



Détection de dispositifs de publication coordonnée dans l'analyse des flux de Twitter



Florence Andréacola



Jean-Marc Francony



Philippe Mulhem,
Lorraine Goeuriot
Georges Quénot



Événementiel Médiatique

- Il fait se rejoindre les pratiques professionnelles des industries de l'information, avec les pratiques amateurs et de *supporterisme*.
- Il constitue une opportunité de mise en lumière pour des acteurs extérieurs à l'événement programmé qui peuvent s'en saisir (*effet d'aubaine*).
- L'objectif est alors d'atteindre une audience plus large et/ou de percoler la sphère médiatique en s'inscrivant dans l'événement et/ou en provoquant l'actualité.

#LONDON2017, #LONDRES2017



- Les compétitions sportives de haut niveau constituent l'archétype de l'événementiel médiatique récurrent qui mobilise dans la durée une très large audience.
- Elles ont toujours été l'occasion d'innovations technologiques et techniques et désormais le reflet d'une communication multicanale de flux.
- Les composantes esthétiques ou dramatiques de la performance athlétique favorisent une communication visuelle abondante.

Une problématique : *quelle information apporte le flux média dans l'analyse du flux de publication et dans la compréhension de l'événement ?*

Publiée 9765 fois dans la collection

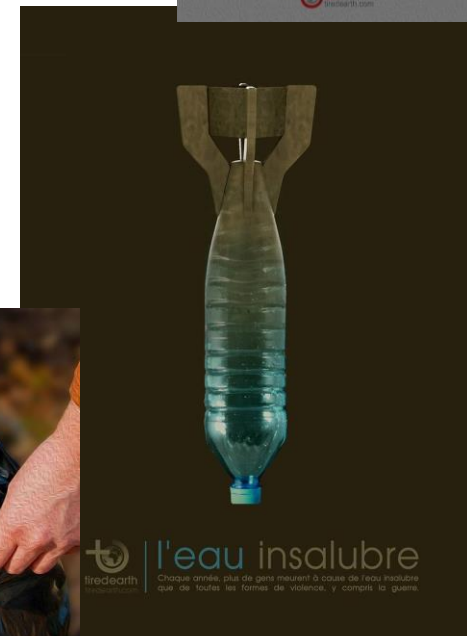
Un phénomène éditorial : TiredEarth

Plus d'de 33% du flux de fichier unique est constitué d'images d'une même collection : TiredEarth

77 images uniques publiées environ 300 fois.

Un militantisme au service d'une cause environnementale.

Un *dispositif* de publication au service d'une campagne d'*influence*.



Missolin963

Créé en mai 2017

13219 Tweets publiés
+/- 33 Tweets journaliers

Fréquence horaire < 5 Tweets

StreamHoraire - Kibana x Common options | Elast... x Karine Missolin (@missolin963) x

Twitter, Inc. [US] | https://twitter.com/missolin963

Applications ArangoDB Web Inter... Garmin Connect Projets - Dashboard Kibana Projets libres de rech... phpMyAdmin Noisi - Improve Foc... discussion_comment... Accueil PapersPublications < PRONOTE, Logiciel < Banque Populaire. Ac... Intranet Google Scholar >>

Accueil Notifications Messages

Recherchez sur Twitter

Tweets 13 k Abonnements 972 Abonnés 244 J'aime 126

Suivre

Karine Missolin
@missolin963
Climate change impacts on our lives extremely.
Inscrit en mai 2017

Tweeter

12.9 k Photos et vidéos

Tweets Tweets & réponses Médias

Karine Missolin @missolin963 · 8 févr.
#environnement #nature #pollution #arbre #ecologie youtube.com/watch?v=fjqWkw...

Est-ce que vous savez que près de 13 kilos de nitrates en excédent qui se déversent dans l'environnement en France chaque seconde, soit 400 000 tonnes par an ?

Suggestions Actualiser · Tout afficher

CORUM @corum_FR Suivre Sponsorisé

Cold Forged @Downb... Suivre

Leaving a Legacy @Leave... Suivre

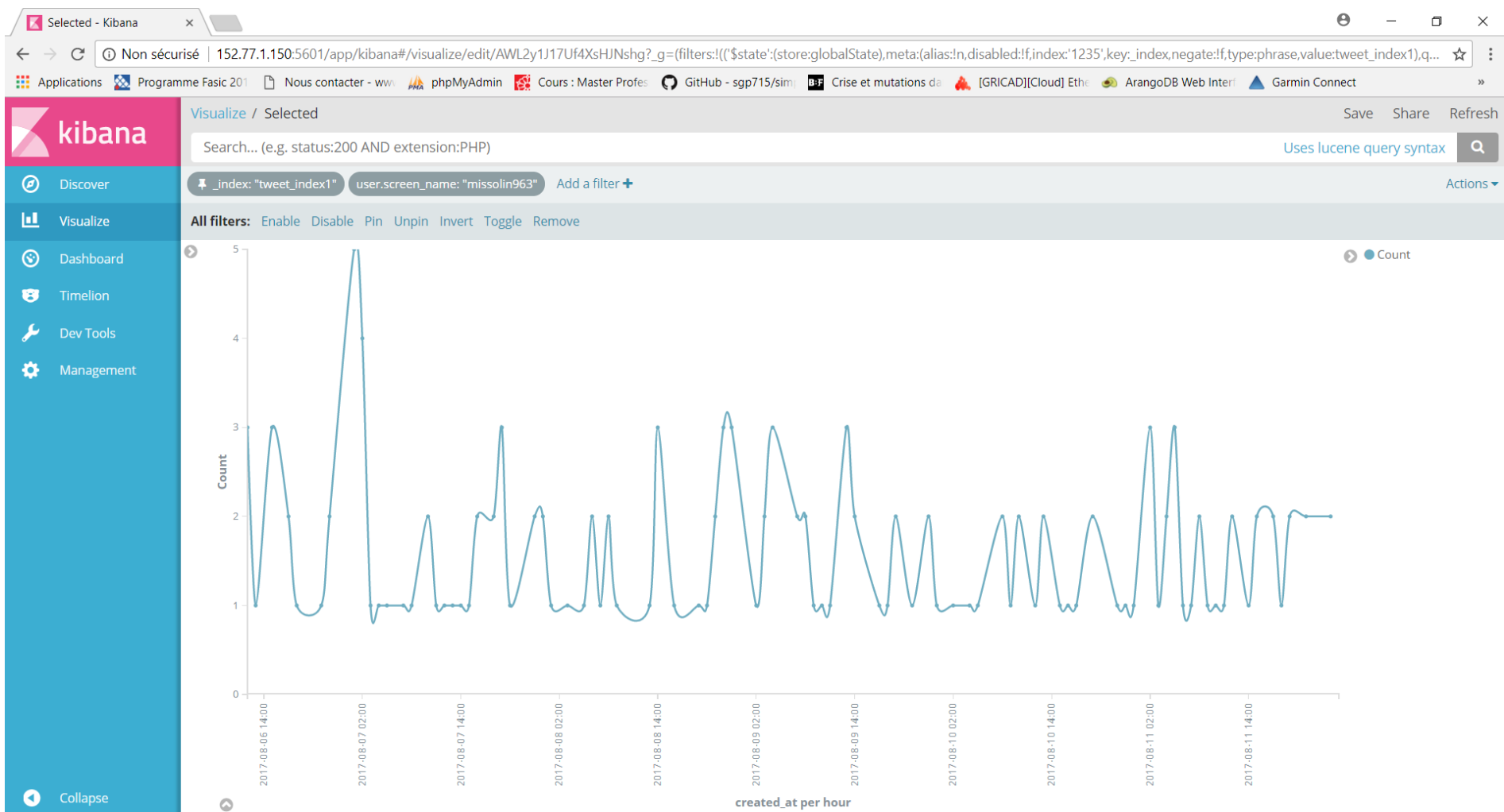
Trouvez vos connaissances

Tendances pour vous Modifier

#loscel 14,3 k Tweets

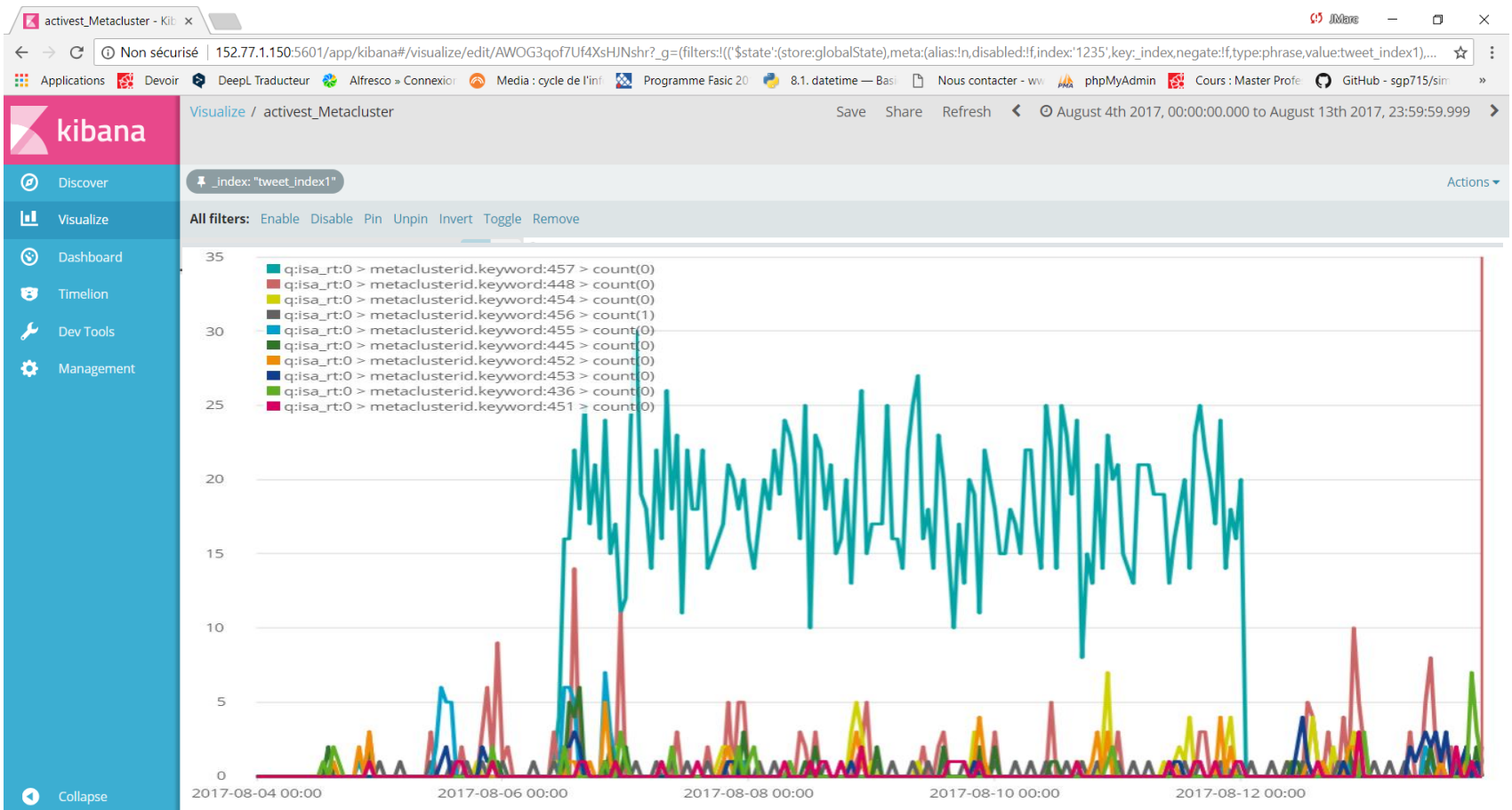
#OMFCGB

Une activité naturalisée



Mais un bruit de fond assourdissant

Des conséquences interprétatives sur le flux événementiel s'il n'est pas détecté.



Dispositif de Publication Concertée

- Il mobilise un ensemble de comptes personnels dont l'activité de publication non contrevenante (a priori) est coordonnée :
 - Il repose sur l'optimisation de la publication individuelle sous contrainte réglementaire (contractuelle).
 - Il vise l'amplification d'un flux de publication tout en restant indécélable en tant que système organisé.
- Il sert la visibilité en intensité et en durée dans les flux d'actualité de Twitter.
- Il répond aux enjeux marketing dans les différents domaines commerciaux, politiques ou idéologiques :
 - de recrutement (*lead generation*),
 - d'image et de notoriété,
 - de transformation de l'audience en action (communication ou autre).

Formalisation des données

- Nous avons donc considéré exclusivement le flux de publications originales (28%) et non sur celui de rediffusion (*retweet*) ;
- Nous avons filtré le flux en ne retenant que les tweets comportant un fichier média (59%) ;
- Pour chacun de ces fichiers média nous avons calculé la signature MD5 afin d'identifier les images identiques.
- La collection résiduelle constitue un graphe bipartite:
Tweet : auteur x image (md5)

Postulat incrémental

- L'objectif poursuivi est celui d'une détection au fil de l'eau et au plus tôt d'un DPC.
- Cela se traduit dans la mise en œuvre d'une algorithmique incrémental suivant un pas temporel pertinent (ici la journée).
- On travaille donc à partir d'un graphe qui incorpore progressivement les relations établies dans la chronologie des journées.

Regroupements par co-clustering

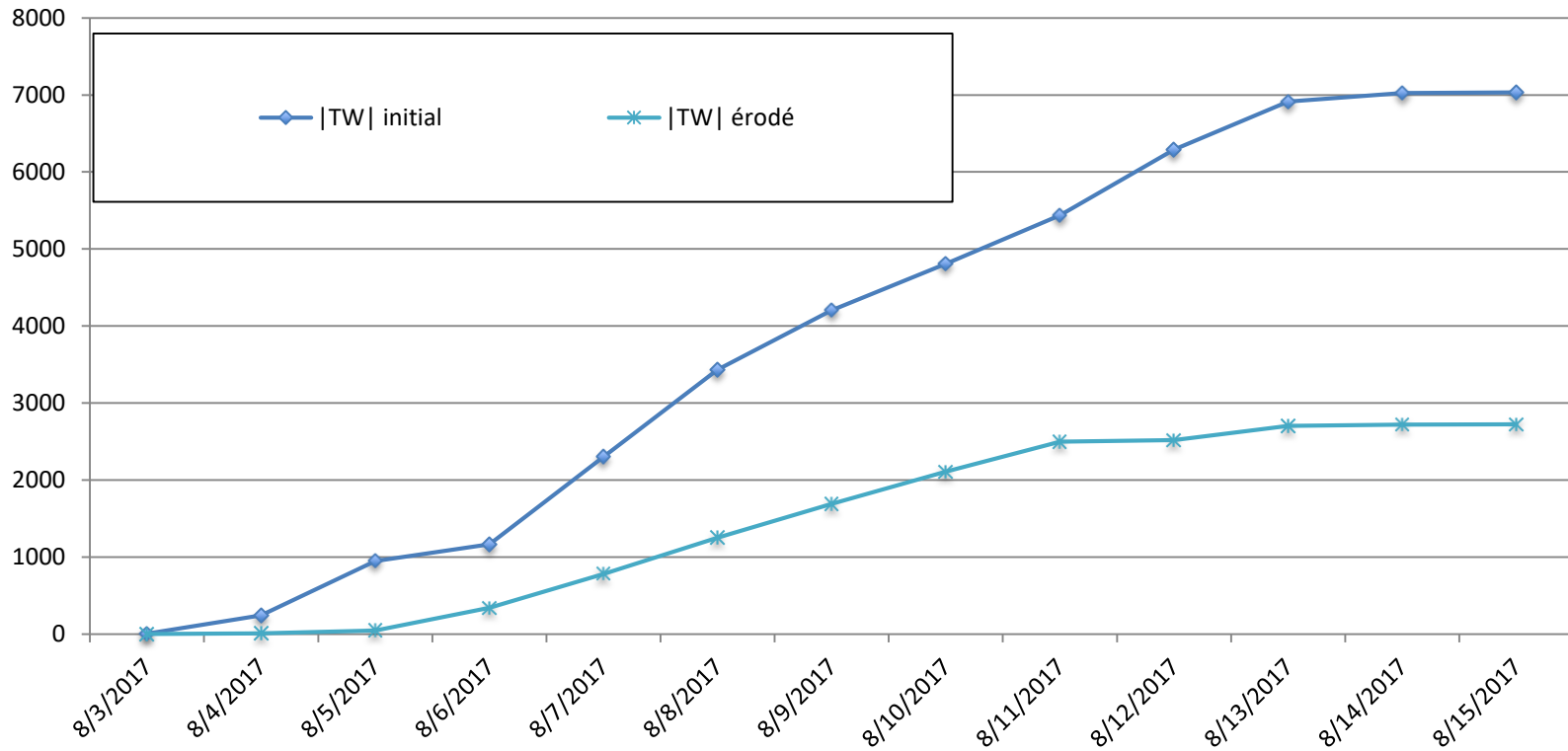
- S'agissant d'un graphe bipartite, nous avons privilégié les approches de co-clustering mis en œuvre suivant les dimensions auteur et images.
- L'approche retenue est celle des modèles de blocs latents [Govaert et Nadif 2003]
- L'inconvénient est que ces modèles sont gourmands en temps de calcul, même sous forme optimisée.
- D'où la mise en place d'une étape d'*érosion*.

Érosion

- L'Objectif de ce traitement est l'élimination rapide des éléments peu connectés (*outliers*) dans le graphe bipartite
- Le principe consiste à retirer itérativement les auteurs et les images ayant zéro ou un seul voisin dans le graphe jusqu'à convergence.
- Compte tenu de l'incrémentation, l'érosion est sans mémoire, c'est-à-dire recalculée sur l'ensemble du graphe historique au jour J.

Érosion - gain

Comparaison nombre de Tweets concernés



Méthode de co-clustering

- L'algorithme de co-clustering imposant de fixer *a priori* un nombre de classes une heuristique d'ajustement est proposée :
 1. Suivant une incrémentation itérative du nombre de clusters demandés.
 2. Évaluer l'évolution du nombre de clusters obtenus en fonction du nombre demandé
 3. Choisir une valeur demandée obtenant une stabilité (seconde valeur de plateau) du nombre obtenu

Formalisation des données

- Ensemble des tweets : TW
- Filtrage de tweets contenant une image ou plus : TWI
 - Un tweet-image twi de TWI est identifié par la signature **MD5** de l'image attachée au tweet émis
- Filtrage des utilisateurs de U qui ont émis au moins une image de TWI : UI
- Graphe bipartite tweet-image/utilisateur $GBTI$ dans $TWI \times UI$

Regroupements par co-clustering

- DPC : regroupe des utilisateurs de UI ayant émis des twi de *TWI* en commun
- Comme on a un graphe bipartite, on utilise des approches de co-clustering sur un graphe binaire, suivant les dimensions des images et des utilisateurs.
 - Le clustering se base sur les 2 dimensions
 - Approche par modèles de blocs latents [Govaert et Nadif 2003]

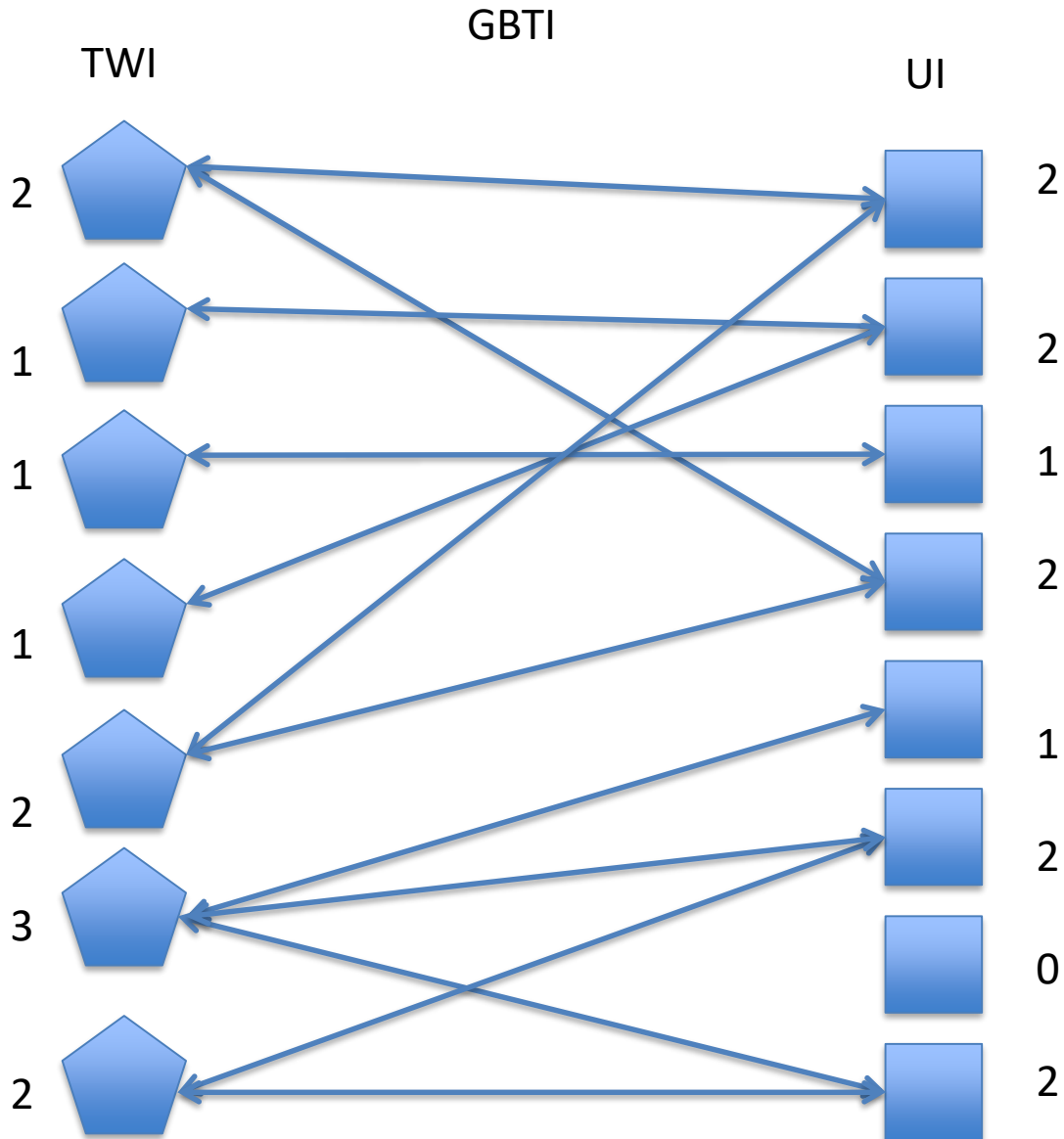
Co-clustering par modèles de blocs latents

- Fonction de densité générale [Govaert et Nadil 2008]

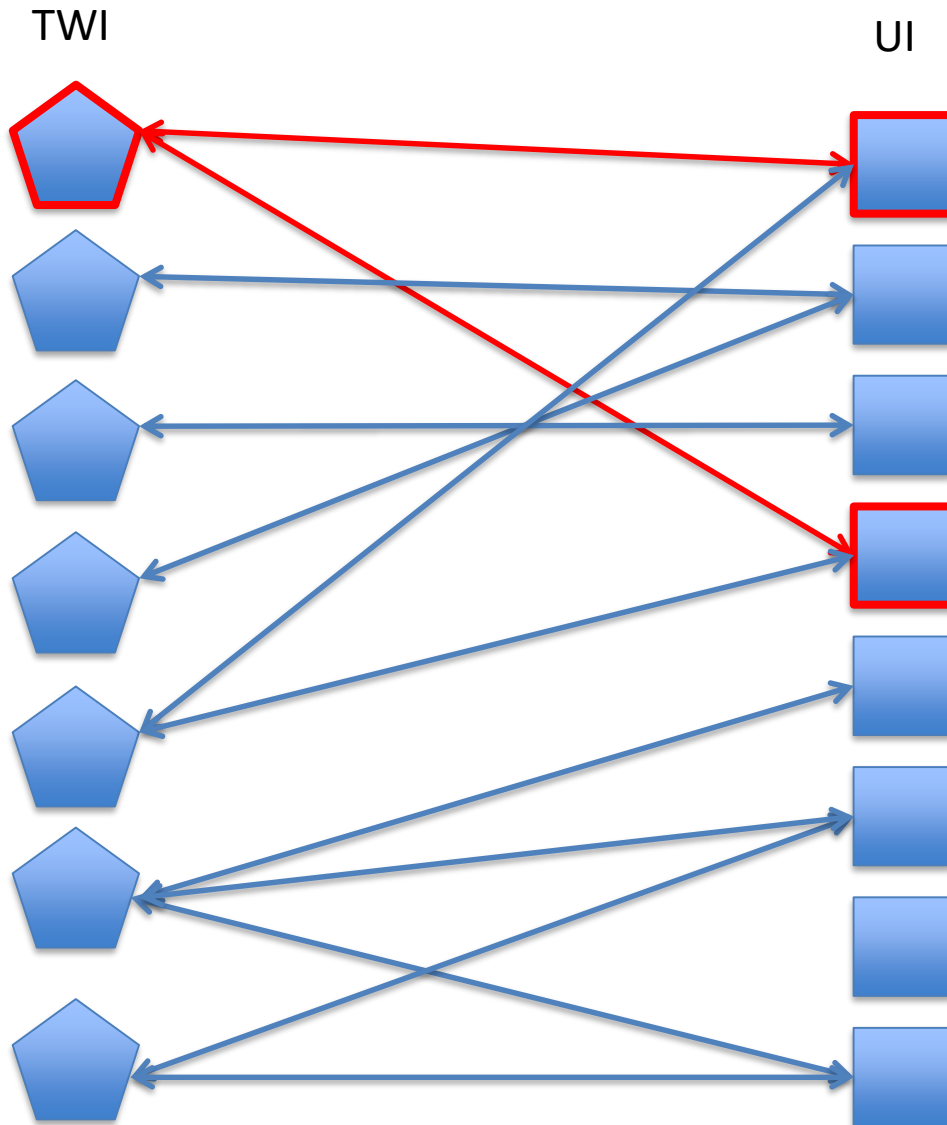
$$f(\mathbf{x}; \theta) = \sum_{(z,w) \in \mathcal{Z} \times \mathcal{W}} \prod_i p_{z_i} \prod_j q_{w_j} \prod_{i,j} \varphi(x_{ij}^j; \alpha_{z_i w_j}).$$

- i : colonnes, j lignes
 - $\mathcal{Z} \times \mathcal{W}$: ensemble de toutes les partitions possibles
 - $p = (p_1, \dots, p_g)$: proba. de colonne pour cluster colonne $[1, d]$
 - $q = (q_1, q_m)$: proba. de ligne pour cluster ligne $[1, m]$
 - Φ : densité "latente" de probabilité pour $x_{i,j}$ sachant les paramètres alpha "conjointes" sur les lignes et les colonnes (données binaires : distributions de Bernoulli)
- Optimisation pour trouver la meilleure partition (z,w) avec paramètres alpha, par une variante de la maximisation de l'espérance appelée CEM

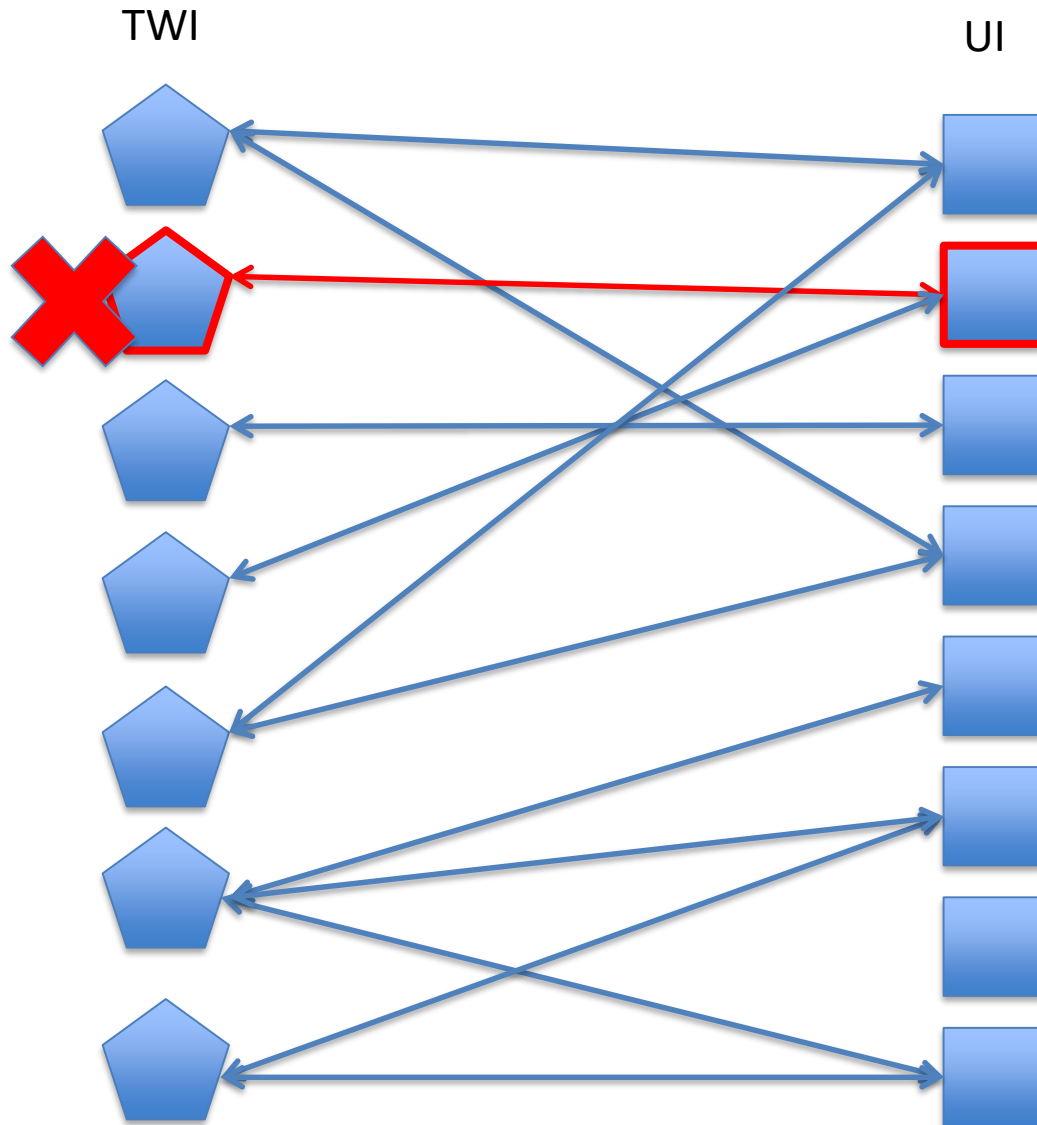
Erosion - Exemple



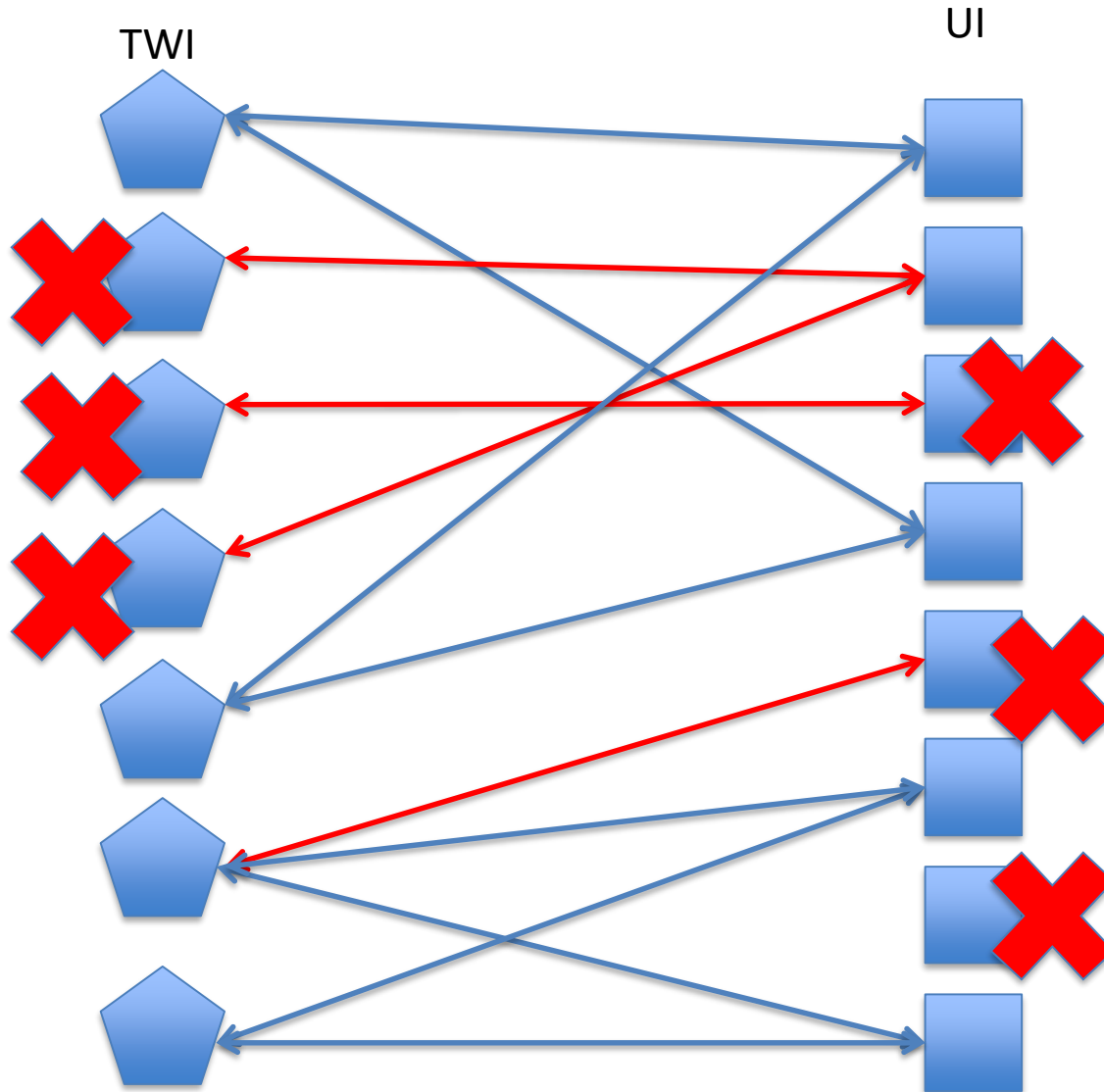
Erosion – Exemple – Boucle 1



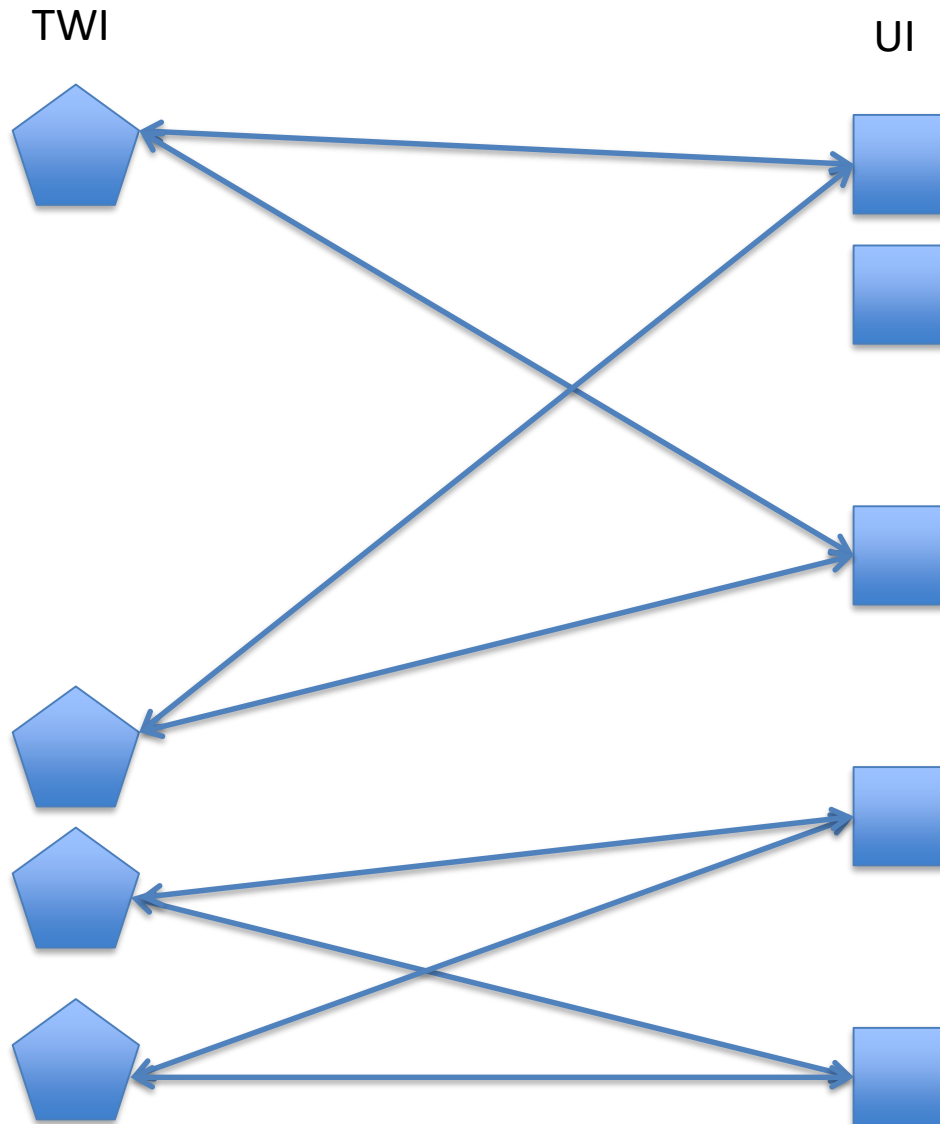
Erosion – Exemple – Boucle 1



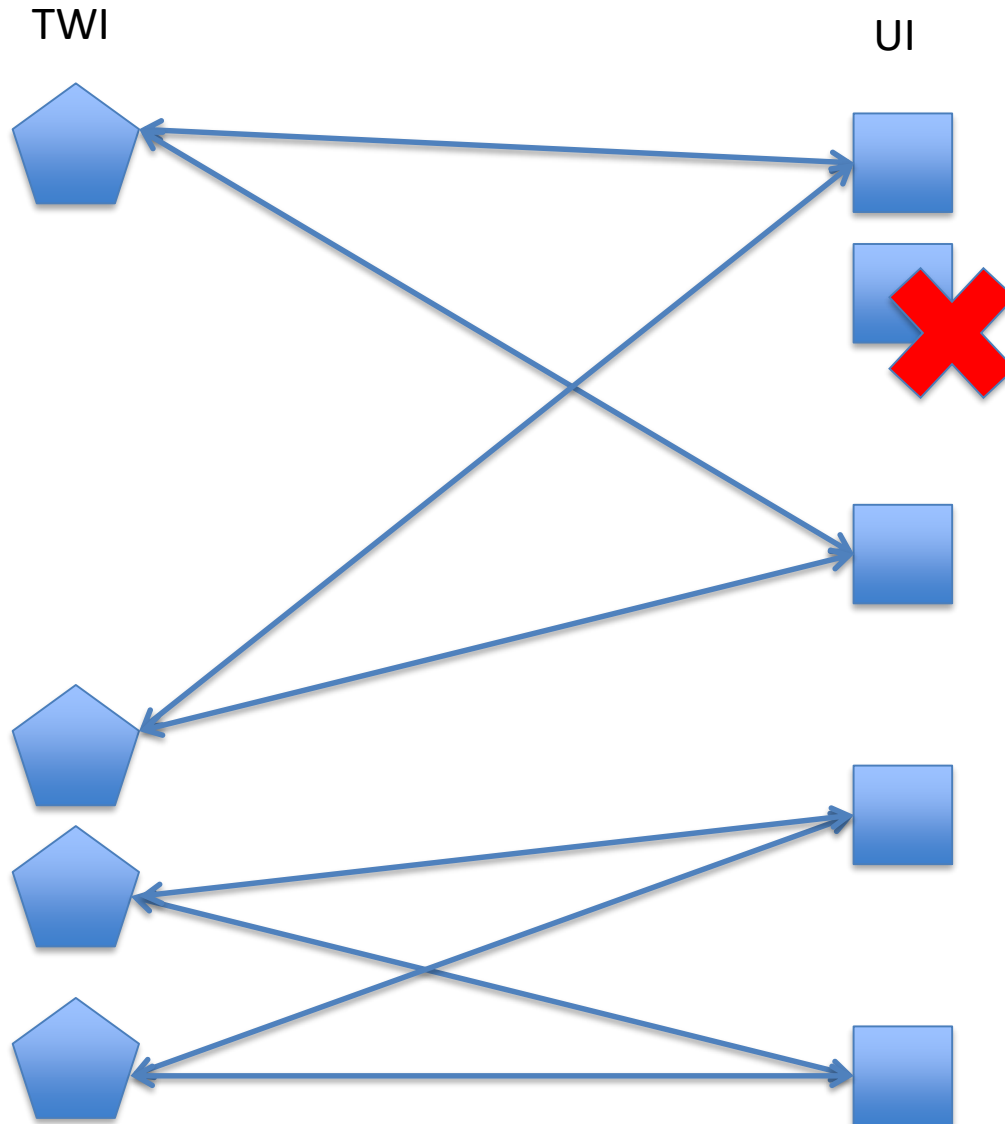
Erosion – Exemple – Boucle 1



Erosion – Exemple – Boucle 1

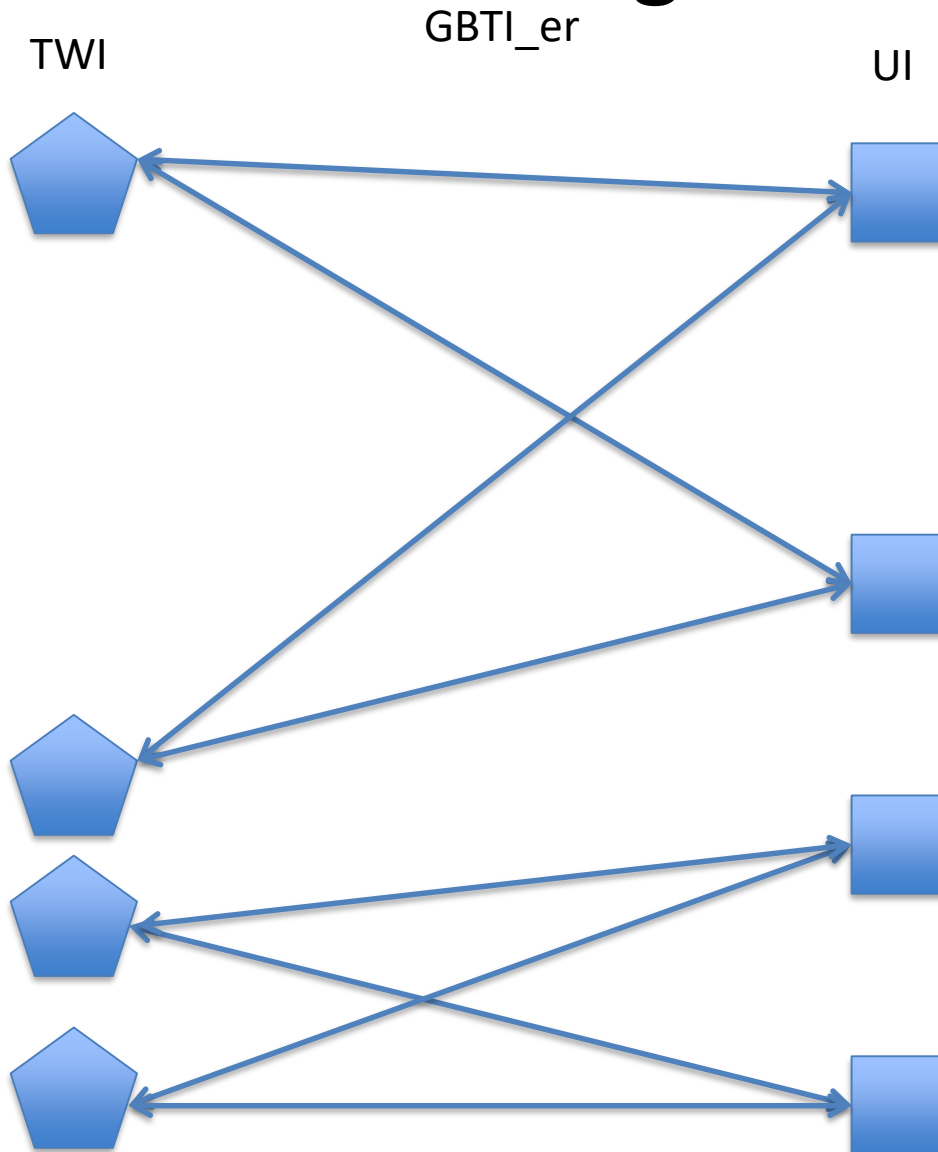


Erosion – Exemple – Boucle 2



Erosion – Exemple – Boucle 2

convergence

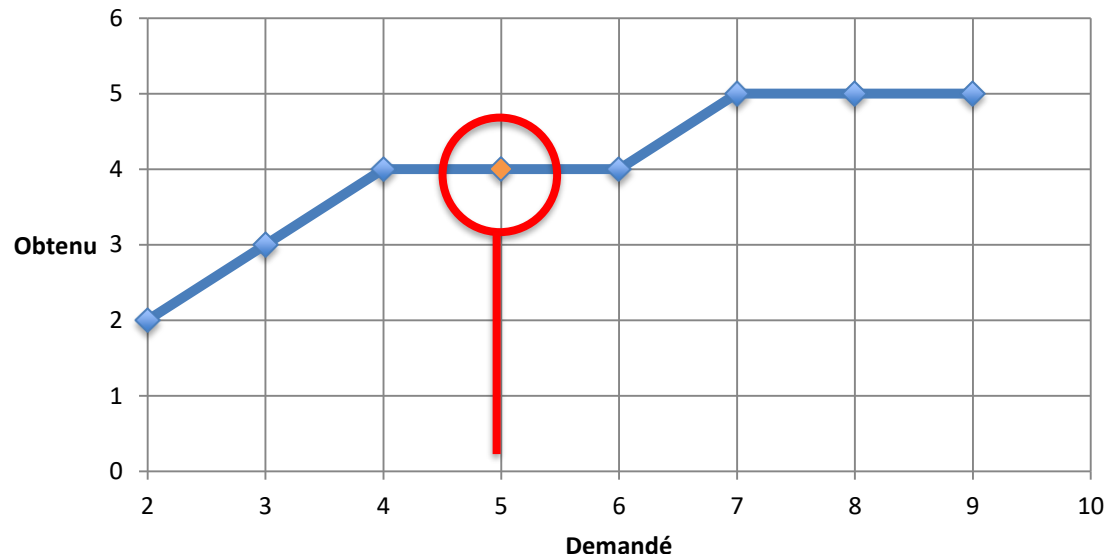


Co-clustering – Nombre de clusters

- Exemple d'évolution

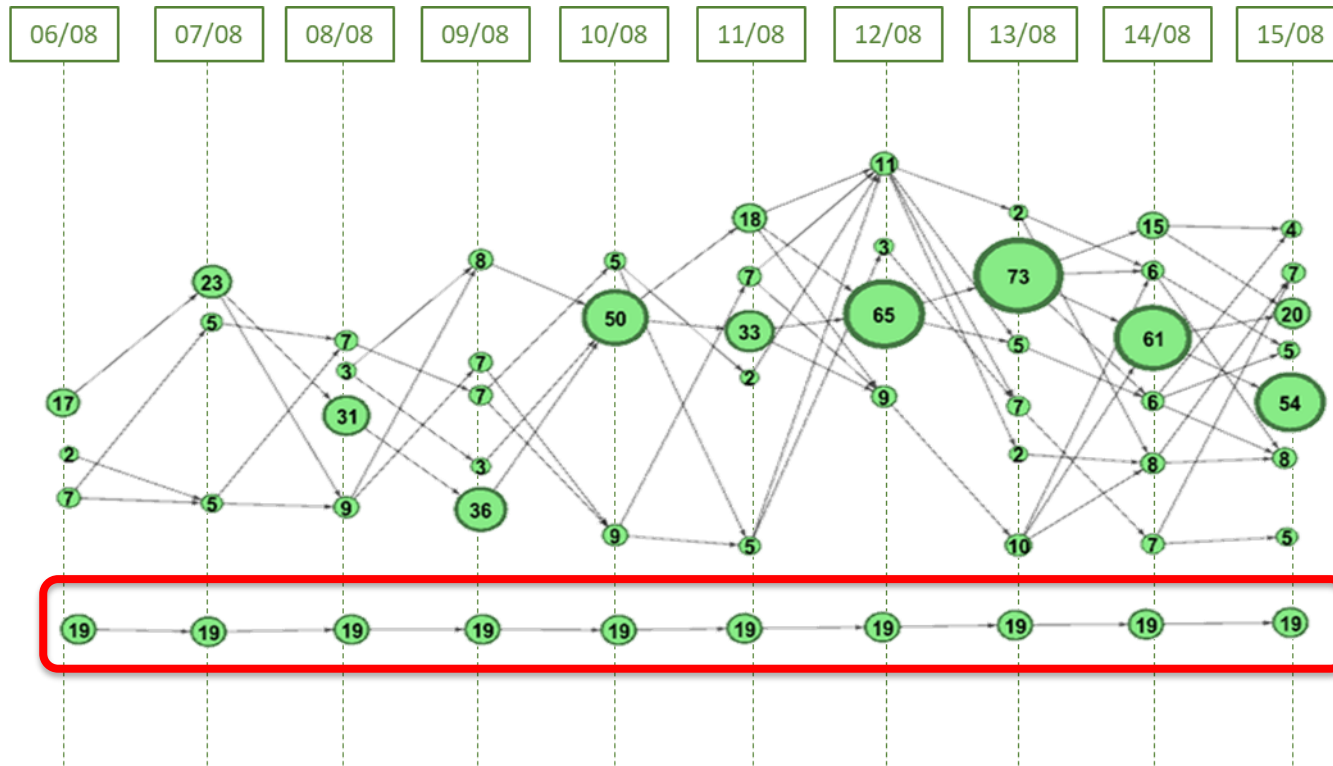
- Jour 6 août 2017

8/6/2017

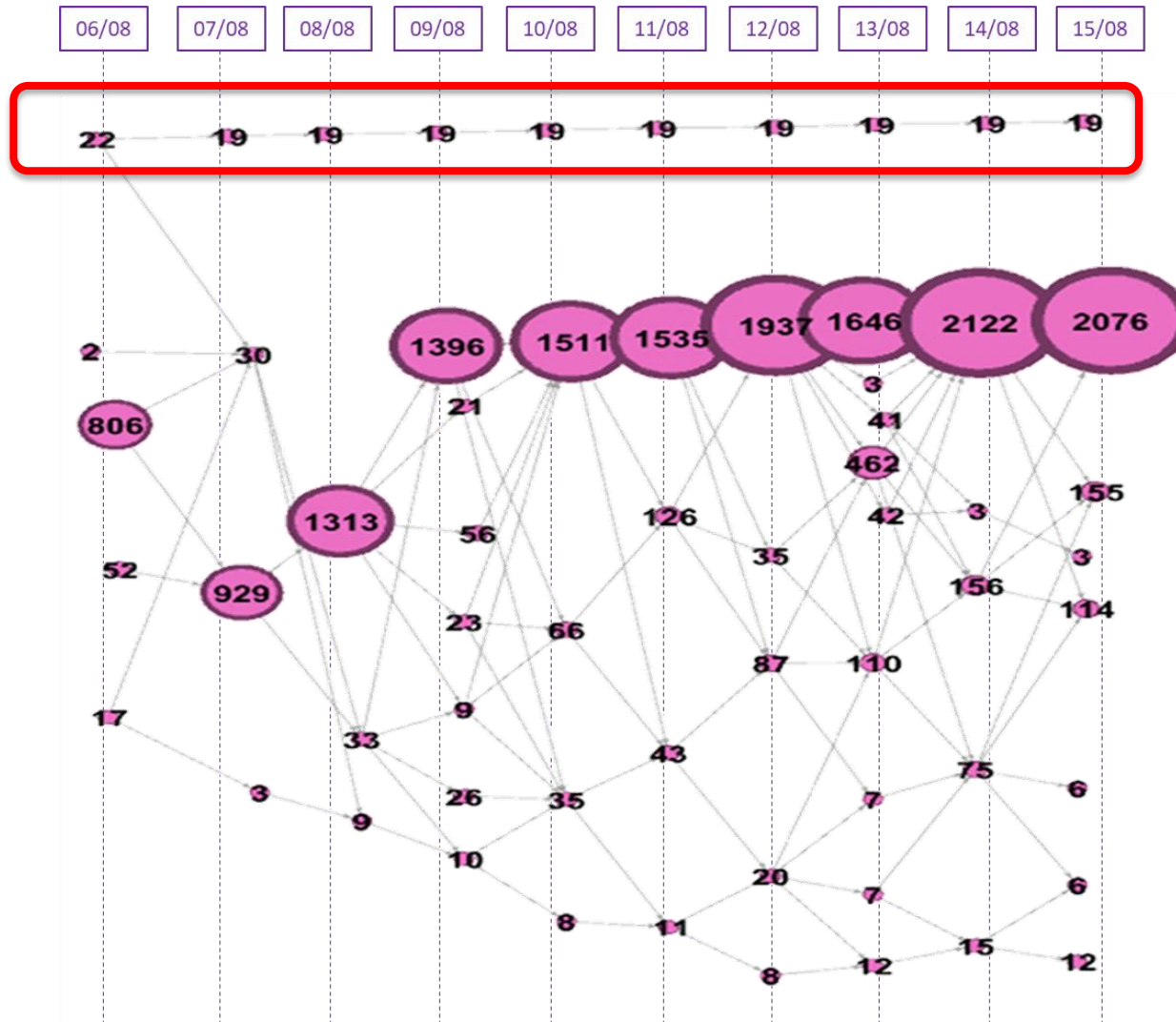


- Décision : 5 clusters demandés pour les utilisateurs

Résultats journaliers – avec érosion



Résultats journaliers – sans érosion

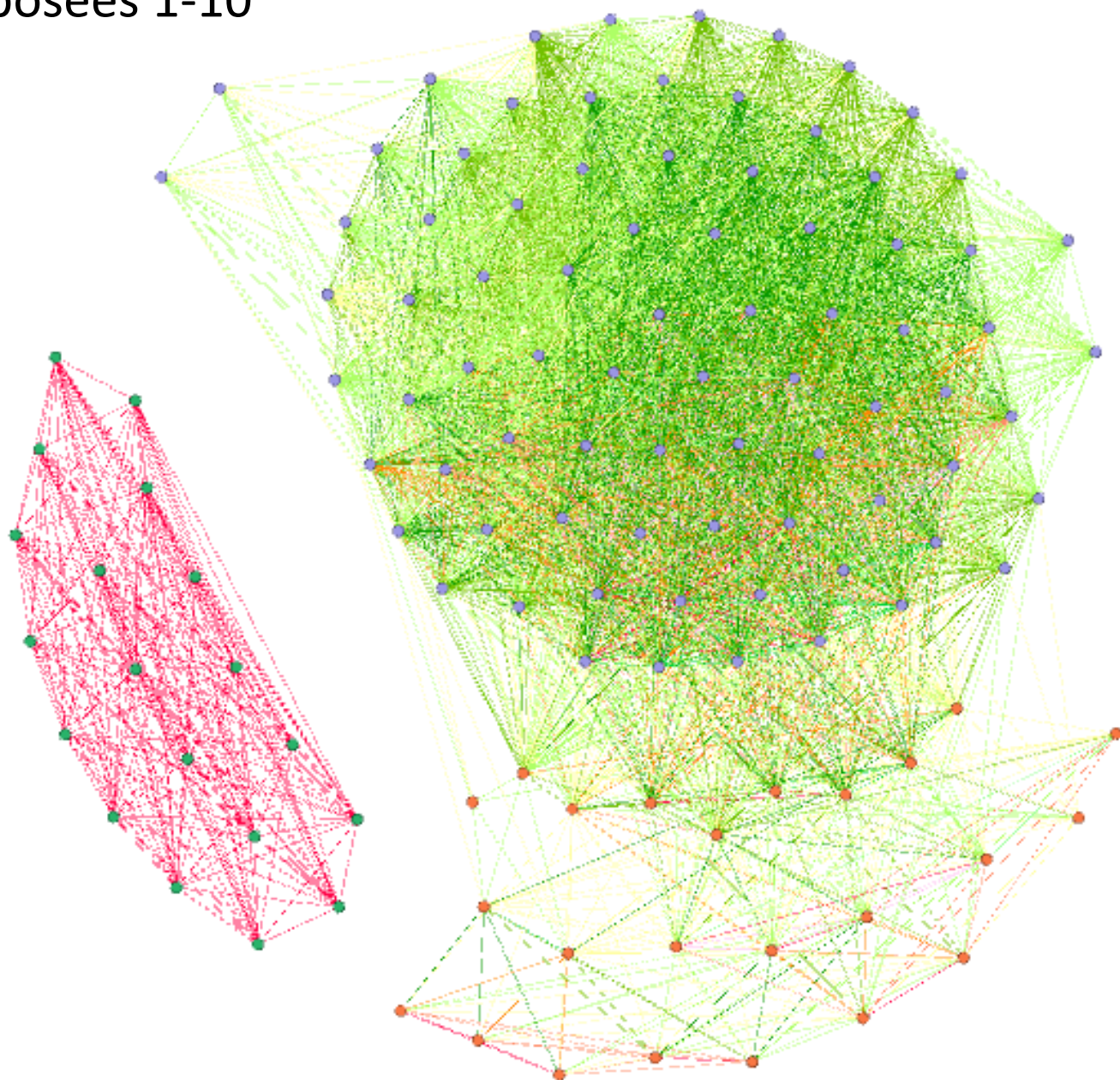


Analyse des trajectoires

- À l'issue du clustering, on affecte le vecteur des attributions de classe journalière aux auteurs.
- On peut alors comparer des trajectoires communes maximales entre auteurs en suivant la chronologie temporelle.
- L'objectif est d'analyser la nature des clusters dans les termes de DPC

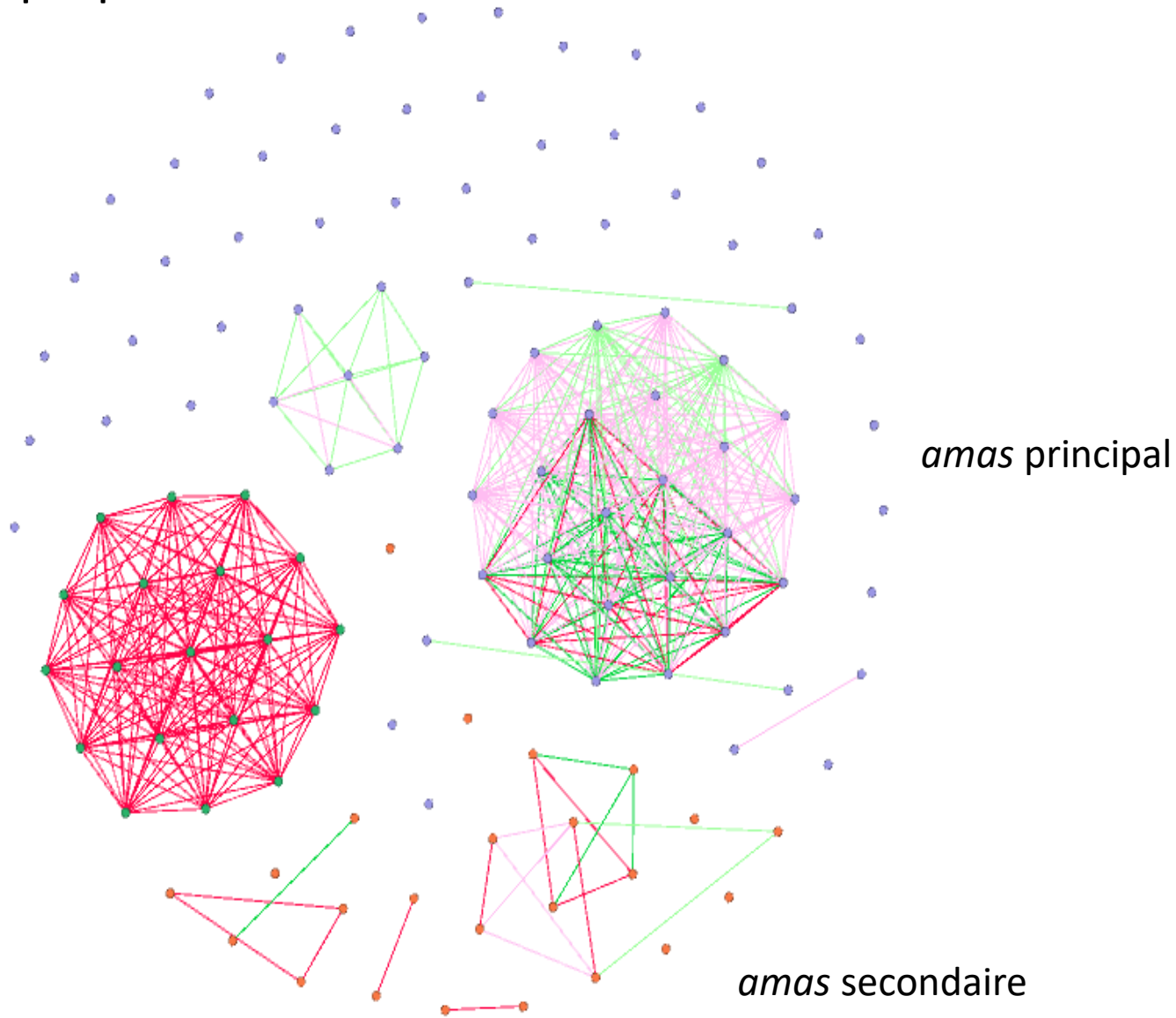
Trajectoires superposées 1-10 gephi

TiredEarth



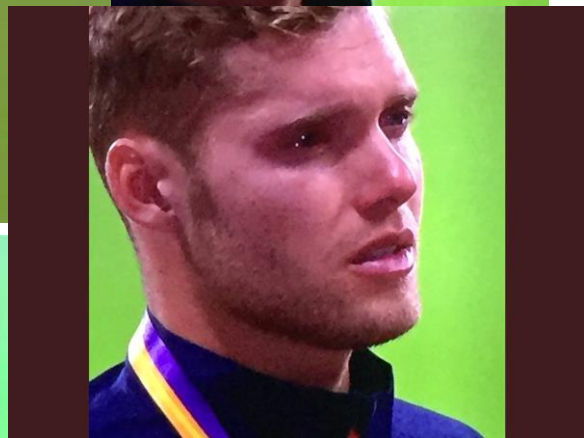
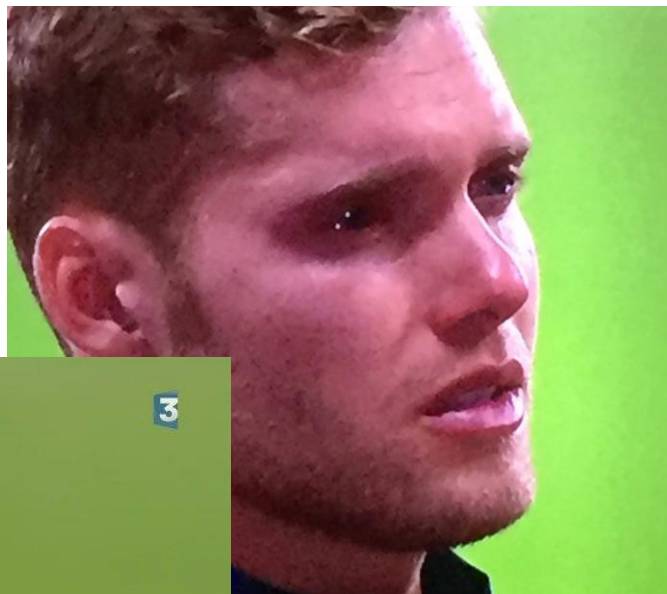
Trajectoires superposées 7-10 gephi

TiredEarth



Conclusion SHS

- En tant que phénomène, l'analyse de DPC est essentielle dans l'interprétation des jeux d'acteurs et des usages sociaux de Twitter dans un contexte événementiel.
- Leur étude contribue à l'analyse de l'action collective et de ses modalités tactiques.
- Elle apporte également un éclairage moins spontané sur la fabrication de l'information de flux.



M. Francony et al.

Conclusion Informatique

- Le flux image introduit une segmentation supplémentaire du flux de publication qu'il est intéressant de prendre en compte dans l'analyse.
- L'unicité de l'image à partir de la clé MD5 doit être relâchée afin d'admettre des altérations de l'image n'affectant pas sa *sémiose* (mais jusqu'où...)
- D'autres clefs telles que les URL (...), pourraient être envisagées de manière complémentaire.