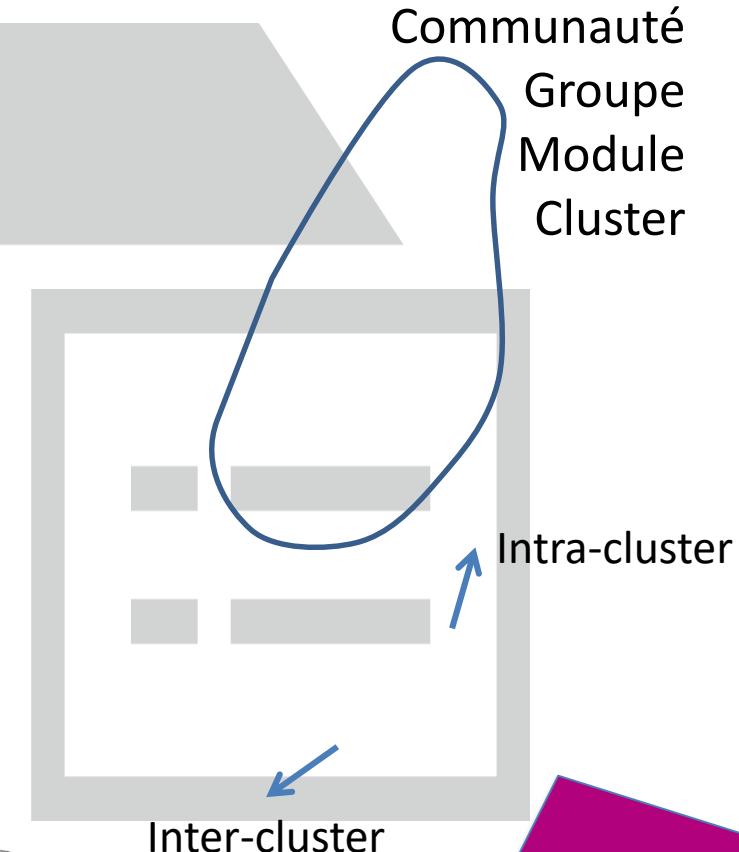
- Employeur
- Employé « BIS »
- Employé « NISS »



Communautés

Coefficient de modularité :

- Mesure de la qualité d'une partition
- $[-1, 1]$, liens inter-cluster vs liens intra-cluster
- Bonne partition (coef. ≈ 1):
liens intra-cluster $>>$ liens inter-cluster
- Optimisation = NP-hard
- Basé sur les liens, pas sur les attributs

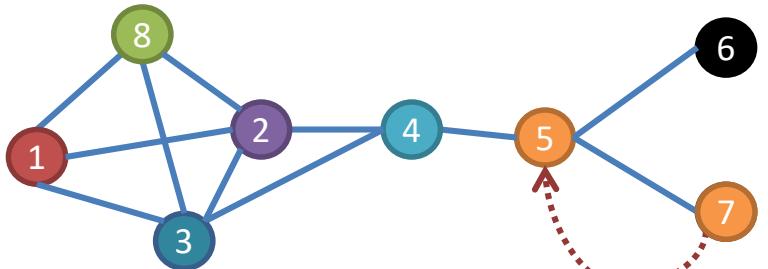


Communautés : intérêts

- Différentes communautés ont souvent des propriétés différentes de l'ensemble
- Permet l'analyse des communautés une par une
- Permet l'identification des « *missing links* » et des « *fake links* »
- Inférer des caractéristiques (attributs)
- Vue d'ensemble avec un « graphe de clusters » (+drill-down)

Algorithme de Louvain (UCL, 2008)

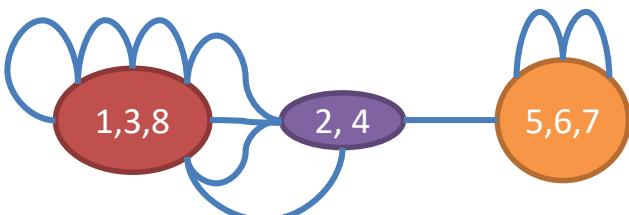
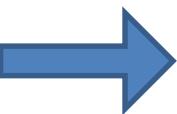
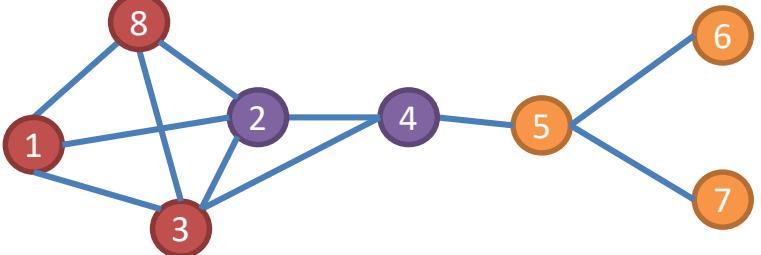
Phase 1



Améliore la « modularity »?
Si oui, migrer



Phase 2



Clustering

- Composantes connexes : groupes « déconnectés »
- Communautés : groupes « sociaux »

Hors présentation :

- Cliques
- Notion de « fortement » ou « faiblement » connecté (réseaux dirigés)

Visualiser un réseau

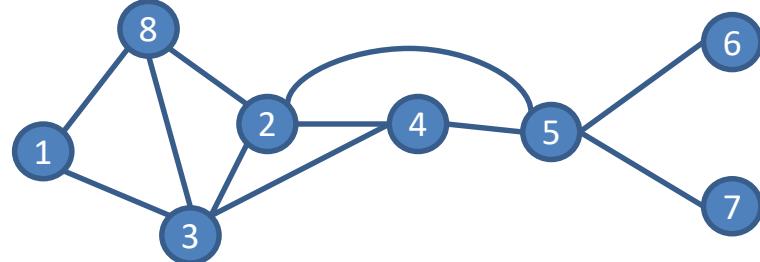
Visualiser

Visualiser un réseau

SPATIALISATION (LAYOUT)

Spatialisation

- Cas simple : graphe **planaire** = qui peut se dessiner sans croisement
- En général :
 - Minimiser les croisements
 - Rapprocher les nœuds connectés
 - Éloigner les nœuds non-connectés



Spatialisation

- Force-directed :
 - Force attractive (spring) entre deux nœuds liés
 - Force répulsive entre tous les nœuds
- Spectral : Utilisation des vecteurs/valeurs propres (*eigenvectors/values*)
- Autres : Hiérarchiques, circulaires, linéaires...

Spatialisation

Yifan Hu

ForceAtlas2

OpenOrd

Fruchterman Reingold

Visualiser un réseau

SIMPLIFICATION

Simplification

- Au-delà de quelques centaines de nœuds, visualiser un graphe est illusoire
- Plusieurs méthodes permettent une approche
 - Découper le graphe
 - Supprimer des nœuds
 - Fusionner des nœuds

Découper le graphe

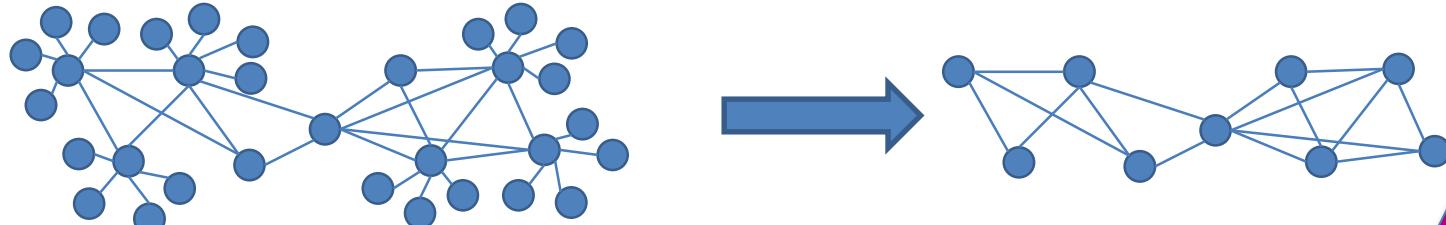
- Composantes connexes
→ Méthode naturelle, pas de perte
- Communautés
→ On perd les relations inter-communauté

Filtrer des noeuds

- Supprimer les **super-connecteurs** (degré élevé)
→ augmente le nombre de composantes connexes



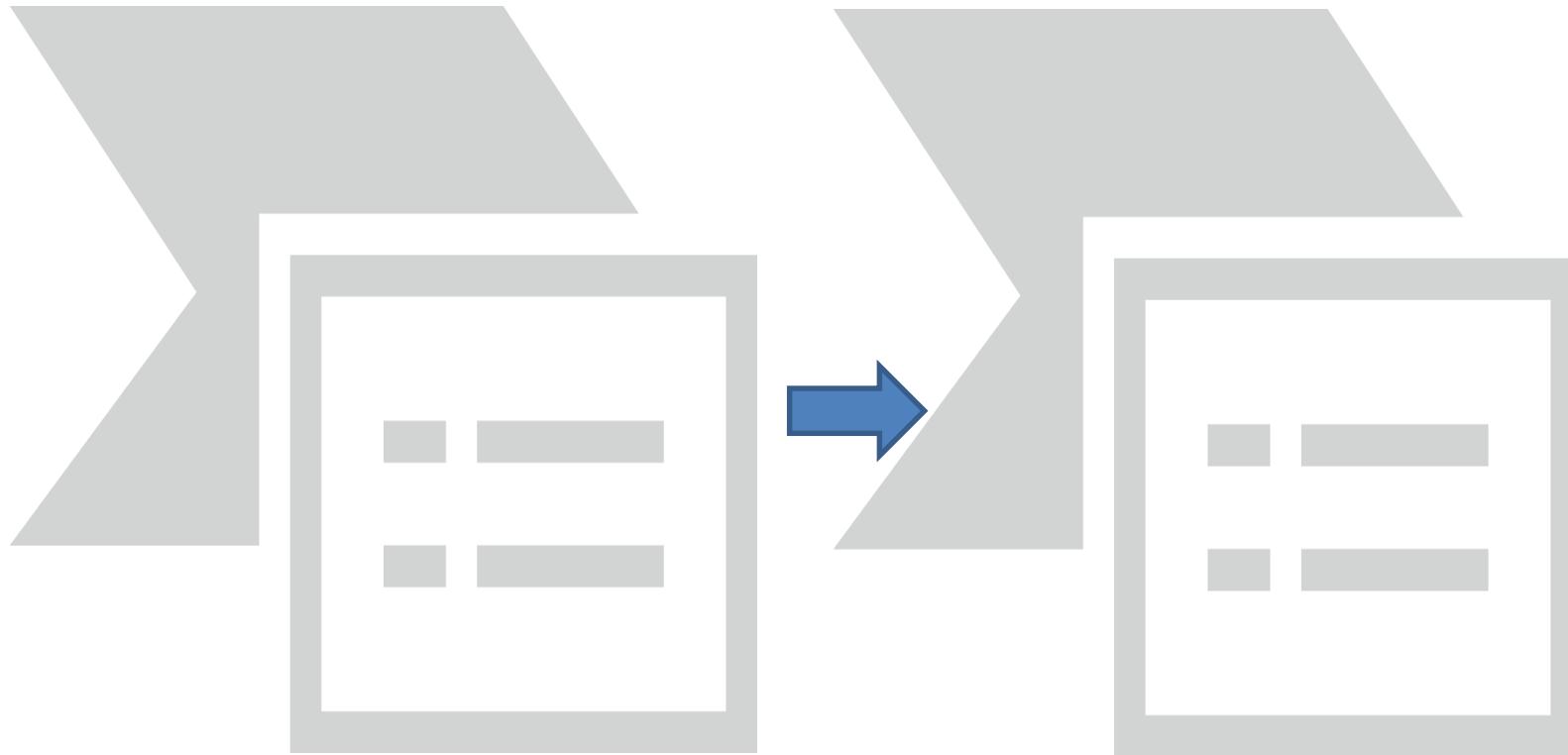
- Supprimer les **feuilles** (degré/K-core faible)
→ focus sur le « backbone/core »



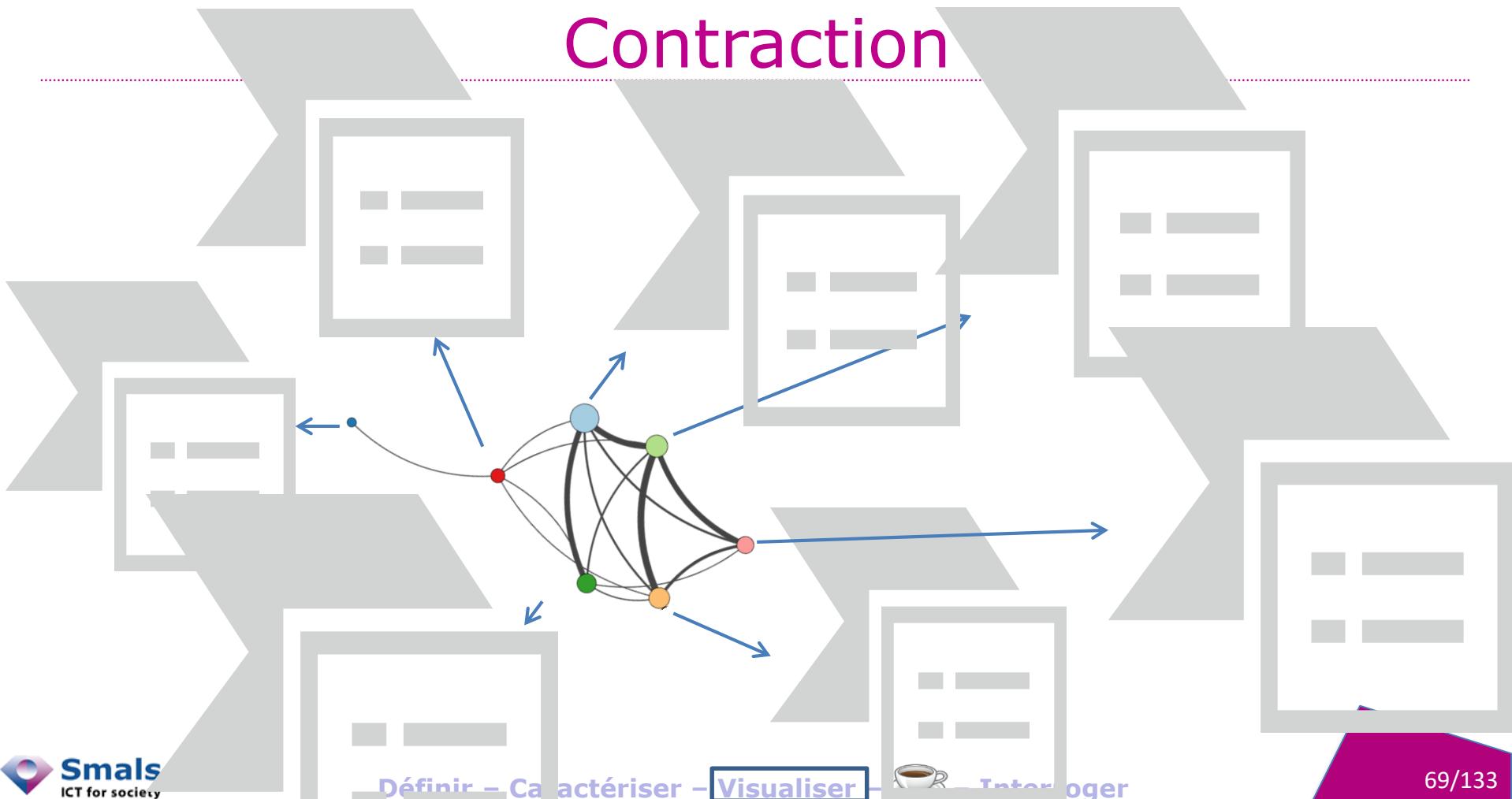
Fusionner des nœuds

- Contraction : fusionner des groupes
- Projection bipartite : conversion d'un graphe biparti en 2 graphes unipartis.

Contraction



Contraction



Contraction : BCE (Administrateurs)

- Administrateurs BCE : Composante géante = 655.158 nœuds (sur $\sim 10^6$)
- On applique Louvain : 743 communautés ($\sim 10''$ en mono-thread)

Contraction : BCE (Administrateurs)

- ▲ Tree
- Star
- Other

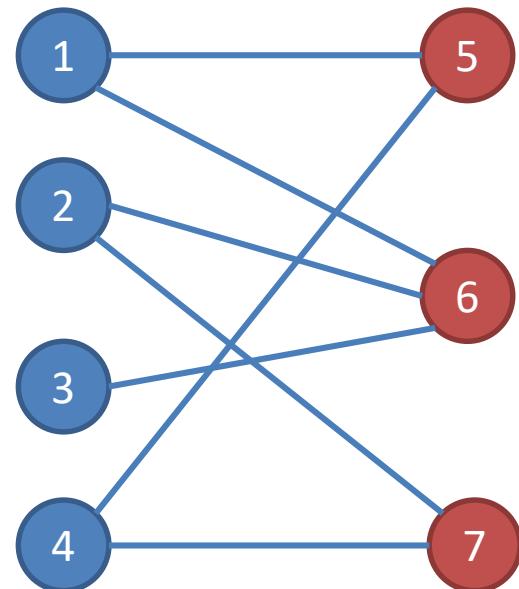
Outgoing edges:

- [10000,inf]
- [1000,9999]
- [100,999]
- [50,99]
- [10,49]
- [5,9]
- 4
- 3
- 2
- 1



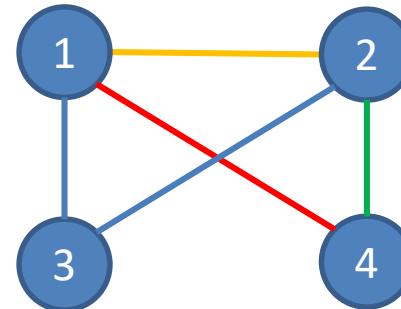
Projection bipartite

Travailleur

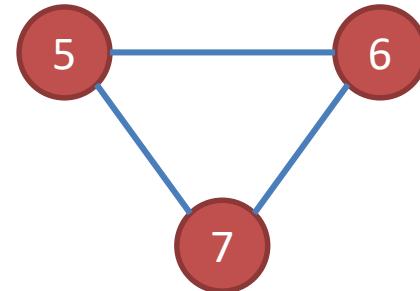


Employeurs

Collègues



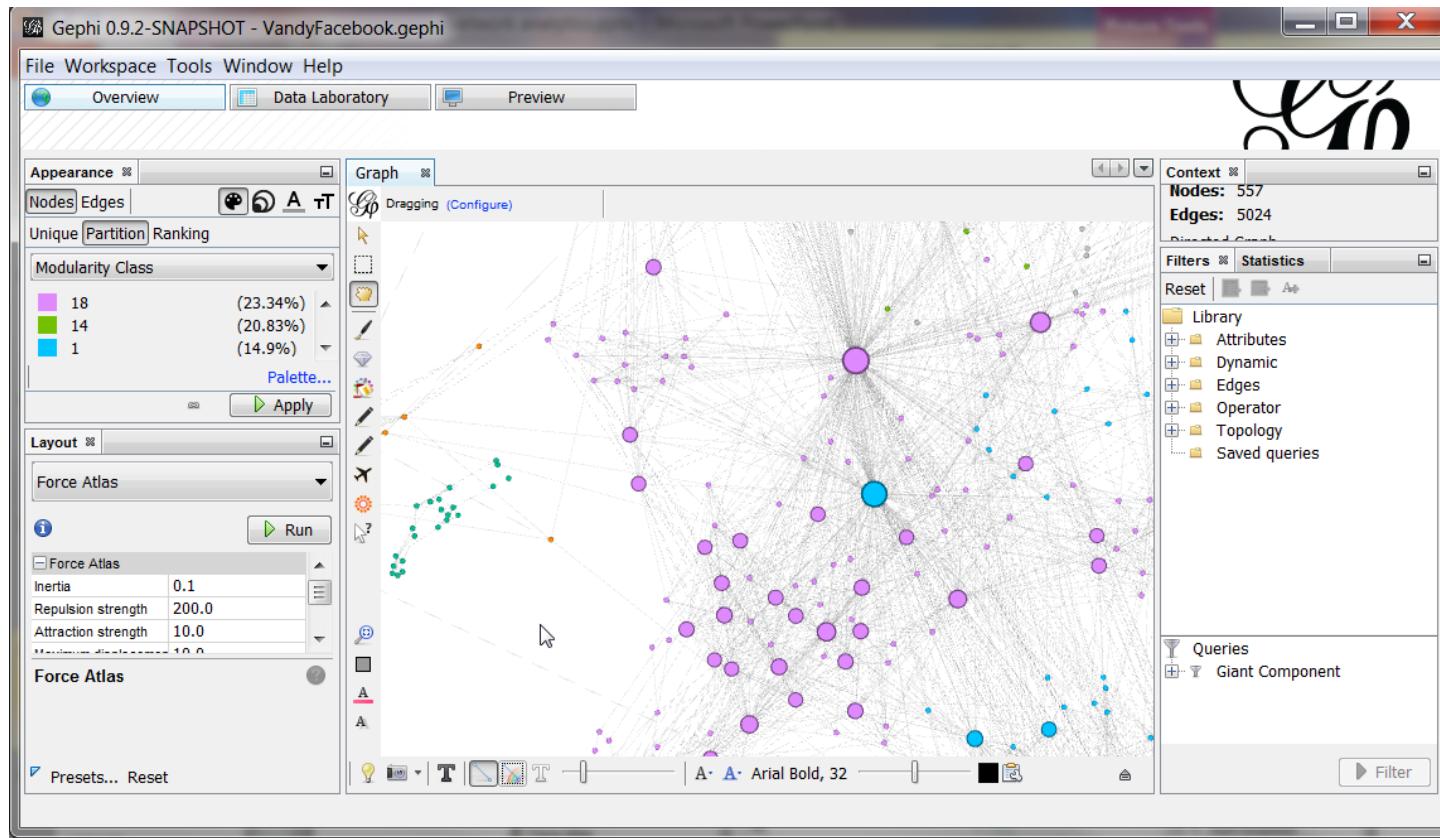
Travailleurs communs



Visualiser un réseau

OUTIL DESKTOP : GEPHI

Gephi



Visualiser un réseau

LIBRAIRIE JS: VIS.JS

vis.js

```
<div id="graph"></div>
<script type="text/javascript">
var nodes = new vis.DataSet([
  {id: 1, label: 'Node 1'},
  {id: 2, label: 'Node 2'},
  {id: 3, label: 'Node 3'},
]);
var edges = new vis.DataSet([
  {from: 1, to: 3},
  {from: 1, to: 2},
  {from: 3, to: 3}
]);
```

```
var container =
document.getElementById('graph');

var data = {
  nodes: nodes,
  edges: edges
};

var options = {};

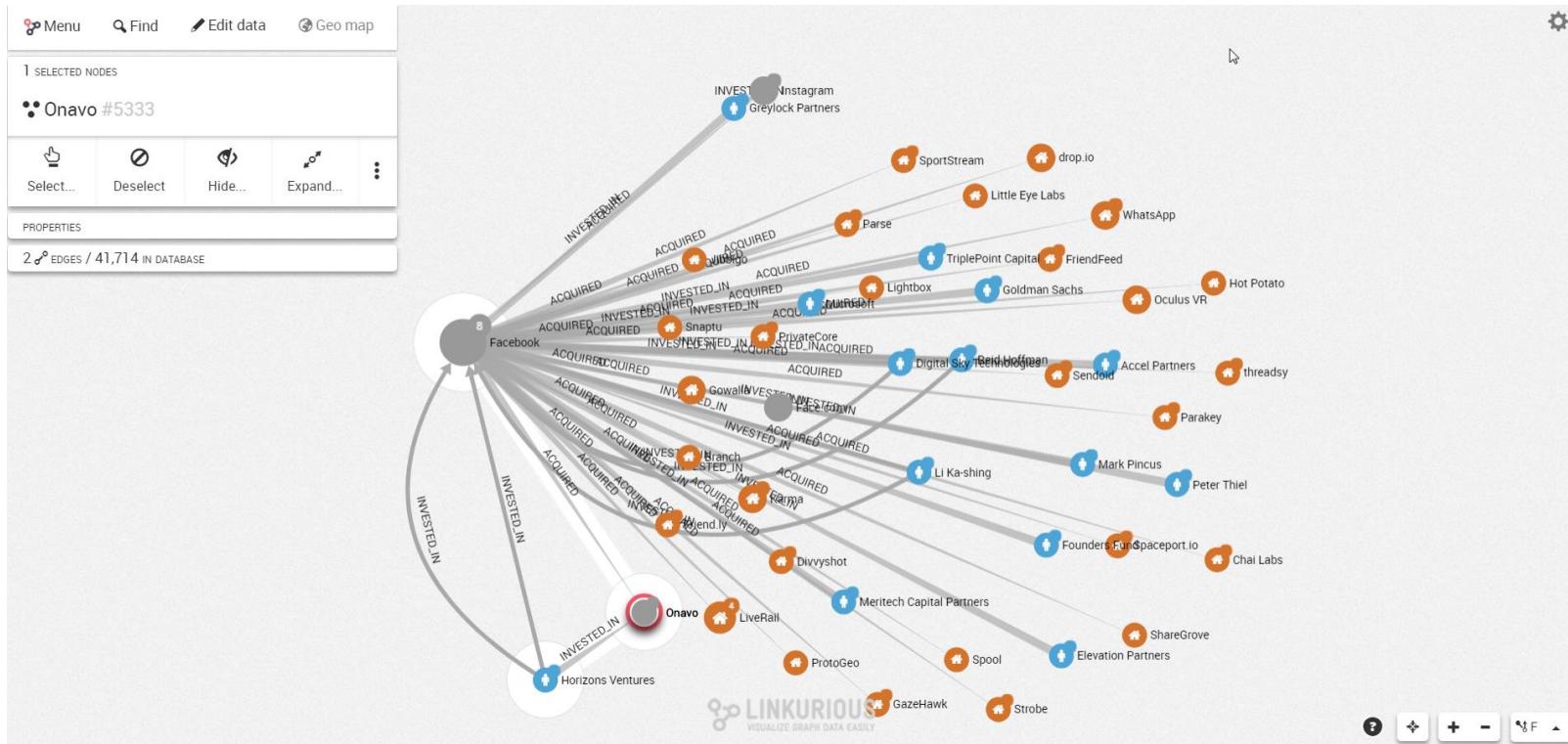
var network =
new vis.Network(container,
  data,
  options);

</script>
```

Visualiser un réseau

OUTIL WEB : LINKURIOUS

Linkurio.us

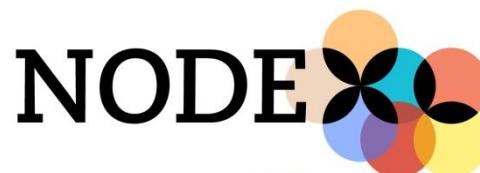


Visualiser un réseau

AUTRES OUTILS

Outils desktop

- Gephi <https://gephi.org> OS
- TouchGraph <http://www.touchgraph.com> €
- Cytoscape <http://www.cytoscape.org> OS
- Maltego <https://www.paterva.com>
- NodeXL <https://nodexl.codeplex.com> OS
- SAS Network Analytics €
- IBM i2 €



Outils Cloud

- Linkurio.us <https://linkurio.us> €
- Keylines <https://cambridge-intelligence.com/keylines> €



Librairies web

- Vis.js <http://visjs.org> OS
- Cytoscape.js <http://js.cytoscape.org> OS
- Sigma.js <http://sigmajs.org> OS
- D3.js <https://d3js.org> OS
- OGMA (Linkurio.us SDK) €



sigmajs



Questions ?



Pause !

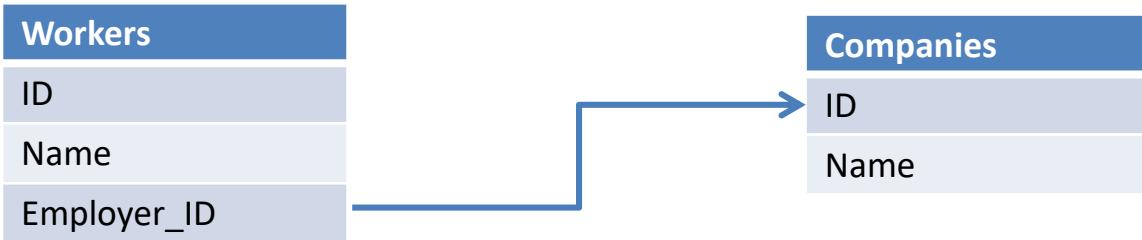




Interroger un réseau

GRAPH DATABASES

RDBMS et relations



Employés de Smals ?

Ce qui intéresse
le développeur

Comment les relations
sont implémentées

```
SELECT Workers.Name
FROM Workers
JOIN Companies
ON Workers.Employer_ID = Companies.ID
WHERE Companies.Name = 'Smals'
```

RDBMS et relations



Employés de Smals ?

Ce qui intéresse
le développeur

Comment les relations
sont implémentées

```
SELECT Workers.Name
FROM Workers
JOIN Works_for
ON Workers.ID = Works_for.Worker_ID
JOIN Companies
ON Works_for.Company_ID = Companies.ID
WHERE Companies.Name = 'Smals'
```

RDBMS et relations

ID	Name
1	Alice
2	Bob
3	Camille
4	Zoé

Liker_ID	Liked_ID
1	3
2	1
2	3
3	2
3	4

```
SELECT p1.Name
```

```
FROM People p1
```

```
JOIN Likes
```

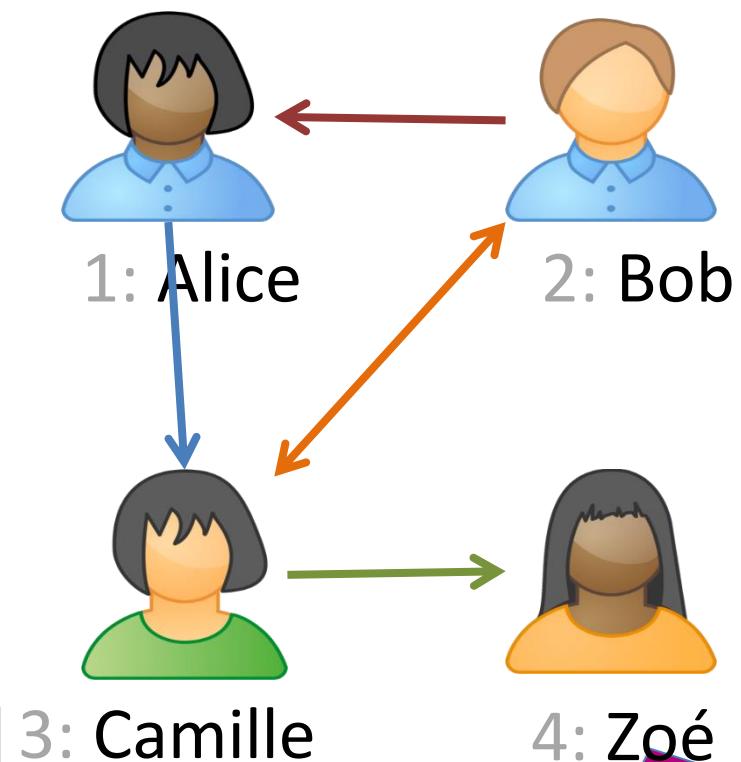
```
ON Likes.Liked_ID = p1.ID
```

```
JOIN People p2
```

```
ON Likes.Liker_ID = p2.ID
```

```
WHERE p2.Name = "Bob"
```

Qui Bob aime ?



RDBMS et relations

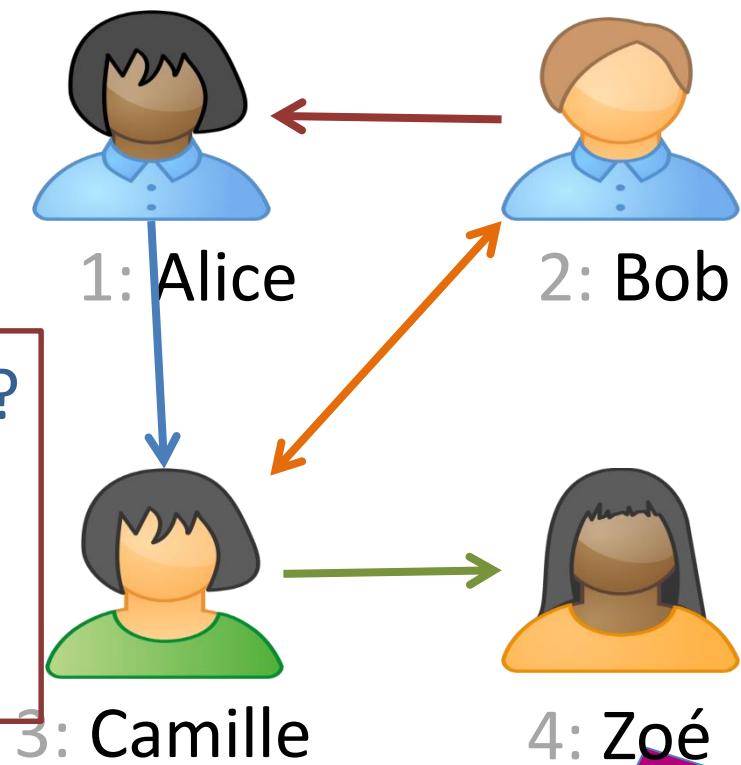
ID	Name
1	Alice
2	Bob
3	Camille
4	Zoé

Liker_ID	Liked_ID
1	3
2	1
2	3
3	2
3	4

```
SELECT p1.Name
```

```
FROM People p1  
JOIN Likes l1  
ON l1.Liked_ID = p1.ID  
JOIN Likes l2  
ON l1.Liker_ID = l2.Liked_ID  
JOIN People p2  
ON l2.Liker_ID = p2.ID  
WHERE p2.Name = "Bob"
```

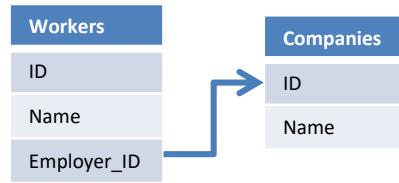
Qui Bob (aime)² ?



RDBMS et relations : limitations

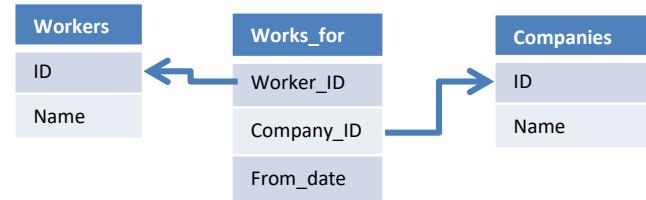
Relations 1-n

Foreign key:



Relations m-n (ou 1-n avec attributs)

Join table:

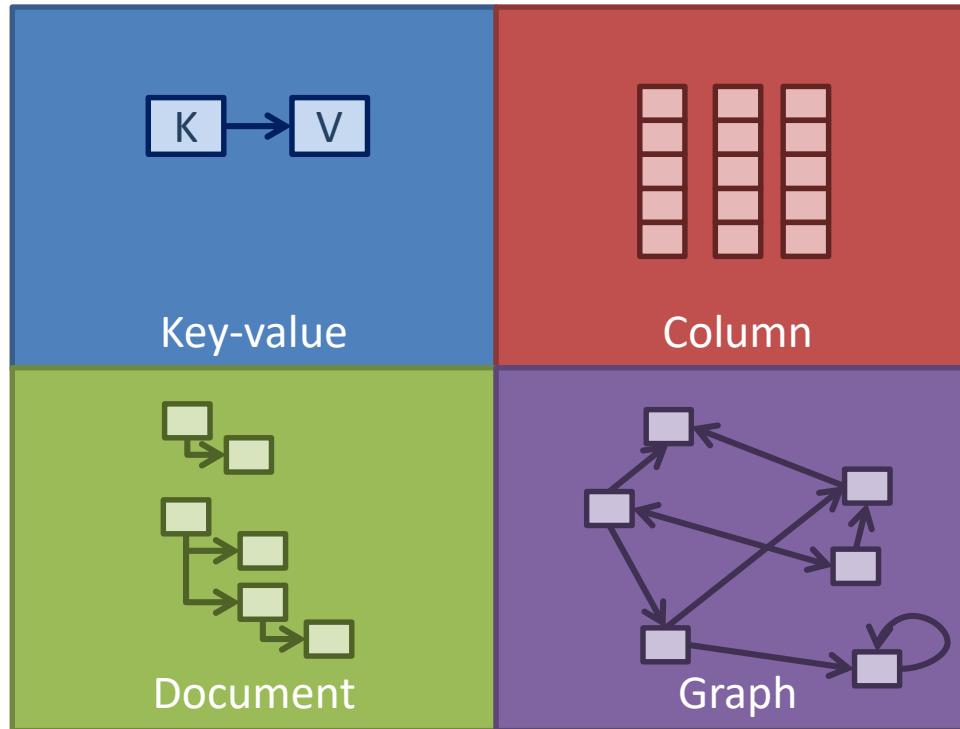


- Détourne le rôle d'un **attribut**
- Pas de type « Key »
- Contrainte d'intégrité possible (externe à la table)
- Structure de la relation à préciser à chaque requête
- Join : Lourd en écriture et en exécution
- Détourne le rôle d'une **table**
- Pas distinguable des autres

Solutions NoSQL



ORACLE
NOSQL DATABASE



Google
BigTable

APACHE
HBASE

AllegroGraph

neo4j

InfiniteGraph
Powered by Objectivity

GraphDB : objectifs

Objectifs des *bases de données orientées graphes* :

Querying language qui sépare la définition **des relations** (à la création) et la recherche de « **matching** »

```
SELECT p1.Name  
FROM People p1  
JOIN Likes l1  
    ON l1.Liked_ID = p1.ID  
JOIN Likes l2  
    ON l1.Liker_ID = l2.Liked_ID  
JOIN People p2  
    ON l2.Liker_ID = p2.ID  
WHERE p2.Name = "Bob"
```

Moteur efficace pour le parcours des **relations**

Cypher (Neo4j)

```
MATCH (:People {Name:"Bob"})  
- [:Likes*2] ->  
(p:People)  
RETURN p.Name
```

GraphDB: moins de code

```
(SELECT T.directReportees AS directReportees, sum(T.count) AS count
FROM (
SELECT manager.pid AS directReportees, 0 AS count
FROM person_reportee_manager
WHERE manager.pid = (SELECT id FROM person WHERE name = "fName lName")
UNION
SELECT manager.pid AS directReportees, count(manager.directly_manages) AS count
FROM person_reportee_manager
WHERE manager.pid = (SELECT id FROM person WHERE name = "fName lName")
GROUP BY directReportees
UNION
SELECT manager.pid AS directReportees, count(reportee.directly_manages) AS count
FROM person_reportee_manager
JOIN person_reportee reportee
ON manager.directly_manages = reportee.pid
WHERE manager.pid = (SELECT id FROM person WHERE name = "fName lName")
GROUP BY directReportees
UNION
SELECT manager.pid AS directReportees, count(L2Reportees.directly_manages) AS count
FROM person_reportee_manager
JOIN person_reportee L1Reportees
ON manager.directly_manages = L1Reportees.pid
JOIN person_reportee L2Reportees
ON L1Reportees.directly_manages = L2Reportees.pid
WHERE manager.pid = (SELECT id FROM person WHERE name = "fName lName")
GROUP BY directReportees
) AS T
GROUP BY directReportees)
UNION
(SELECT T.directReportees AS directReportees, sum(T.count) AS count
FROM (
SELECT manager.directly_manages AS directReportees, 0 AS count
FROM person_reportee_manager
WHERE manager.pid = (SELECT id FROM person WHERE name = "fName lName")
UNION
SELECT reportee.pid AS directReportees, count(reportee.directly_manages) AS count
FROM person_reportee_manager
JOIN person_reportee reportee
ON manager.directly_manages = reportee.pid
WHERE manager.pid = (SELECT id FROM person WHERE name = "fName lName")
GROUP BY directReportees
) AS T
GROUP BY directReportees)
```

```
SELECT depth1Reportees.pid AS directReportees,
count(depth2Reportees.directly_manages) AS count
FROM person_reportee_manager
JOIN person_reportee L1Reportees
ON manager.directly_manages = L1Reportees.pid
JOIN person_reportee L2Reportees
ON L1Reportees.directly_manages = L2Reportees.pid
WHERE manager.pid = (SELECT id FROM person WHERE name = "fName lName")
GROUP BY directReportees
) AS T
GROUP BY directReportees)
UNION
(SELECT T.directReportees AS directReportees, sum(T.count) AS count
FROM(
SELECT reportee.directly_manages AS directReportees, 0 AS count
FROM person_reportee_manager
JOIN person_reportee reportee
ON manager.directly_manages = reportee.pid
WHERE manager.pid = (SELECT id FROM person WHERE name = "fName lName")
GROUP BY directReportees
)
SELECT L2Reportees.pid AS directReportees, count(L2Reportees.directly_manages)
AS count
FROM person_reportee_manager
JOIN person_reportee L1Reportees
ON manager.directly_manages = L1Reportees.pid
JOIN person_reportee L2Reportees
ON L1Reportees.directly_manages = L2Reportees.pid
WHERE manager.pid = (SELECT id FROM person WHERE name = "fName lName")
GROUP BY directReportees
) AS T
GROUP BY directReportees)
```

MATCH

```
(boss) - [:MANAGES*0..3] -> (sub) ,  
(sub) - [:MANAGES*1..3] -> (report)
```

WHERE boss.name = « John Doe »

RETURN sub.name **AS** Subordinate,

count(report) **AS** Total

GraphDB : moins de code

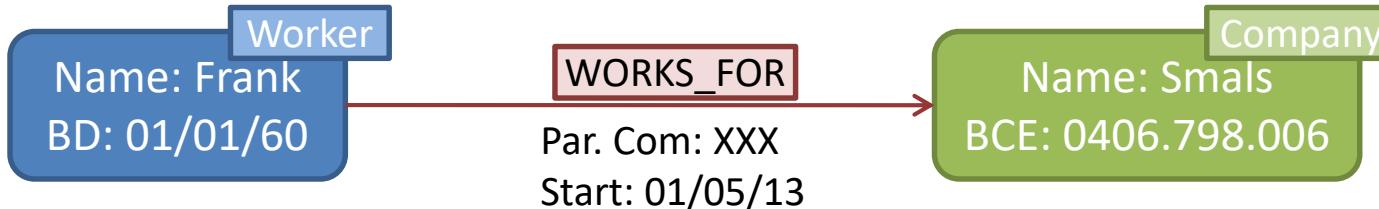
"We found Neo4j to be literally thousands of times faster than our prior MySQL solution, with queries that require 10-100 times less code. Today, Neo4j provides eBay with functionality that was previously impossible."



Interroger un réseau

NEO4J

Requête Neo4j



Nœud de type « Worker »

Nœud de type « Company »

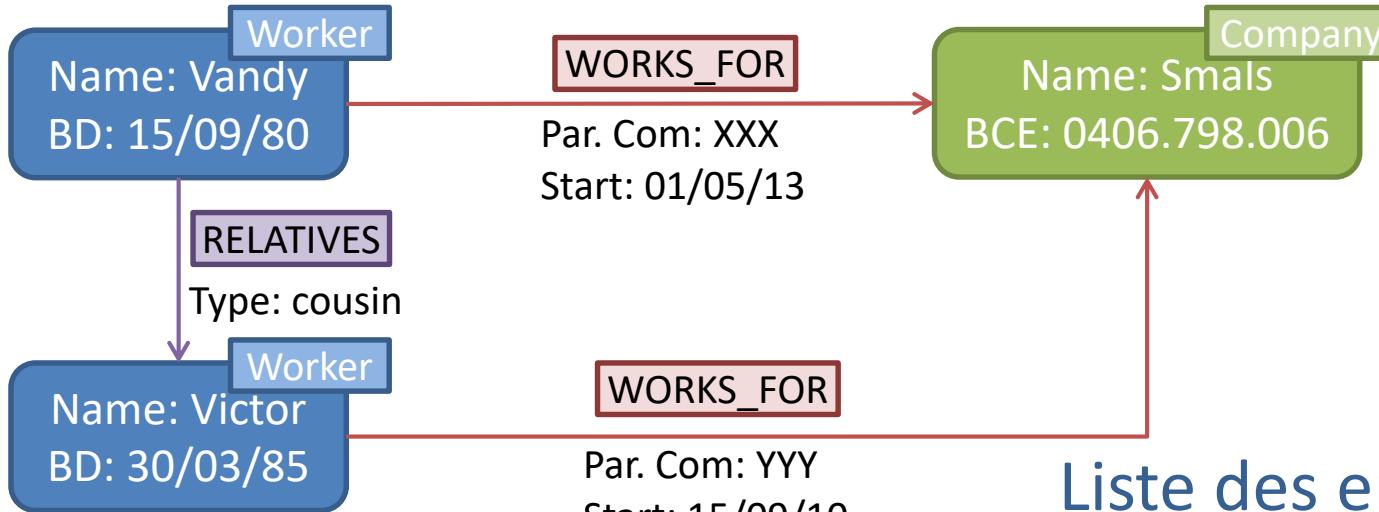
Relation de type « WORKS_FOR »

```
MATCH (w:Worker)  
WHERE c.Name = "Smals"  
RETURN w.Name, w.BD, r.Start
```

Variante :

```
MATCH  
  (w:Worker)  
  - [r:WORKS_FOR] ->  
  (c:Company {Name:"Smals"})  
RETURN ...
```

Requête Neo4j



Liste des employés de
Smals apparentés

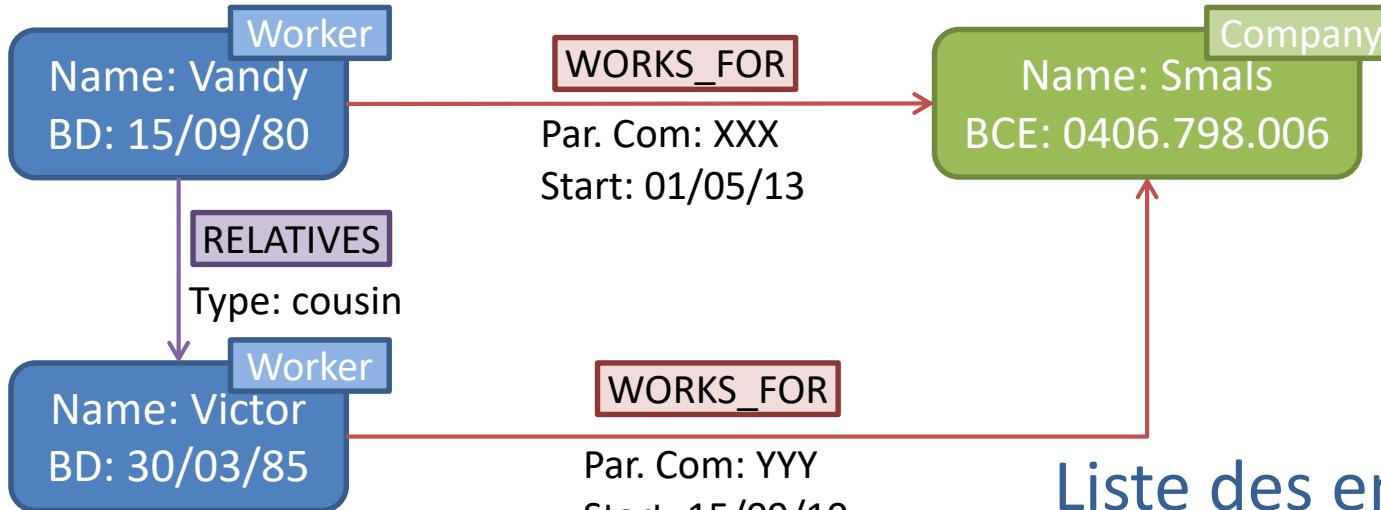
MATCH

```
(w1:Worker) - [:WORKS_FOR] -> (c:Company {Name:"Smals"}),  
(w2:Worker) - [:WORKS_FOR] -> (c),  
(w1) - [:RELATIVES] -> (w2)
```

RETURN

w1.Name, w2.Name

Requête Neo4j



Liste des employés de
Smals apparentés

MATCH

```
(w1:Worker) - [:WORKS_FOR] -> (c:Company {Name:"Smals"})  
<- [:WORKS_FOR] - (w2:Worker) - [:RELATIVES] -> (w1)
```

RETURN w1.Name, w2.Name

Fonctionnement

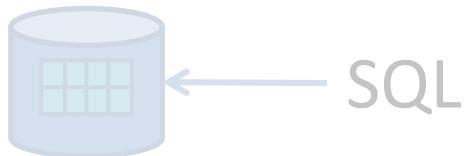
	RDBMS	Graph DB
schema	<pre> graph LR subgraph RDBMS [RDBMS Schema] direction TB W[Workers
ID
Name] --- WF[Works_for
Worker_ID
Company_ID
From_date] WF --- C[Companies
ID
Name] end </pre>	<pre> graph LR subgraph GraphDB [Graph DB Schema] direction TB W1[Worker
Name: Vandy
BD: 15/09/80] -- WORKS_FOR --> C1[Company
Name: Smals
BCE: 0406.798.006] end </pre>
WHERE ...	Search in Companies (with index): $O(\log([\text{Table size}]))$	Search in nodes (with index) : $O(\log([\text{DB size}]))$
Relation	JOIN : <ul style="list-style-type: none"> - Search in Works_for : $O(\log([\text{Table size}]))$ - Search in Workers : $O(\log([\text{Table size}]) \times [\text{nb res}])$ 	Follow pointers in node : $O([\text{nb of employees}])$

Interroger un réseau

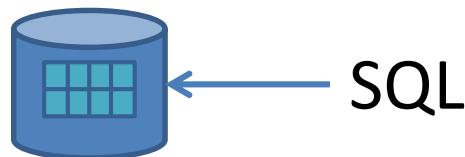
GRAPHDB : LA SOLUTION ?

RDBMS ~~et~~ GraphDB ?

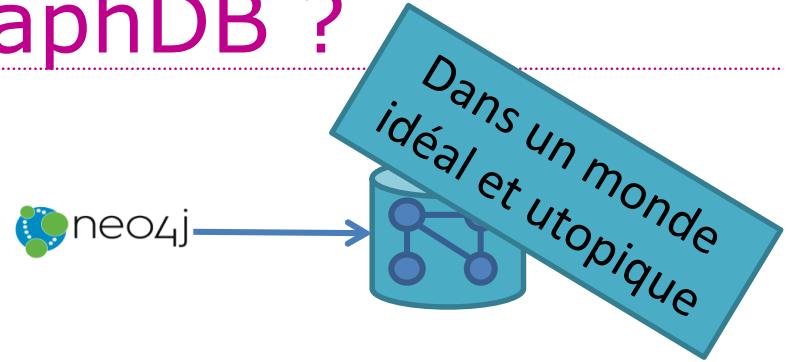
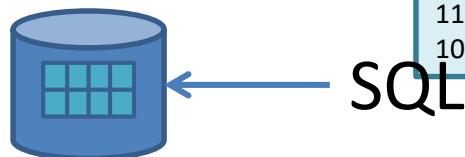
1



2



3



GraphDB : Avantages

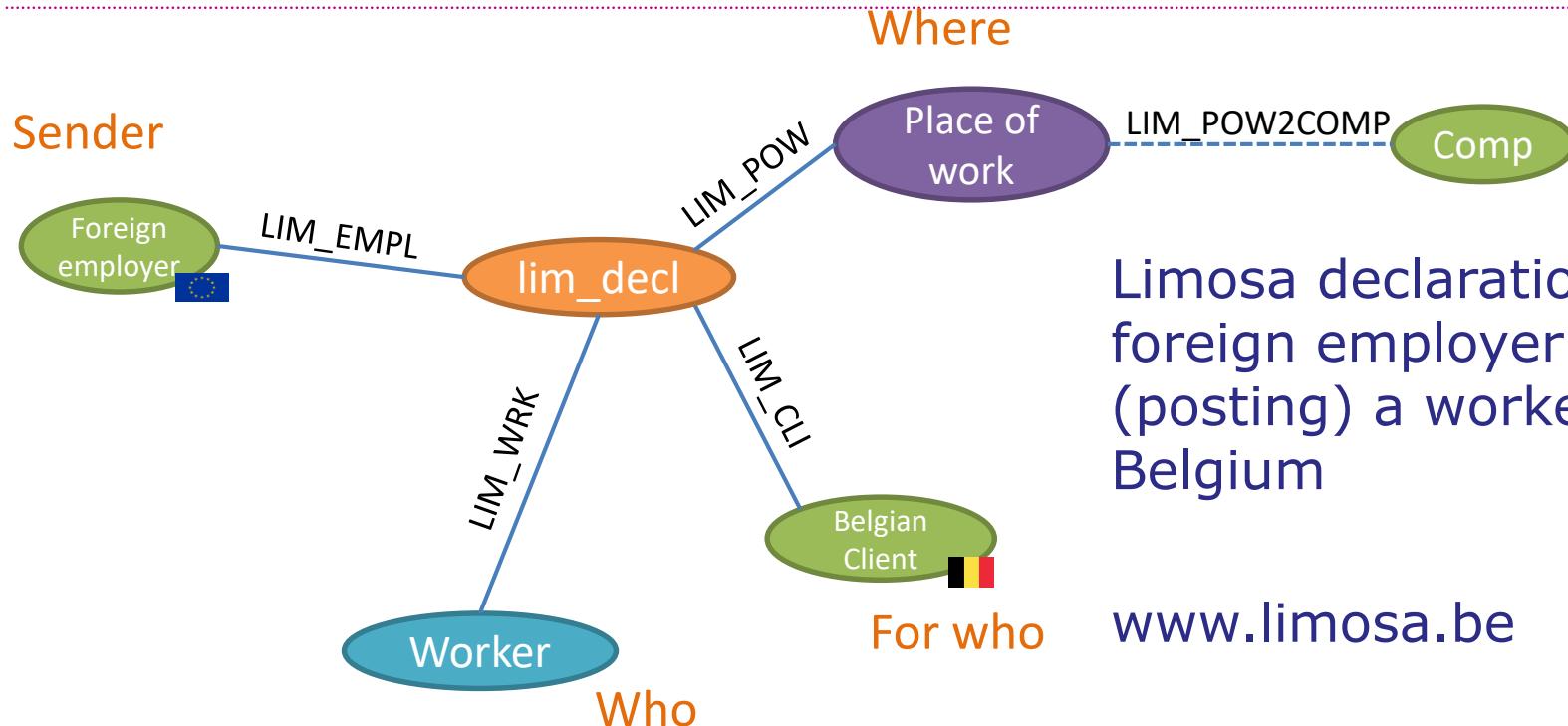
- Langage plus proche du **modèle**
- Pas de **clés** à gérer, ni de contrainte d'intégrité
- Pas de structures de **relations** à chaque requête
- Très efficace pour le **parcours** de relations
- Beaucoup de **champs d'application** : recommandation, fraude, infra, MDM, KM...

GraphDB : limitations

- Pas la **maturité** des RDBMS (robustesse, haute dispo., communauté...)
- Pas de **standard**
- Survie à l'**effet de mode** ? (cf OO DB)
- Pas optimal pour des **recherches, agrégation, batch process, transactionnel**
- Pas une alternative, plutôt un **complément** (ne résout pas les mêmes problèmes)

USE CASE: SOCIAL DUMPING

Limosa: detachment declaration

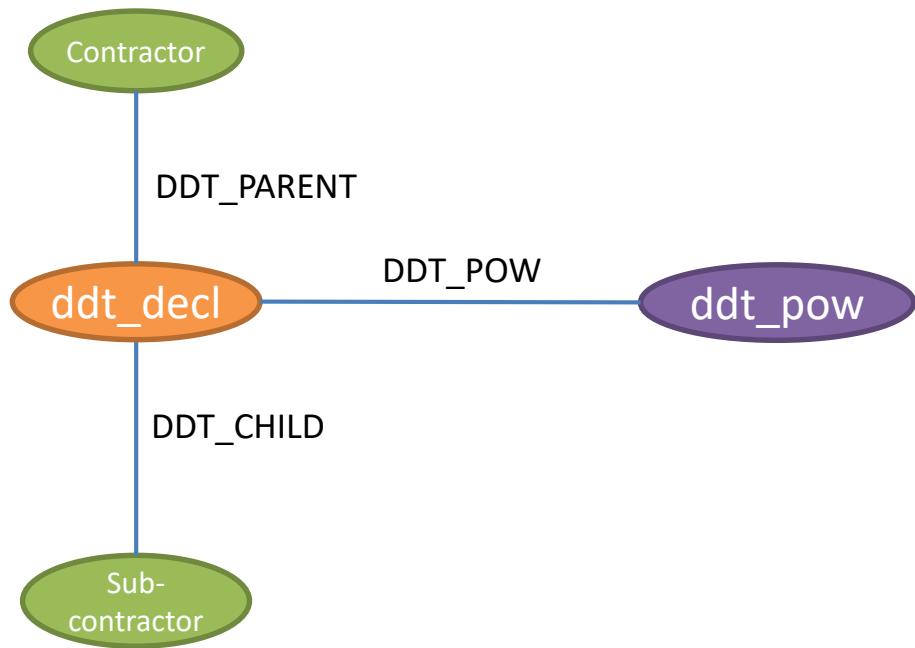


Limosa declaration: A foreign employer sending (posting) a worker in Belgium

www.limosa.be

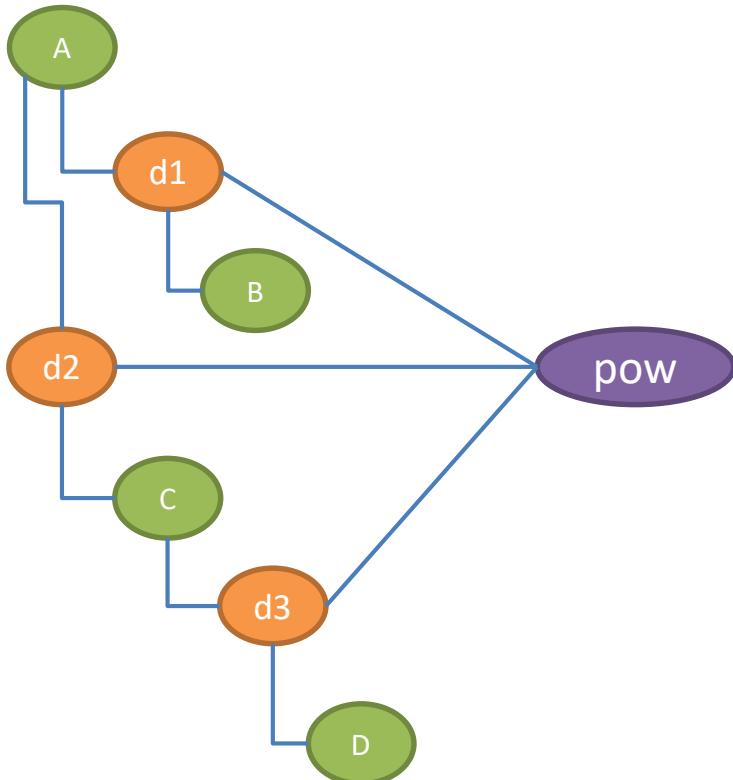
Limosa : Landenoverschrijdend Informatiesysteem ten behoeve van MigratieOnderzoek bij de Sociale Administratie
Syst. d'infor. transfrontalier en vue de la recherche en matière de migration auprès de l'adm. sociale

DDT: sub-contractors



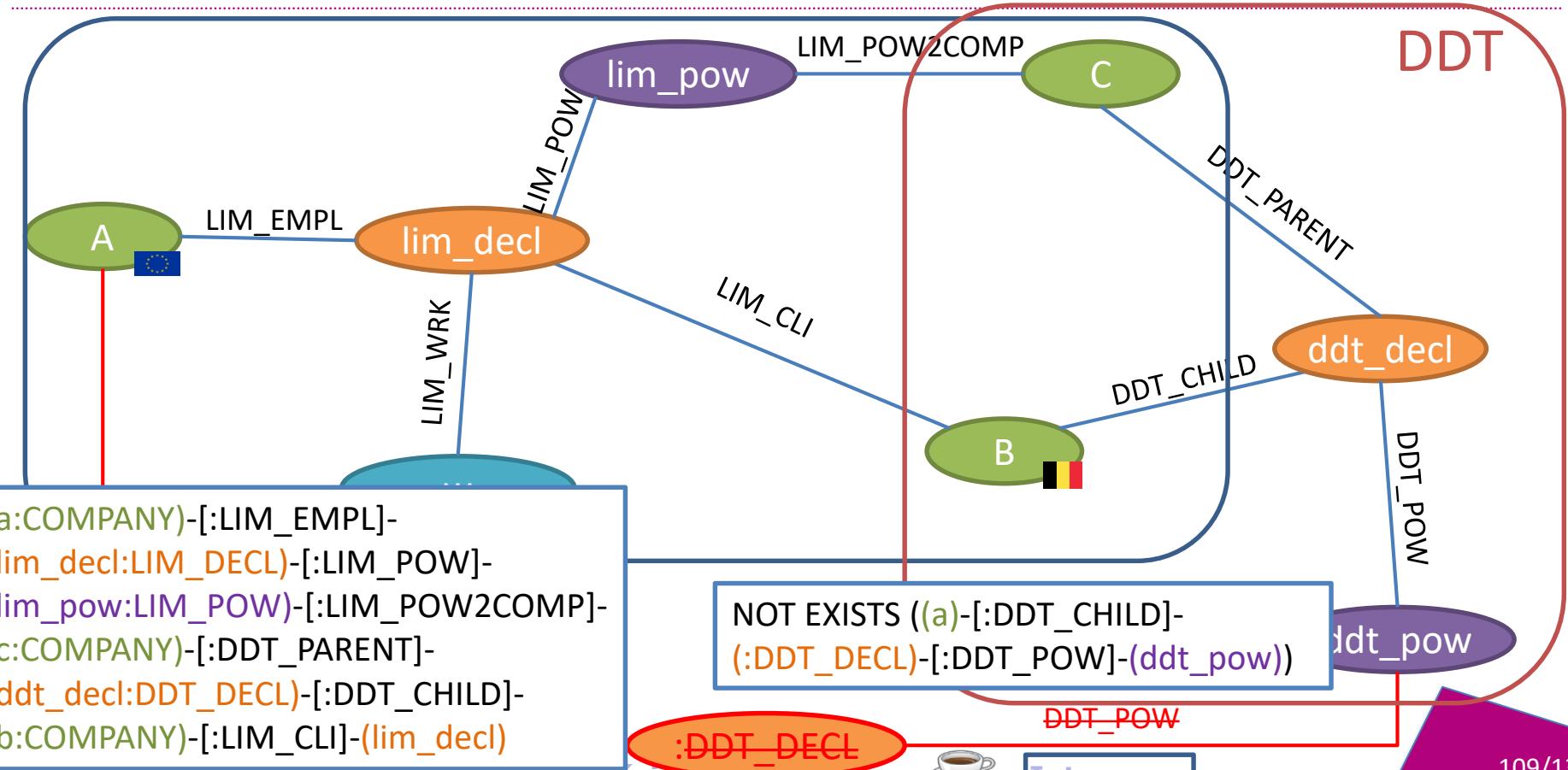
DDT: Déclaration De Travaux
Aangifte van werken

DDT: sub-contractors



On a specific worksite (PoW), we have a chain/tree of sub-contractors

Full schema



Full query

MATCH

```
(a:COMPANY)-[:LIM_EMPL]-(lim_decl:LIM_DECL)-[:LIM_POW]-  
(lim_pow:LIM_POW)-[:LIM_POW2COMP]-(c:COMPANY)-[:DDT_PARENT]-  
(ddt_decl:DDT_DECL)-[:DDT_CHILD]-(b:COMPANY)-[:LIM_CLI]-(lim_decl),  
(ddt_decl)-[:DDT_POW]-(ddt_pow:DDT_POW)
```

WHERE

```
NOT EXISTS ((a)-[:DDT_CHILD]-(:DDT_DECL)-[:DDT_POW]-(ddt_pow)) ...
```

RETURN

```
DISTINCT a, b, c,  
COUNT(DISTINCT lim_decl),  
COUNT(DISTINCT ddt_decl)
```

```
ORDER BY COUNT(DISTINCT lim_decl) DESC
```

Implementation

- We copied in Neo4j 17M nodes, 45M relationships :
 - 800.000 companies, 1.650.000 (foreign) workers
 - 13 millions declarations
 - 185.000 duplication groups
 - 250.000 streets
 - ...
- 9 node labels, 15 relationship types
- All companies triplets matching 1st pattern: ~10 seconds (325 results)

In SQL/RDBMS ?

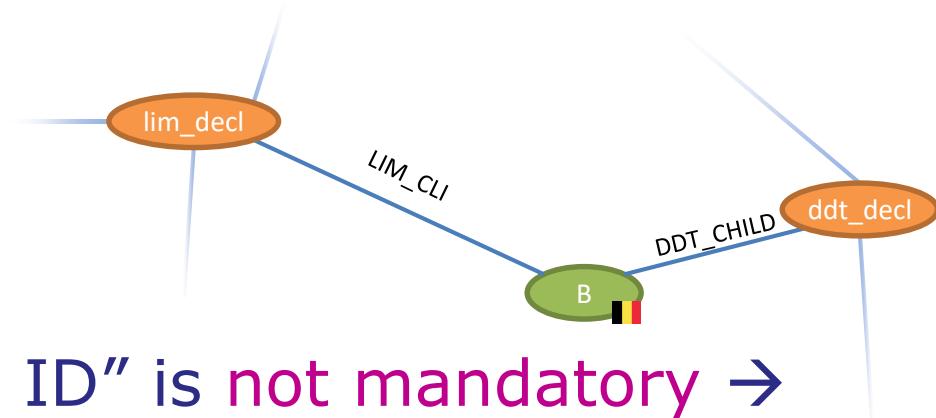
- In SQL/RDBMS,
 - Any node and any relationship is a record : ~65 M records
 - Any type/label corresponds to a table : 24 tables
- Any $(n:\text{Type1})-[:REL1]-(m:\text{Type2})$ corresponds to :

```
FROM Type1
JOIN REL1 ON Type1.key = REL1.type1_key
JOIN Type2 ON Type2.key = REL1.type2_key
```

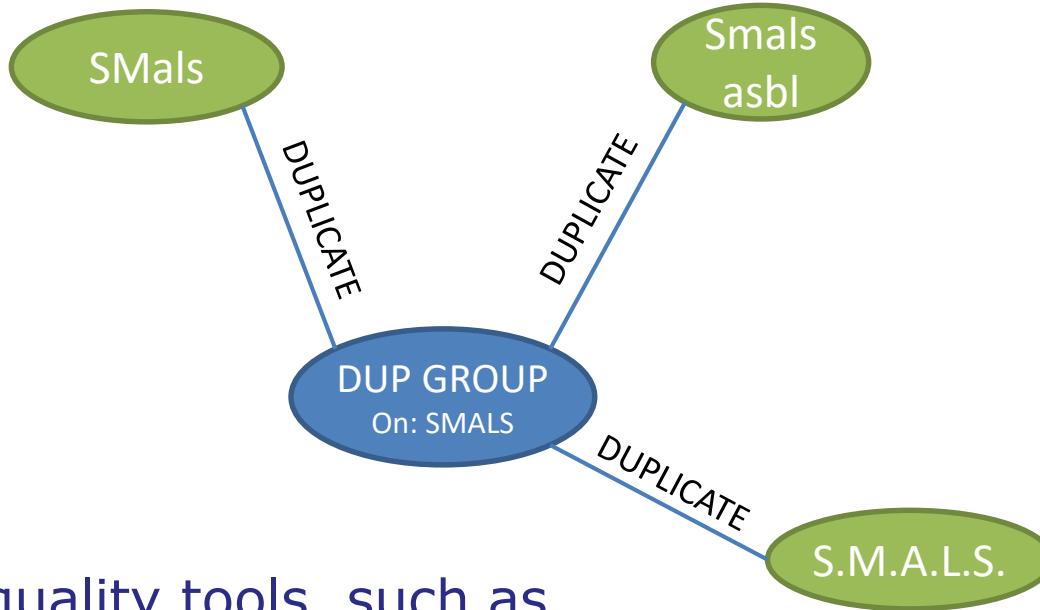
→ Almost 20 joins!
- → Query far more complex to write
- → Much longer to run

Data quality problem

- This only works if « B » has been declared with the **same ID** in both systems !
- Sometimes, the “official ID” is **not mandatory** → name + address is OK
- We have to deal with this!

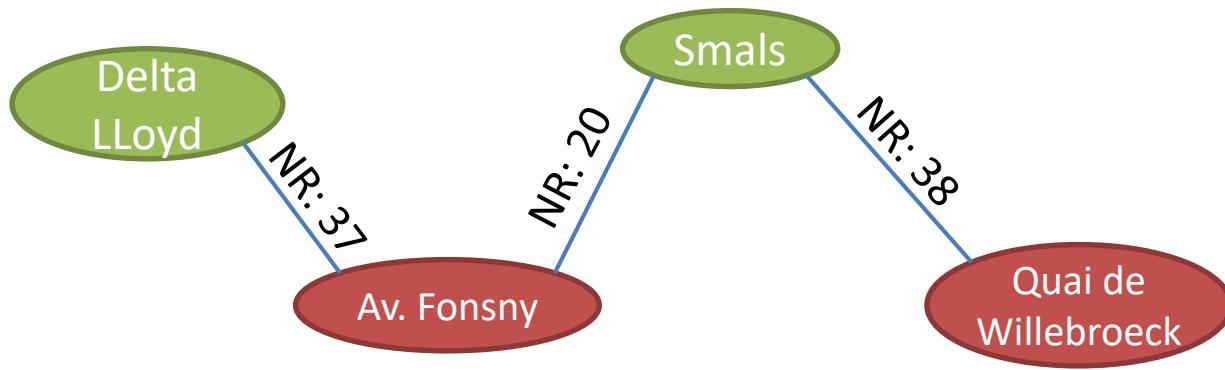


Duplicates



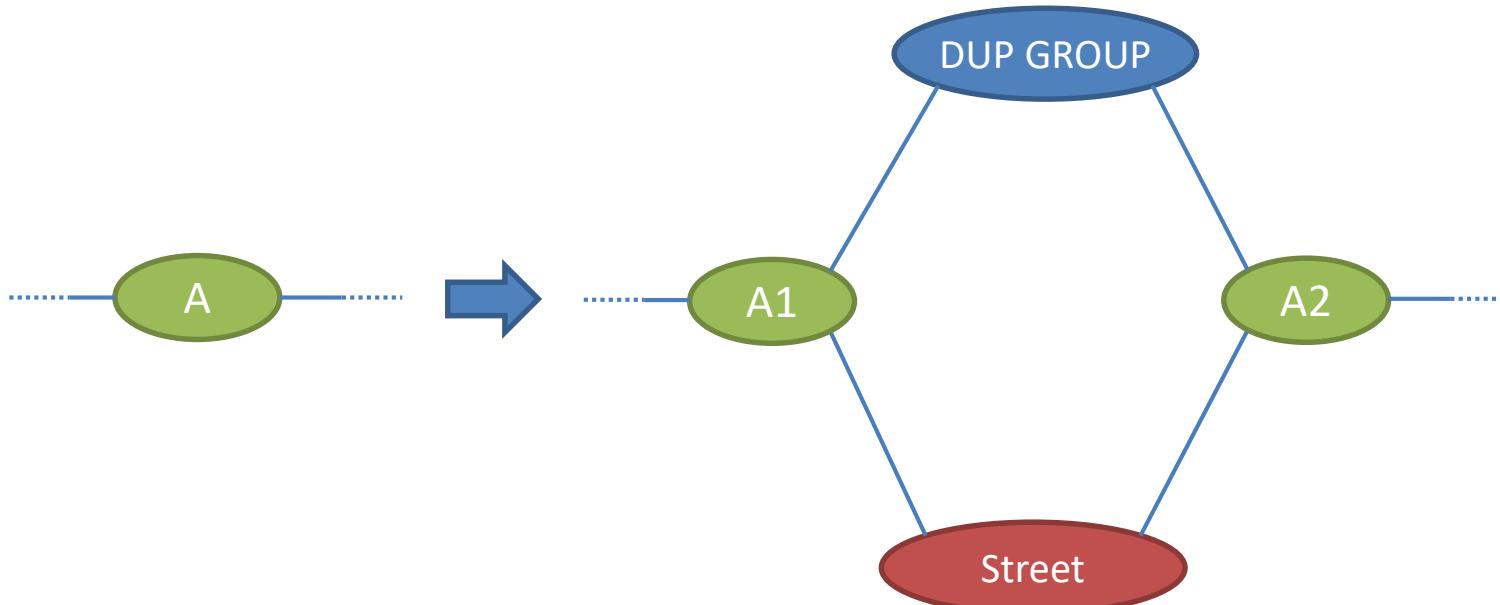
Thanks to data quality tools, such as Trillium (IntoDQ), OpenRefine...

Addresses

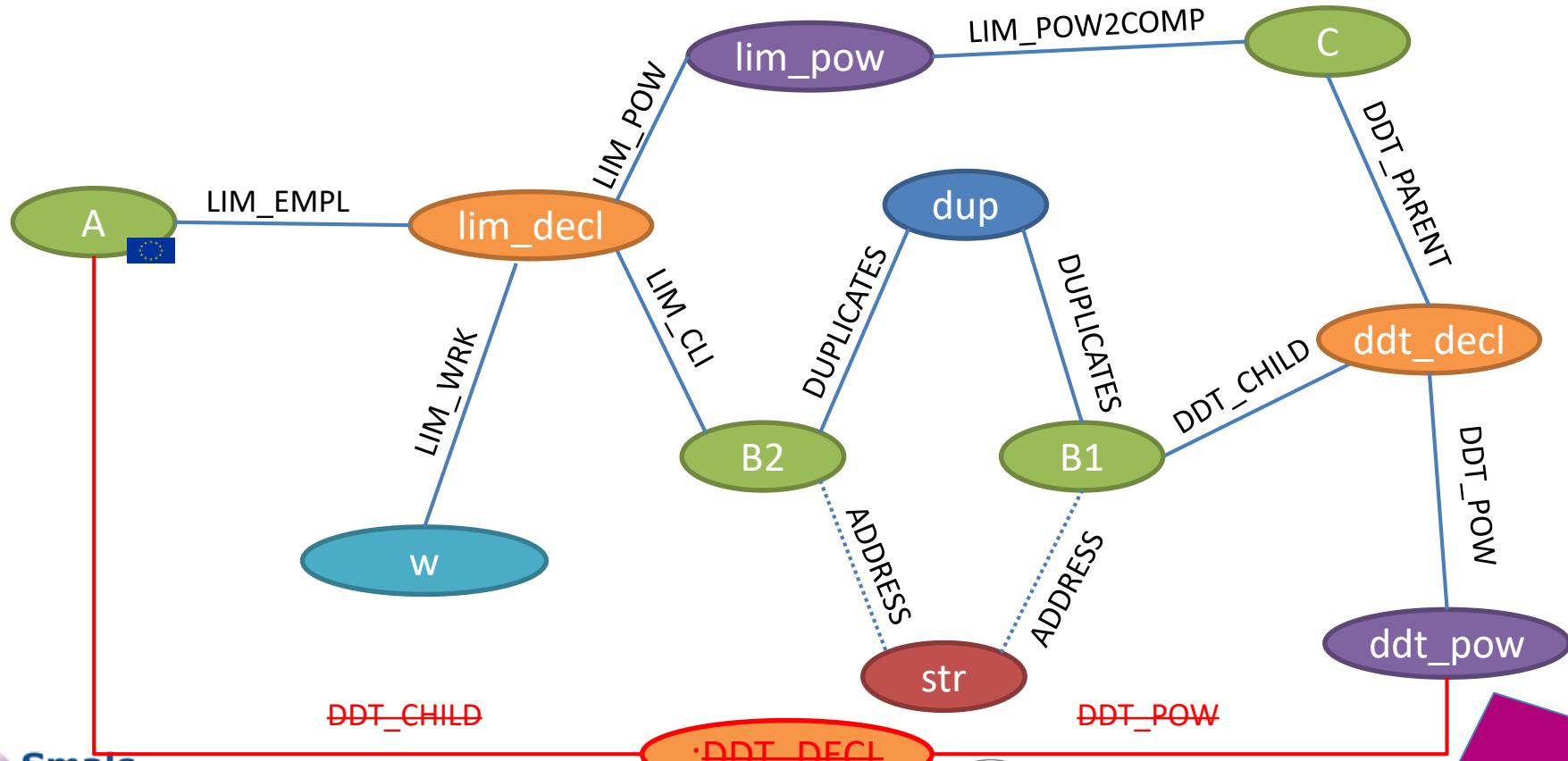


A standardization is needed! Thanks to data quality tools

Duplicates combinations



Schema with B duplicated



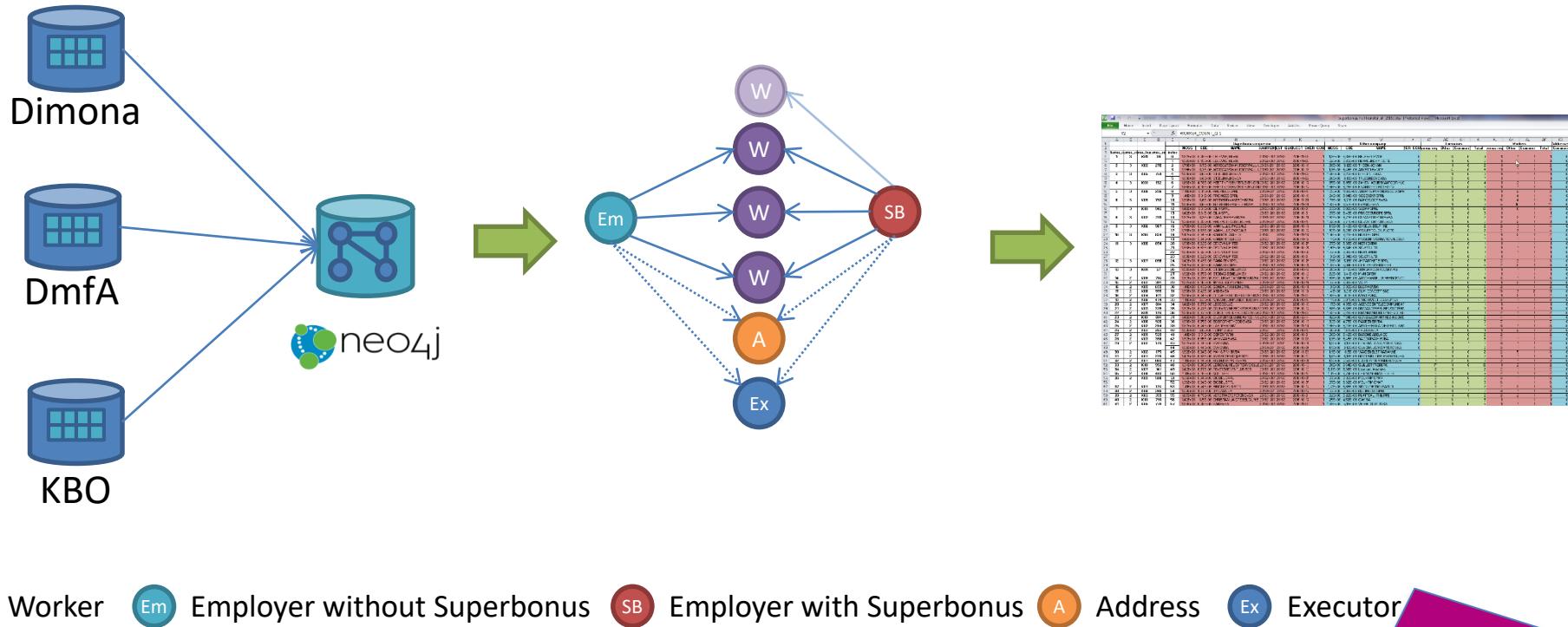
Interroger un réseau

USE CASE : SUPERBONUS

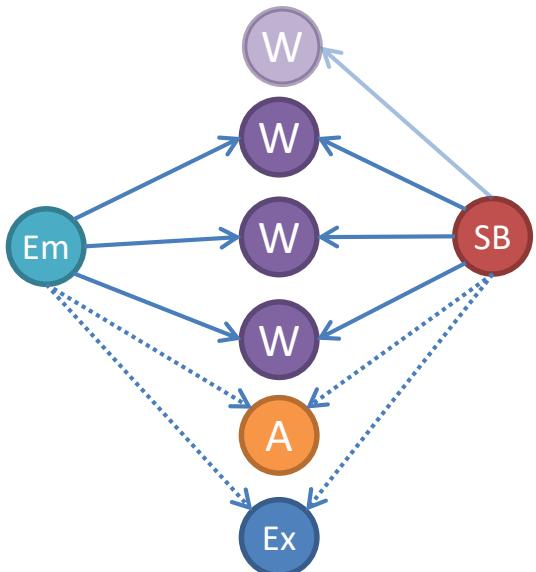
Superbonus : contexte

- Avec le TaxShift (2016), des avantages sont donnés aux **nouveaux employeurs**
- Beaucoup d'employeurs existants **se recréent** pour profiter de ces avantages
- Le « **Network Analytics** » permet d'identifier beaucoup de ces cas

Superbonus



Full transfert



$\text{COUNT}(W) > \text{Em.Worker_Cnt_2015} * 80\%$
 $\text{COUNT}(W) > \text{SB.Worker_Cnt_2016} * 80\%$

```
MATCH      (c_sb:SB) -- (w:Worker) -- (c_old:Company)
WHERE     c_sb <> c_old
OPTIONAL MATCH (c_sb) -- (e:Executor) -- (c_old)
OPTIONAL MATCH (c_sb) -- (a:Address) -- (c_old)
WITH      c_sb, c_old,
          COUNT(DISTINCT w) AS w_count
WHERE
          w_count > c_sb.WORKER_CNT_16 * 0.8 AND
          w_count > c_old.WORKER_CNT_15 * 0.8
RETURN    c_sb, c_old, w_count
```

Worker

Smals
ICT for society

Employer without SB

Employer with SB

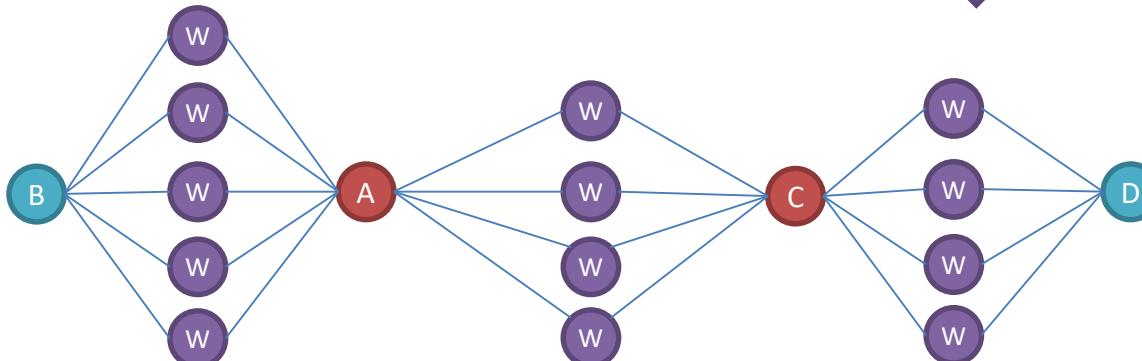
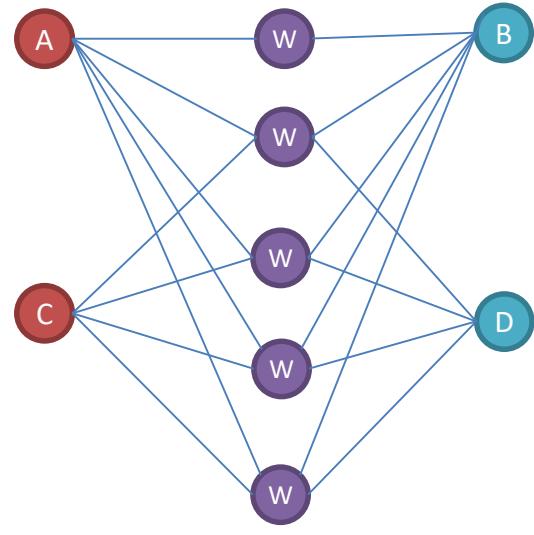
Address Executor

Table vs visualisation

SB req.	Other	Common workers
A	B	5
A	C	4
C	D	4

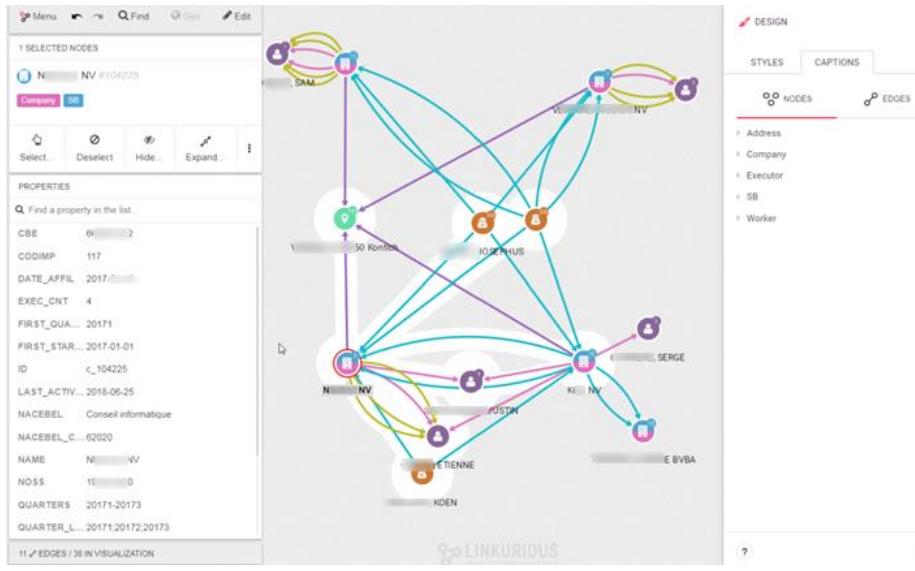


Une table ...
Deux schémas



Linkurious

- Outil de visualisation d'une graph database
- Utilisable par des end users



VISUALIZATION

12 NODES / 609249 IN DATABASE

18 EDGES / 1418444 IN DATABASE

STYLES

CAPTIONS

NODES

EDGES

Company

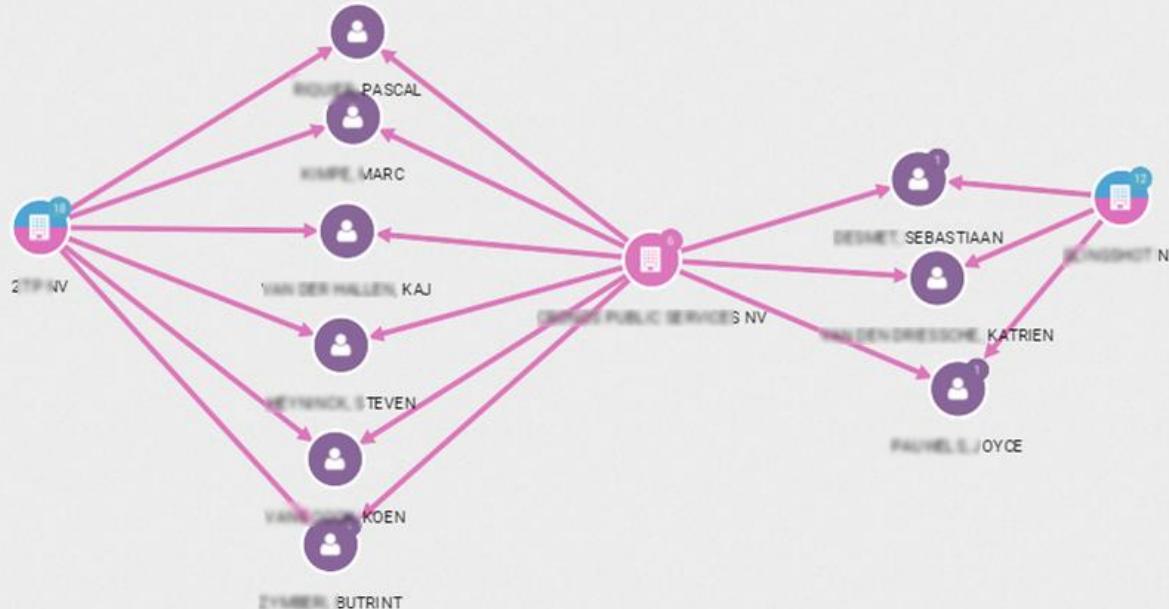
100%

SB

100%

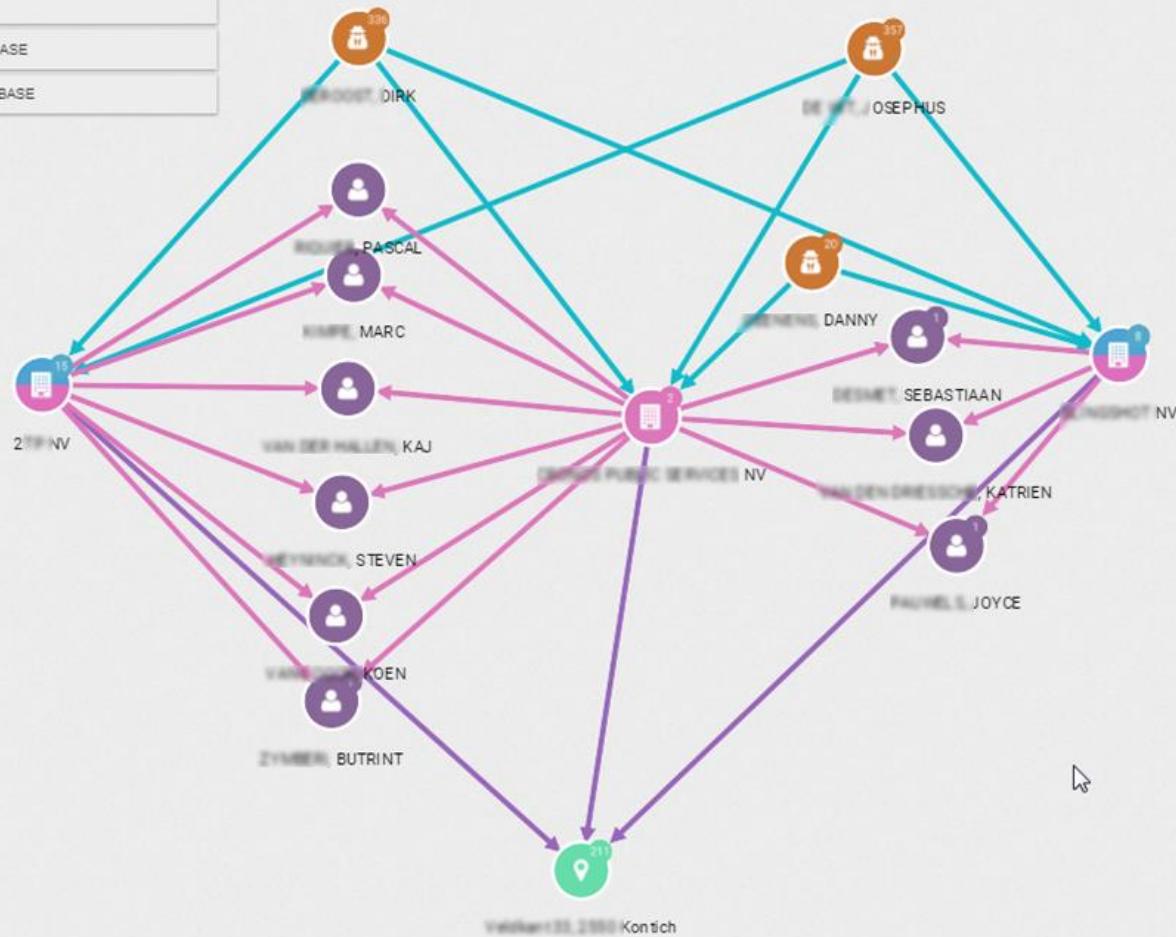
Worker

100%



VISUALIZATION

16 NODGES / 609249 IN DATABASE

29 v² EDGES / 1418444 IN DATABASE

STYLES

CAPTIONS

NODES

EDGES

Address

100%

Company

100%

Executor

100%

SB

100%

Worker

100%



Conclusions

En résumé

Caractériser (théorie) : tout un réseau, un nœud (centralité), deux nœuds (distance), un groupe de nœuds (clustering)

Visualiser : Desktop, Javascript, Cloud

Manipuler (coding) : Envir. classiques ou Big Data

Interroger (graph DB) : focus sur les relations

Application : fraude, gestion des infras/logicielle, recommandation, intelligence, KM...

Conclusions

- De la place pour les graph analytics/graph DB dans nos institutions
- Pas uniquement pour la fraude !
- Pour les data-scientists/BI :
 - nouvelle façon de réfléchir, centré sur les relations
- Pour les DBA/développeurs :
 - nouvelle famille de DB, complémentaires
- Pour les analystes/business :
 - inspiration pour de nouveaux use cases

Conclusions : expérience personnelle

- Préparation des données :
 - R/Python + igraph/graphx
- Première analyse :
 - Gephi/TouchGraph
- Production :
 - GraphDB (Neo4j...)
- Pour le client :
 - Interactif (futur) : librairie JS / outils « à la » Linkurio.us
 - Statique : Gephi/TouchGraph, igraph → Excel, PDF...

Références : www.smalsresearch.be

- <https://www.smalsresearch.be/author/berten/>
- Un fraudeur ne fraude jamais seul ([blog 1](#), [blog 2](#))
- Graph DB vs RDBMS ([blog 1](#), [blog 2](#))
- Gérer les doublons dans une Graph DB ([blog](#))
- Le marché du travail salarié en Belgique : une analyse réseau ([blog 1](#), [blog 2](#), [blog 3](#))
- 7 raisons d'utiliser une Graph Database ([blog](#))

Références

- **Graph Databases**, *Robinson & all*, O'Reilly 2015
- **Fraud Detection: Discovering Connections with Graph Databases**, *Neo4j*, Sodawksi & Rathle
- **Fraud Analytics Using Descriptive, Predictive, and Social Network Techniques: A Guide to Data Science for Fraud Detection**, *Bart Baesens, Veronique Van Lasselaer, Wouter Verbeke*, Willey, 2015



Vandy Berten
02/787.57.32
vandy.berten@smals.be



More on Smals Research :
Website : www.smalsresearch.be
Blog : www.smalsresearch.be/blog
Twitter : [@SmalsResearch](https://twitter.com/SmalsResearch)

