

Séminaire 2018



Lundi 26 mars 2018



RENUROHM - Les données numériques (twitter) : un tournant pour l'analyse des relations hommes-milieux ?

Application au Parc national des Calanques

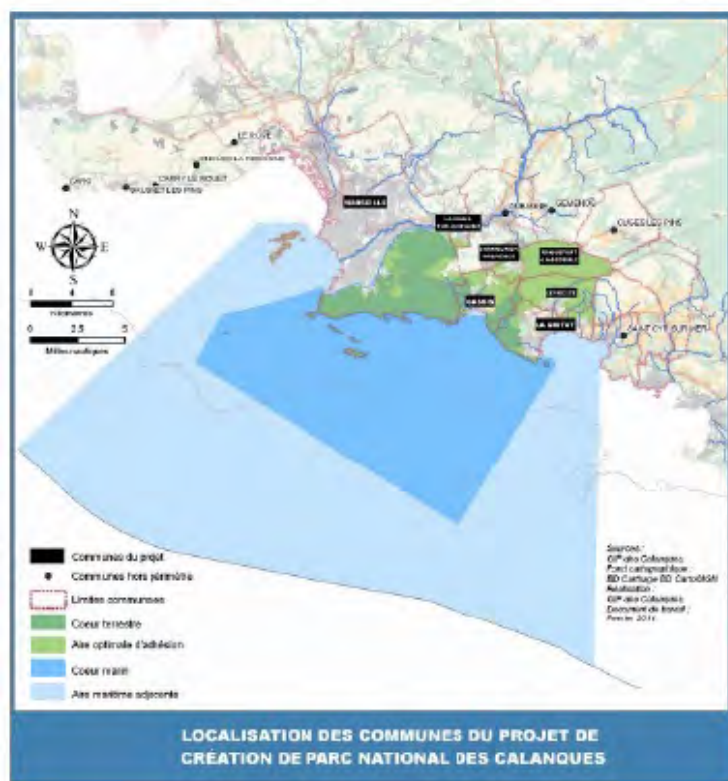
Siqi Fan, Philippe DEBOUDT, Amel FRAISSE, Eric KERGOSIEN

<http://renurohm.univ-lille.fr>



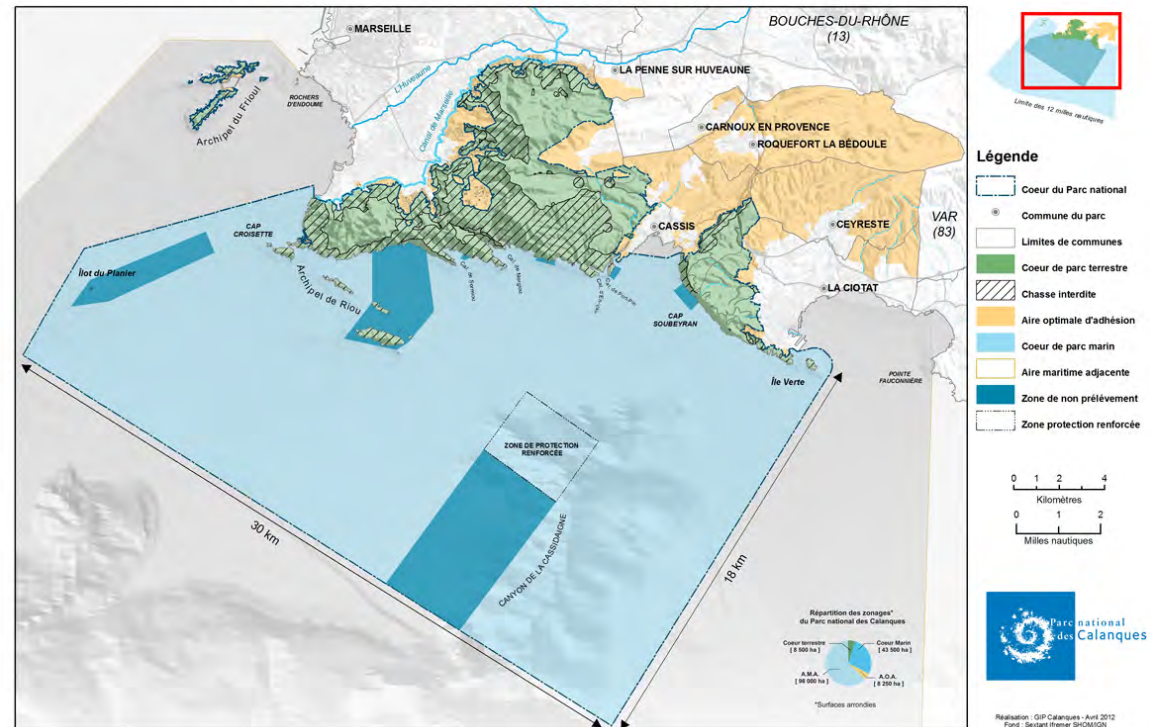
1. LE TERRAIN ET RECHERCHES ANTERIEURES SUR LE TERRITOIRE

2007



2012

Périmètres du Parc national des Calanques



1. LE TERRAIN ET RECHERCHES ANTERIEURES SUR LE TERRITOIRE

De 2008 à 2012 :

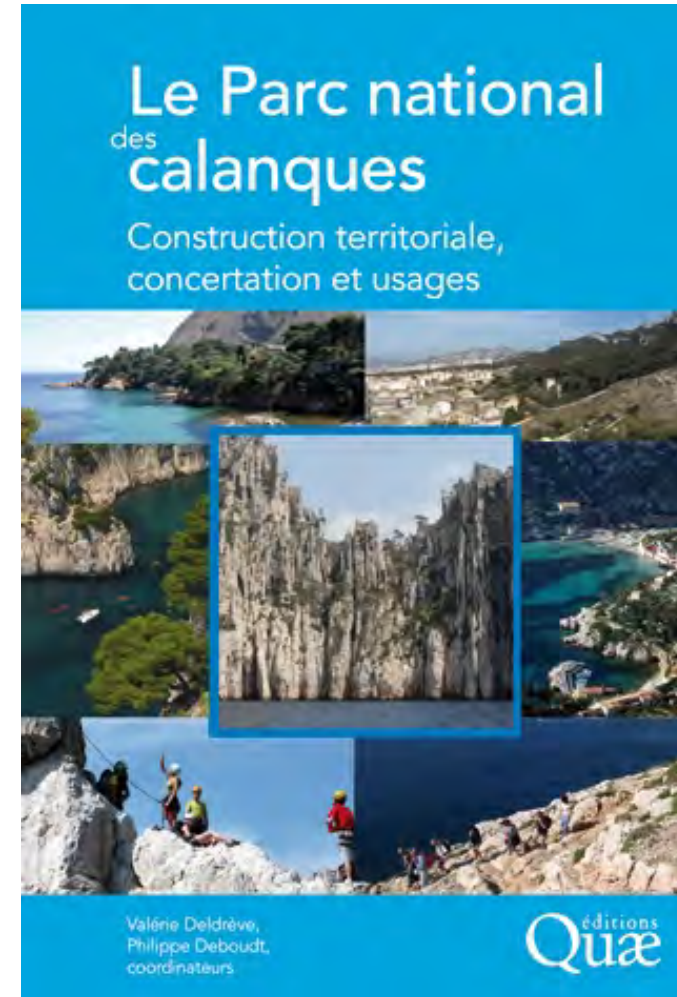
- Concertation mise en œuvre pour réaliser la charte du parc national (Deldrève et Deboudt, 2012).

→ difficulté d'articuler les enjeux globaux et locaux dans la construction d'un tel projet de territoire qui a principalement bénéficié aux usagers traditionnels et s'est accompagnée d'une exclusion des enjeux urbains et maritimes

Volonté :

Volonté de réinterroger ces résultats en proposant une méthodologie mobilisant des corpus de données numériques issues des réseaux sociaux (twitter)

Hypothèse : Internet et réseaux sociaux : prise de parole des « sans voix » (Arsène, 2013)



Programme concertation décision environnement



2. OBJECTIFS ET METHODOLOGIE

Objectifs du projet RENUROHM: <http://renurohm.univ-lille.fr>



- Explorer la faisabilité d'une nouvelle forme d'interdisciplinarité sciences de l'information-communication et géographie de l'environnement
- Démarrer la constitution d'un corpus de données numériques en lien avec l'OHM
- Mettre en place une méthodologie semi-automatique combinant une approche qualitative et une approche quantitative pour l'analyse des communications médiées (twitter)
- Approche reproductible sur d'autres territoires

3. OBJECTIFS ET METHODOLOGIE

- Questions posées :
 - Quelles sont les thématiques et les sujets abordés dans les tweets ?
 - Quels sont les couples d'hashtag les plus utilisés ?
 - Quelles sont les évolutions observées selon différentes temporalités ?
 - Quelles sont les relations entre ces acteurs ? (*Qui parle à Qui ?*)
 - Quels sont les acteurs qui s'expriment sur les lieux/sujets en lien avec le Parc national des Calanques ? (*Qui parle de Quoi ?*)

A terme : Est-ce que les discours, débats, projets contenus dans ces réseaux numériques alimentent, rejoignent ou s'opposent à ceux développés dans l'espace physique ?

3. OBJECTIFS ET METHODOLOGIE

- Problématiques associées :
 - Comment identifier le corpus de messages courts Twitter relatif à la thématique des recherches ;
 - Comment ensuite identifier automatiquement les informations pertinentes à partir de messages courts non standard (faute, pas de vraies phrases dans les tweets, etc.) ;
 - Comment extraire et analyser ces informations : combinaison méthode automatique et analyse experte pour le choix des données, leur traitement et la validation les résultats obtenus

4. DONNEES DE DEPART

Two tweets are shown. The first tweet is by François GARREAU (@GARREAU75) dated July 6th, mentioning #Boues #rouges de #Gardanne and #pollution. The second tweet is by ChristianFaivre (@bleusuper) dated June 28th, mentioning #boues rouges and #justice. A video thumbnail is included in the second tweet, showing a landscape with red mud. Three red boxes on the right are numbered 1, 2, and 3, with lines pointing to the author's name/handle, the hashtags, and the mentioned user in the second tweet respectively.

1. ChristianFaivre : nom de l'auteur ; @bleusuper : nom de l'auteur identifié ;
2. #boues / #justice : hashtags (qui permet de trouver ce tweet plus facilement que ce qui n'a rien marqué)
3. personne(s) indiquée(s) qui a lien avec le contenu du tweet

4. DONNEES DE DEPART : PHASE DE COLLECTE

2 périodes temporelles choisies pour la collecte des données :

- **2007-2011** : processus de création du Parc national des Calanques et de l'organisation du processus de concertation pour élaborer la charte du Parc national.
 - Identification de **groupes d'acteurs favorables ou opposés** à la création d'un **Parc national** dans le territoire des Calanques ;
- **2012-2017** : année de création du Parc national (en 2012) jusqu'aux années récentes
 - Période marquée par **plusieurs conflits d'usages** et notamment celui fortement médiatisé, provoqué par **des rejets de boues rouges issues de la production d'alumine par l'usine ALTEO** de Gardanne, dans les espaces du cœur maritime du Parc national.

4. DONNEES DE DEPART : PHASE DE COLLECTE

Collecte des données : phase qualitative

1. Expertise de géographes spécialistes du territoire d'études : liste d'acteurs ayant participé au processus de création du Parc national des Calanques (associations, conseils de quartiers, personnalités politiques, et citoyens)

2. Identification sur le portail Web des acteurs listés d'une liste de mots clés que nous avons utilisés sous forme de **hashtags** pour collecter les tweets

Hashtags	Description
#ParcNationalCalanques	Parc National des Calanques
#PNCal	Parc National des Calanques
#BouesRouges	Boues Rouges
#Marseille	Marseille
#CreationPNCalanques	Création du Parc National des Calanques
#Alteo	Entreprise
#bouerouge	Boues Rouges
#LittoralSudMarseille	Comité santé
...

Utilisation API Search de Twitter pour collecter et filtrer les messages qui contiennent nos hashtags : 40771 tweets collectés.

5. TRAITEMENT DES DONNEES : METHODOLOGIE

Etape 2 : Approche statistique

Hypothèse : si un mot m est fortement corrélé à un hashtag h de notre liste alors ce mot est un candidat pertinent pour notre analyse de relations territoriales.

- Pour chaque hashtags, extraction de l'ensemble des mots qui lui est associé dans le corpus de tweets :

1. Mesure Information mutuelle introduite par (Fano, 1961) afin de mesurer l'association entre un mot m du corpus et un hashtag h : pour chaque couple de variables aléatoires (m, h) , nous mesurons leur degré de dépendance au sens probabiliste. L'information mutuelle est donnée par la formule suivante :

$$IM(h, m) = \log_2 \left(\frac{freq(h, m)}{freq(h) \cdot freq(m)} \right)$$

Avec

- $freq(h, m)$ est le rapport entre le nombre de tweets contenant le mot m et le hashtag h ($|Th, m|$) et le nombre total de tweets ($|T|$)
- $freq(h)$ est le rapport entre le nombre total de tweets contenant le hashtag h ($|Th|$) et le nombre total de tweets
- $freq(m)$ est le rapport entre le nombre total de tweets contenant le hashtag h ($|Tm|$) et le nombre total de tweets

5. TRAITEMENT DES DONNEES : METHODOLOGIE

Etape 3 : Visualisation et analyses

Plateforme R utilisée pour ces traitements (version R3.4 du logiciel Rstudio) :

- langage R pour réaliser la fouille de textes
- Ggplot2 les visualisations des résultats.
- Mots ordonnés par ordre croissant selon le degré d'association : du plus pertinent au moins pertinent.
- Suppression des tweets intégrant #marseille et calanques : 12 hashtags sélectionnés après validation experts
- Filtrage des couples pertinents selon un seuil défini avec les experts : > 2 occurrences du couple dans le corpus traité.

	B	C	D	E	F
	MotRequest	hashtag	IM	Ocurrence	
1	ParcNationalCalanques_	cg13	9.75444836289823		1
2	ParcNationalCalanques_	contratdebaie	9.75444836289823		1
3	ParcNationalCalanques_	enqu	9.75444836289823		1
4	ParcNationalCalanques_	eurom	9.75444836289823		
5	ParcNationalCalanques_	evidence	9.75444836289823		
6	ParcNationalCalanques_	frioul	9.75444836289823		
7	ParcNationalCalanques_	guyteissier	9.75444836289823		
8	ParcNationalCalanques_	marseillegrandest	9.75444836289823		
9	ParcNationalCalanques_	nihous	9.75444836289823		
0	ParcNationalCalanques_	parcnationalcalanc	9.75444836289823		
1	ParcNationalCalanques_	promessesquinengag	9.75444836289823		
2	ParcNationalCalanques_	schpountz	9.75444836289823		
3	ParcNationalCalanques_	sormiou	9.75444836289823		
4	ParcNationalCalanques_	teissier	9.75444836289823		
5	ParcNationalCalanques_	ump	9.75444836289823		
6	ParcNationalCalanques_	unfollowsinonrien	9.75444836289823		
7	PNCAL_	eaart	7.88844566828897		
8	PNCAL_	jamansd	7.88844566828897		
9	PNCAL_	logia	7.88844566828897		

	B	C	D
	MotRequest	hashtag	IM
2	ALTEO_	jobs	0.307293754210503
1	ALTEO_	montreal	0.307293754210502
1	ALTEO_	job	0.307293754210503
9	ALTEO_	canada	0.307293754210503
2	ALTEO_	google	0.307293754210502
0	ParcNationalCalanques_	parcnationalcalanc	9.75444836289823
0	Parc National des Calanques	marseille	5.61052274754212
0	Boues rouges_	hongrie	4.71421241370507
2	ALTEO_	serp	0.307293754210502
9	ALTEO_	seo	0.307293754210503
3	ALTEO_	emploi	0.307293754210503
9	ALTEO_	smm	0.307293754210503
1	ALTEO_	career	0.307293754210503
3	ALTEO_	fresher	0.307293754210503
9	ALTEO_	ecommerce	0.307293754210502
4	ALTEO_	facebook	0.307293754210503

5. TRAITEMENT DES DONNEES : RESULTATS LIES AUX THEMES

Etape 3 : Visualisation et analyses des thématiques

Nuage de hashtags – mots clés thématiques pré-sélectionnés par les experts et mondes lexicaux associés avant 2012

- le cas des couples pour le thème « environnement »

	Avant 2012	Après 2012
Tweets	2064	16683
Couples	345	3133
ALTEO	342	1966
(ALTEO, bouesrouges)	0	382
(ALTEO, pollution)	0	117
(ALTEO, environnement)	0	63
Bouesrouges	111	3180
(Bouesrouges, bouesrouges)	4	484
(Bouesrouges, pollution)	4	101
(Bouesrouges, environment)	4	100
Parc National des Calanques	175	1683
(PNCAL, bouesrouges)	0	67
(PNCAL, pollution)	0	20
(PNCAL, environment)	4	27

Création du parc
Pollution boues r

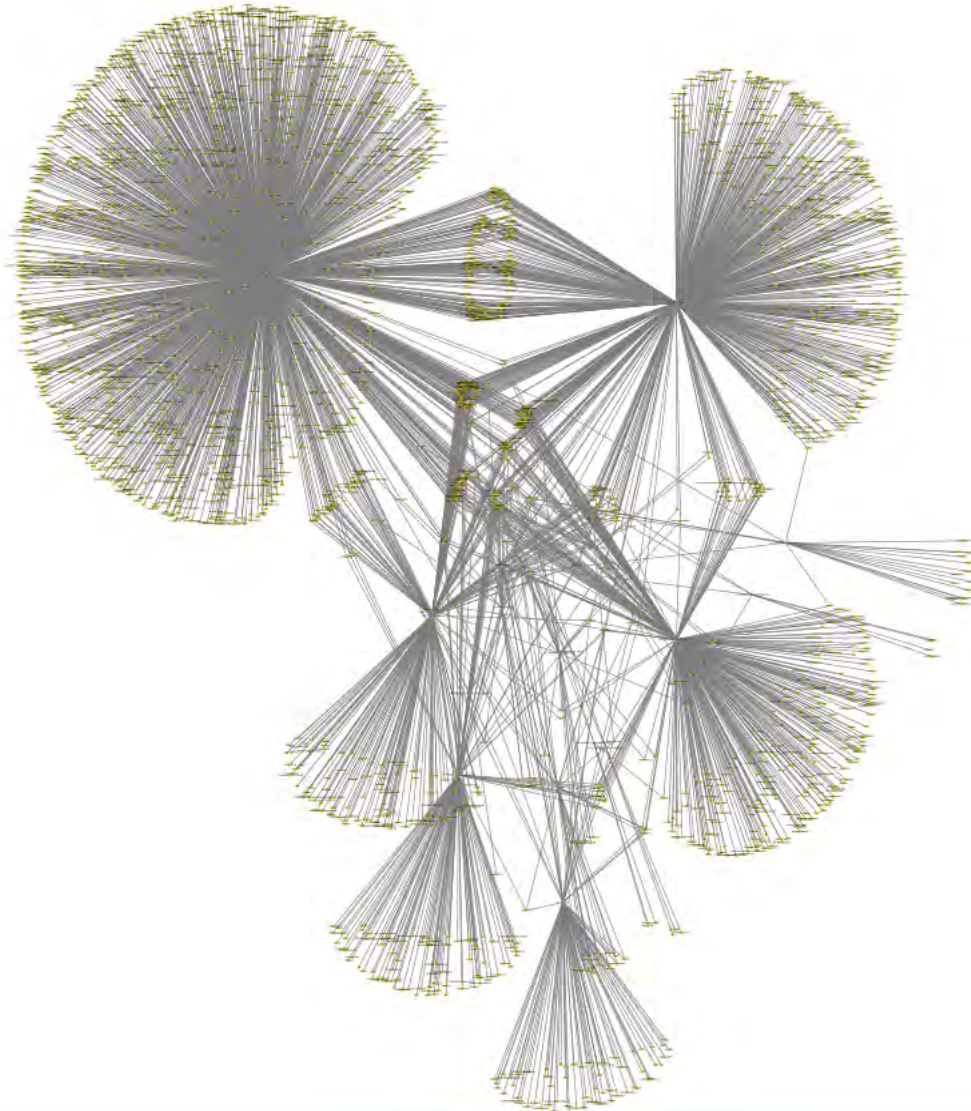
→ Relance de la collecte avec les nouveaux HashTags identifiés

5. TRAITEMENT DES DONNEES : RESULTATS LIES AUX THEMES

Etape 3 : Visualisation et analyses des thématiques

Nuage de hashtags – mots clés thématiques pré-sélectionnés par les experts et mondes lexicaux associés à partir de 2012

- le cas des couples pour le thème « **environnement** »



5. TRAITEMENT DES DONNEES : RESULTATS LIES AUX THEMES

Etape 3 : Visualisation et analyses des thématiques

Nuage de hashtags – mots clés thématiques pré-sélectionnés par les experts

- le cas des couples pour le thème « environnement »

- Evolution des thématiques dans le temps

de 2007 à 2011

Hashtags Experts	Hashtags
Alteo	#jobs, #emploi, #career, etc.
Parc National Calanques	#parcnationalcalanques, #marseille, #nature, #environnement
Boues rouges	#hongrie, #pollution

de 2012 à 2017

Hashtags Experts	Hashtags
Alteo	#bouesrouges, #emploi, #pollution, #calanques, #marseille, #environnement
Parc National Calanques	#calanques, #marseille, #environnement, #ArchipelDeRiou, #parcnationaldescalanques, #bouesrouges, #nature
Boues rouges	#bouesrouges, #calanques, #marseille, #alteo, #pollution

5. TRAITEMENT DES DONNEES : RESULTATS LIES AUX ACTEURS

Qui parle sur qui ? Et sur Quoi?

- Avant 2012

B	C	D	E	F
UserID	UserName	Mention	Nombre	
80310555	un_job_de_suite	algebrik	37	
80310555	un_job_de_suite	algebrik2010emploi	17	
19266151	Alexandre Penyauski	agence	8	
19266151	Alexandre Penyauski	alteo	8	
99609168	Alteo	jdnebusiness	6	
1,8E+08	Gengembre Dominique	scoopit	5	
2,74E+08	7AY4	idipodong	5	
4471191	François Cazals	agence	4	
4471191	François Cazals	alteo	4	
4471191	François Cazals	bitchlemagazine	4	
4471191	François Cazals	intempestif	4	
17406148	l_loirs	agence		
17406148	l_loirs	alteo		
21040181	Allan Teo	ashleygagateo		
47288166	Andr?<f0><U+009F><U+	agence		
47288166	Andr?<f0><U+009F><U+	alteo		
99609168	Alteo	gregfromparis		
1,14E+08	Louis Dollo	addthis		

- À partir de 2012

UserID	UserName	mention
446939787	<U+B808><U+B385><U+C740> <U+BE0C><U+C6>	bot
446939787	<U+B808><U+B385><U+C740> <U+BE0C><U+C6>	alteo
1108658778	l'Encre de Mer	encredemer
481170516	khaldi	royalsegolene
384206154	Didier R	parccalanques
481170516	khaldi	nichelerivasi
874404774	Vialidad Bs. As.	mininfra
874404774	Vialidad Bs. As.	baprovincia
874404774	Vialidad Bs. As.	vialidadba
2936740168	Anarchozy	manuelvalls
1945218524	InsoumisdeBerre13	cnpcf
212778422	SiDD VLAK <f0><U+009F><U+0092><U+0080>	datpiff
874404774	Vialidad Bs. As.	solotransito
1556975869	El interior existe!	rutarota
384206154	Didier R	prefet13
436922925	Dani	royalsegolene
436922925	Dani	didierreault
932234984	PCF 13	cnpcf
7.912099e+17	@antaresdul3	youtube
384206154	Didier R	parccalanque
436922925	Dani	marsactu
104883773	Marsactu	marsactu
267938974	bernard jegou	lemarin
287557798	Andre Dechene	royalsegolene
436922925	Dani	nichelerivasi

- Nécessité d'étudier en détails

les liens entre acteurs :

- Analyse fine des échanges

→ Identification des communautés d'acteurs

5. TRAITEMENT DES DONNEES : RESULTATS LIES AUX ACTEURS

Qui parle sur qui ? Et sur Quoi?

- Avant 2012

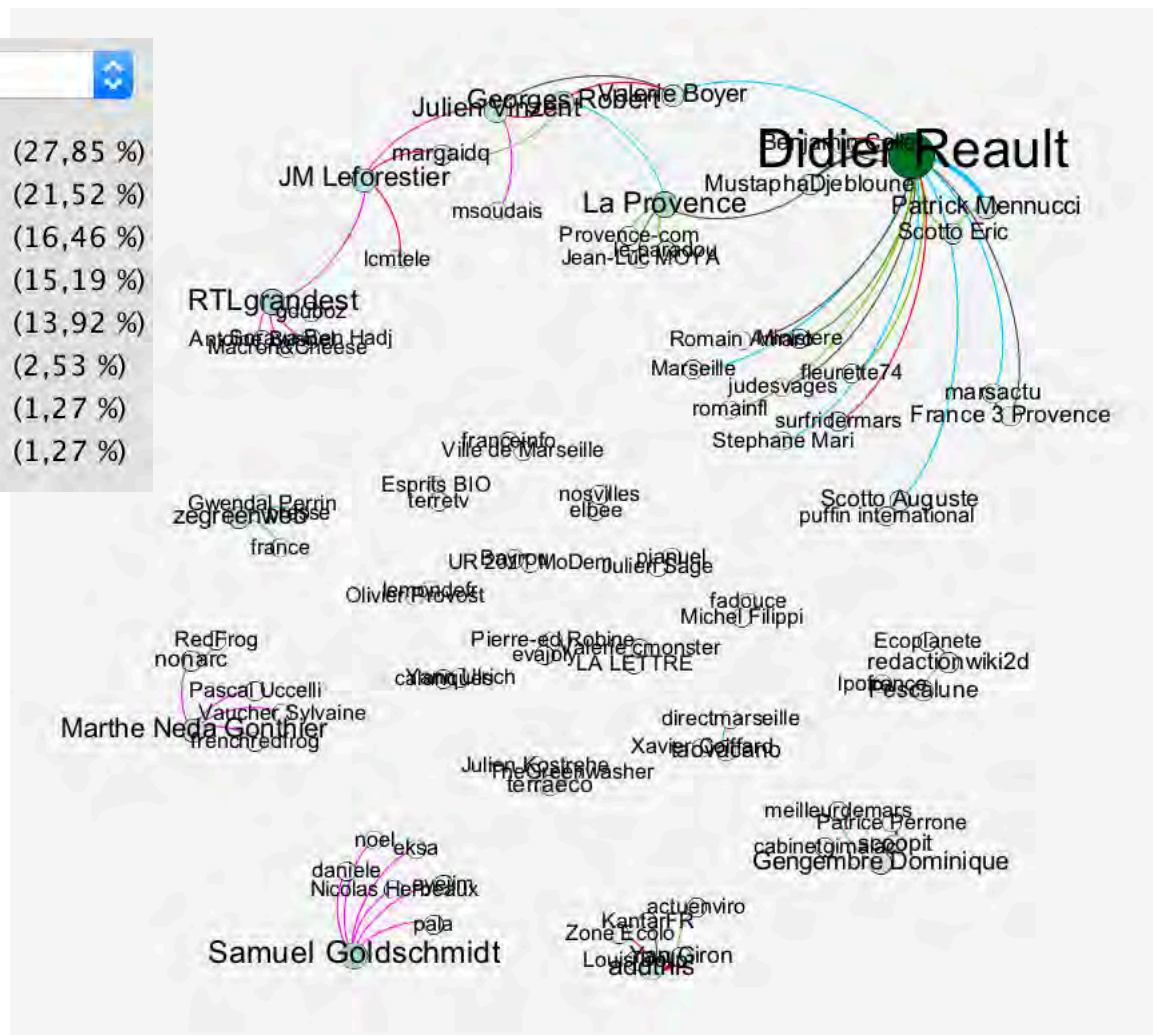
B	C	D	E	F
UserID	UserName	Mention	Nombre	
80310555	un_job_de_suite	algebrik	37	
80310555	un_job_de_suite			
19266151	Alexandre Penyauski			
19266151	Alexandre Penyauski			
99609168	Alteo			
1,8E+08	Gengembre Dominique			
2,74E+08	7AY4			
4471191	François Cazals			
4471191	François Cazals			
4471191	François Cazals			
4471191	François Cazals			
17406148	l_loirs			
17406148	l_loirs			
21040181	Allan Teo			
47288166	Andr?			
47288166	Andr?			
99609168	Alteo			
1,14E+08	Louis Dollo			
1,01E+08	XXXXXXXXXX			

A	B	C	D	E	F	G
1	source	target	Theme	TEXTE TWEET		
2	Gengembre Dominique	scoopit	Projet Amenagement	parc national des calanques @scoopit http://bit.ly/IO4		
3	Gengembre Dominique	scoopit	Projet Amenagement	Parc national des Calanques : un projet en attente d'am		
4	Gengembre Dominique	scoopit	Parc National	La crZation du parc national des Calanques repoussZe au		
5	Gengembre Dominique	scoopit	Parc National	Parc National des Calanques : le Club Alpin Marseille Pro		
6	Gengembre Dominique	scoopit	Projet Amenagement	Un pas de plus vers le Parc national des Calanques Ä @		
7	Louis Dollo	addthis	Enquete	KAIRN Parc National des Calanques : Žcrivez aux enqu		
8	Louis Dollo	addthis	Parc National - protestation	Le Parc National des Calanques : Moi je suis contre! : htt		
9	Louis Dollo	addthis	Parc National - protestation	Ils disent non au Parc National des Calanques http://ww		
10	Yan Giron	addthis	Parc National - protestation	P cheurs et plaisanciers protestent en mer contre le pa		
11	Yan Giron	addthis	Parc National - protestation	Marseille : contre la parc national des calanques - ENVIR		
12	Yan Giron	addthis	Parc National - protestation	P cheurs et plaisanciers protestent en mer contre le pa		
13	Didier Reault	Patrick Mennucci	Parc National - pour	Guy #Teissier est bien sur a fond sur l'ordi et la crZation		
14	Didier Reault	Patrick Mennucci	Parc National - pour	@patrickmennucci š quelques encablures du #ParcNatio		
15	Didier Reault	Patrick Mennucci	Parc National - pour	Merci du soutien ! Les reco de l'EP seront ŽtudiZes. Le cc		
16	Patrice Perrone	scoopit	Parc National - pour	MARSEILLE POUR UN PARC NATIONAL DES CALANQUES		
17	Patrice Perrone	scoopit	Enquete	Parc national des Calanques : l'enqu te publique est ou		
18	MustaphaDjebloune	Didier Reault	Parc National - protestation	@laprovence je suis š l'UMP mais je ne cautionne pas les		
19	MustaphaDjebloune	Didier Reault	Enquete	- @DidierReault: Parc national des Calanques : l'enqu		
20	puffin international	Scotto Auguste	Parc National - pour	@scottoauguste c'est a un sondage ?! Et je peux voter		
21	puffin international	Scotto Auguste	Parc National - pour	@scottoauguste je viens de voter pour le #ParcNational		
22	Xavier Coiffard	taovacano	Pollution	@taovacano La revelete c'est des aiguilles un peu partou		
23	Macron&Cheese	RTLgrandest	Pollution	Les boues rouges s'Zchappent de la digue.		
24	taovacano	directmarseille	Parc National - pour	Enfin une bonne nouvelle !: RT @DirectMarseille: "Le pa		
25	Patrick Mennucci	Didier Reault	Election	RŽunion du groupe Faire Gagner Marseille a la Mairie du		
26	gduboz	RTLgrandest	Pollution	RIP le pantalon de @rtlgrandest victime des boues rouge		
27	Soraya Ben Hadj	RTLgrandest	Pollution	Trois mois apr s, les victimes des boues rouges toujour		
28	Julien Vinzent	msoudais	Pollution	@msoudais ne serait pas surpris que personne ne reprei		
29	Julien Vinzent	Valerie Boyer	Enquete	Allez surtout ^ l'enqu te publique RT @valerieboyer13		
30	KantarFR	addthis	Enquete	Bonjour ! Focus sur l'opinion des riverains sur le projet d		
31	TheGreenwasher	terraeco	Pollution	Revue de presse du week-end !: sur @terraeco > des bou		
32	RedFrog	nonarc	Pollution	RT @nonarc: CoulZes de boues rouges - Appel urgent de		
33	Marthe Neda Gonthier	nonarc	Pollution	V @pascaluccelli @VaucherSylvaine @frenchredfrog @r		
34	Marthe Neda Gonthier	frenchredfrog	Pollution	V @pascaluccelli @VaucherSylvaine @frenchredfrog @r		

- Nettoyage
- Analyse qualitative des experts pour le filtrage des tweets pertinents

5. TRAITEMENT DES DONNEES : RESULTATS LIES AUX ACTEURS & THEMES

- Qui parle sur qui et sur quoi (2007 – 2011)

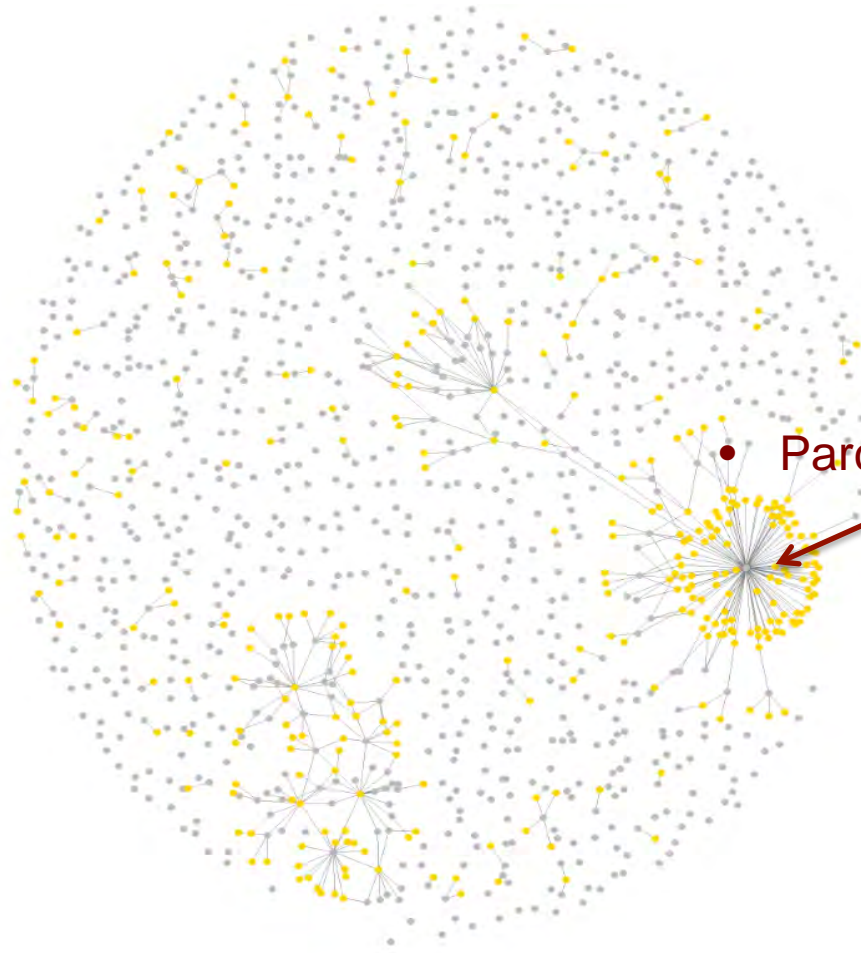


- Analyses :
 - Présence forte des politiques et des médias (journalistes)
 - Acteurs participants au processus de construction du parc absents!

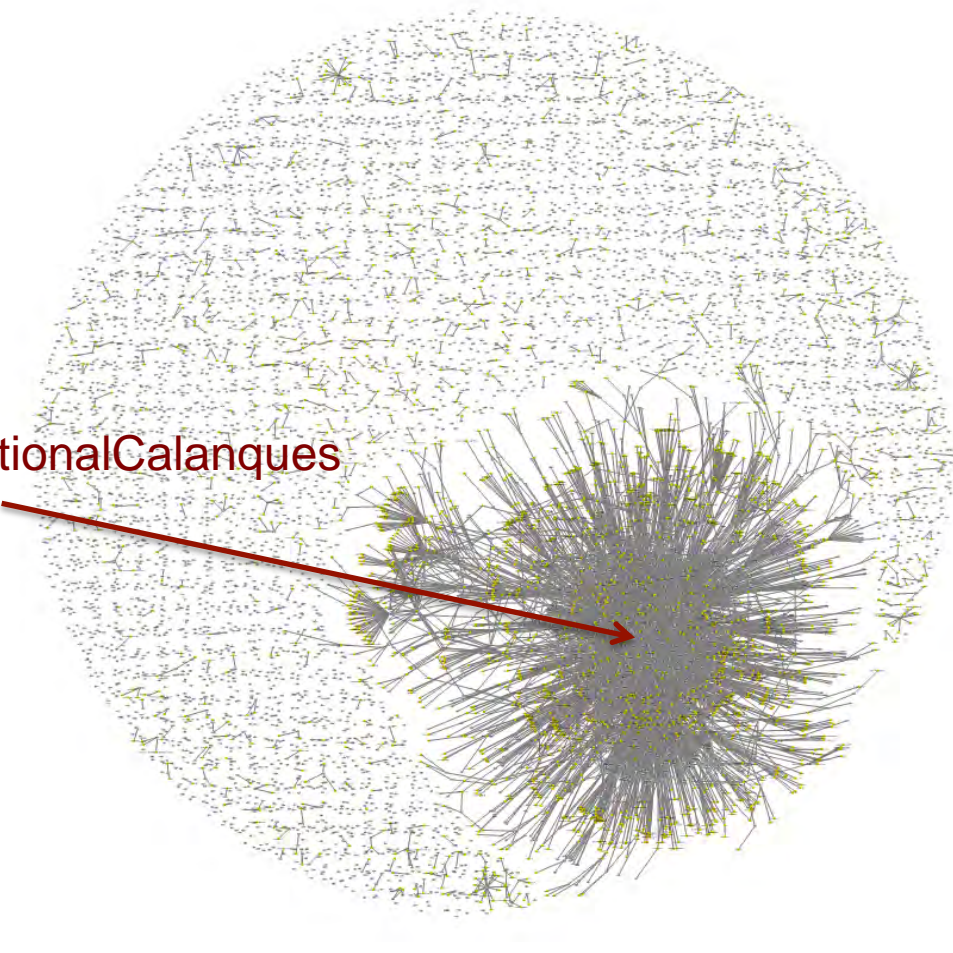
5. TRAITEMENT DES DONNEES : RESULTATS LIES AUX ACTEURS & THEMES

Qui parle sur qui ? Et sur Quoi?

- Avant 2012



- À partir de 2012



- ParcNationalCalanques

- 20% du corpus traité sur la période 2012 – 2017
- Pistes : augmentation de la mention liée à la pollution (Boues rouges et Alteo).
- Analyse à suivre sur les acteurs qui s'expriment depuis 2012 (sans-voix? politiques?)

6. PREMIERES CONCLUSIONS

- Une **approche quanti – quali intéressante** pour le traitement et l'analyse des données provenant des réseaux sociaux
- Un usage des réseaux sociaux de plus en plus important sur la thématique environnement en lien avec le parc national des calanques
 - Des thématiques qui ressortent : #pollution, #nature, #bouesrouges...
 - 2007 – 2011 :
 - Des acteurs politiques et journalistes qui s'expriment avant 2012. Et les sans-voix?
 - Des thématiques autour de la construction du parc, de la propagande pour le parc, de la pollution contrairement à ce qui était dit par les politiques,
 - Une évolution des thématiques dans le temps : forte présence de l'acteur industriel Alteo dans le contenu des tweets

6. PERSPECTIVES

- Finaliser les analyses des données sur la période 2012 – 2017
- Approches de fouille de textes pour préciser :
 - Le marquage des entités spatiales à partir de messages courts (Zenasni *et al.*, 2016 ; 2017)

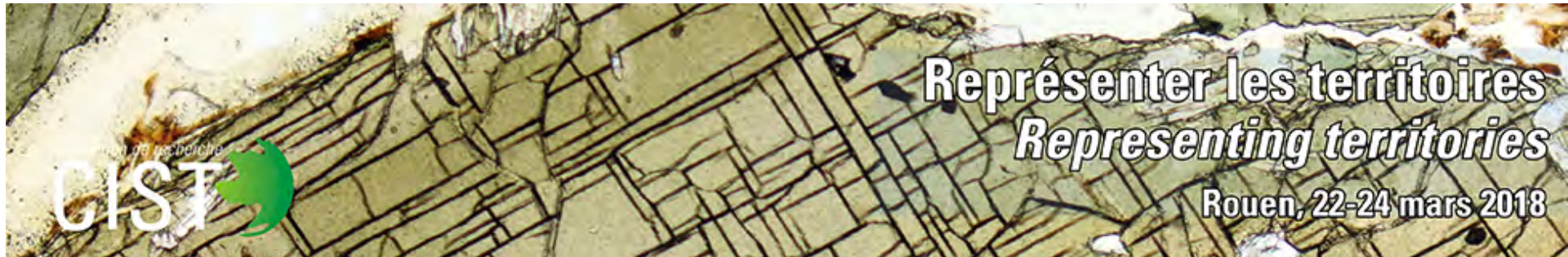


Références bibliographiques :

- Arsène S., 2013, *Vers une recomposition des pouvoirs : Internet et réseaux sociaux*, CERISCOPE Puissance
- Berthelot M.-A., Severo M., Kergosien E., 2016, , Cartographier les acteurs d' un territoire : une pproche appliquée au patrimoine industriel textile du Nord-Pas-de-Calais, In 3ème colloque international du CIST (CIST 2016), pp.6, Grenoble.
- Deboudt P., 2015, L'aménagement du littoral à l'épreuve des inégalités environnementales, *Annales des mines, Responsabilité et Environnement*, n°79, p.83-89.
- Deboudt P., Deldrève V., 2015, Inégalités et concertation « encadrée » : le projet du parc national des calanques, *in* L. Mermet et D. Salles (dir.), *Environnement et transition écologique*, De Boeck éd., coll. Ouvertures Sociologiques, p. 151-166.
- Deldrève V., Deboudt Ph. (dir.), 2012, *Le parc national des calanques : construction territoriale, concertation et usages*, QUAE, 231 p.
- Deboudt Ph. (éd.), 2010, *Inégalités écologiques, territoires littoraux, développement durable*, Presses Universitaires du Septentrion, 409 p.
- Fraisse A., Paroubek P., 2014, *Twitter as a Comparable Corpus to build Multilingual Affective Lexicons*, *in* proceedings of the 7th International Workshop on Building and Using Comparable Corpora at LREC 2014 (BUCC 2014), pages 17-21. May 26-31, 2014., Reykjavik, Iceland.
- Kergosien E., B. Laval, Roche M., Teisseire M., Are opinions expressed in land-use planning documents? In *International Journal of Geographical Information Science*, vol. 28, issue 4, jan. 2014. [Rank A, IF: 1.61 in 2012)
- Pak A. , Paroubek P., Fraisse F., Francopoulo G., 2014, *Normalization of Term Weighting Scheme for Sentiment Analysis. Book Chapter, Human Language technology Challenges for Computer Science and Linguistics. Series: Lecture Notes in Artificial Intelligence*, Springer, Vol. 8387. ISBN 978-3-319-08957-7. Vetulani, Zygmunt, Mariani, Joseph (Eds.). May 27, 2014.
- Zenasni S., Kergosien E., Roche M., Teisseire M., 2016, Extracting new Spatial Entities and Relations from Short Messages, In the 8th International ACM Conference on Management of Digita EcoSystems (MEDES'2015), pp. 8, Hendaye (France).

6. Références bibliographiques : <http://renurohm.univ-lille.fr>

Contacts : philippe.deboudt@univ-lille.fr, eric.kergosien@univ-lille.fr



DES DÉCOUVERTES

*Faites de rencontres
& d'émotions*

DES CONNAISSANCES

*Au service de la nature
& des hommes*

DES ACTIONS

*Au côté des acteurs
du Parc national*

