



# Détection automatique d'évènements à partir de Twitter

Houssem Eddine DRIDI

*Soutenance de thèse, 2 octobre 2014*



**Housseem dridi**

Voir ma page de profil

33

TWEETS

164

ABONNEMENTS

20

ABONNÉS

Écrire un nouveau Tweet...

Suggestions · Actualiser · Tout afficher



**Mohamed Ali Chebaane** @MedAl... ✕

Suivi par Eya Turki et d'autres

Suivre



**Yassine Brahim** @Yassine\_Brahim ✕

Suivre

Parcourir les catégories · Trouver vos amis

Tendances · Modifier

#galaartis

#PCPointer

#ThingsNotToDoOnAFirstDate

#bcdebate2013

#cdnpoli

Oakville

CBC

Bernabeu

Madrid

## Tweets



**Oueld Ifriqiya** @Khamousss

1 min

Le FLN, un État dans l'État. #Algérie

Ouvrir



**Tunisia News** @tunisnews

1 min

#Tunisie #Tunisia Tunisie - Politique : Soumaya Ghannouchi s'attaque à la France [bit.ly/14Q4tKu](http://bit.ly/14Q4tKu)

Ouvrir



**Tunisia News** @tunisnews

1 min

Urgent : Report de la grève des transporteurs de carburants : Tunisie : Actualités : Société : ... [goo.gl/8wS0F](http://goo.gl/8wS0F) #Tunisie #Tunisia

Ouvrir



**SAT-7 Arabic** @sat7ar

1 min

[youtu.be/qEzWEEUC07w?a](http://youtu.be/qEzWEEUC07w?a) via @YouTube : برنامج كرمالك عراق - القداصة

Afficher le média



**koora** @KooraTunisie

2 min

Mi-temps : RM-DOR 0-0 (1-4)

Ouvrir



**radioexpressfm** @radioexpressfm

3 min

L'interview de Hedi Ben Abbes dans Tac au Tac [fb.me/IHx5mvzb](http://fb.me/IHx5mvzb)

Ouvrir



@khawlamagrebi الاخوان عدو الاتمان

4 min

علماء: فيروس التهاب الكبد الفيروسي عمره أكثر من 82 مليون سنة [fb.me/PRHluygX](http://fb.me/PRHluygX)

#أوووووووووف

Ouvrir



**Housseem dridi**

Voir ma page de profil

33

TWEETS

164

ABONNEMENTS

20

ABONNÉS

Écrire un nouveau Tweet...

Suggestions · Actualiser · Tout afficher



**Mohamed Ali Chebaane** @MedAl... ✕

Suivi par Eya Turki et d'autres

Suivre



**Yassine Brahim** @Yassine\_Brahim ✕

Suivre

Parcourir les catégories · Trouver vos amis

Tendances · Modifier

#galaartis

#PCPointer

#ThingsNotToDoOnAFirstDate

#bcdebate2013

#cdnpoli

Oakville

CBC

Bernabeu

Madrid

## Tweets



**Oueld Ifriqiya** @Khamousss

1 min

Le FLN, un État dans l'État. #Algérie

Ouvrir



**Tunisia News** @tunisnews

1 min

#Tunisie #Tunisia Tunisie - Politique : Soumaya Ghannouchi s'attaque à la France [bit.ly/14Q4tKu](http://bit.ly/14Q4tKu)

Ouvrir



**Tunisia News** @tunisnews

1 min

Urgent : Report de la grève des transporteurs de carburants : Tunisie : Actualités : Société : ... [goo.gl/8wS0F](http://goo.gl/8wS0F) #Tunisie #Tunisia

Ouvrir



**SAT-7 Arabic** @sat7ar

1 min

[youtu.be/qEzWEEUC07w?a](http://youtu.be/qEzWEEUC07w?a) via @YouTube برنامج كرمالك عراق - القدس

Afficher le média



**koora** @KooraTunisie

2 min

Mi-temps : RM-DOR 0-0 (1-4)

Ouvrir



**radioexpressfm** @radioexpressfm

3 min

L'interview de Hedi Ben Abbes dans Tac au Tac [fb.me/IHx5mvzb](http://fb.me/IHx5mvzb)

Ouvrir



@khaw/amagreb1 الاخوان عدو الاتمان

4 min

علماء: فيروس التهاب الكبد الفيروسي عمره أكثر من 82 مليون سنة [fb.me/PRHluygX](http://fb.me/PRHluygX)

#أوووووووف

Ouvrir



**Housseem dridi**

Voir ma page de profil

33

TWEETS

164

ABONNEMENTS

20

ABONNÉS

Écrire un nouveau Tweet...

Suggestions · Actualiser · Tout afficher



**Mohamed Ali Chebaane** @MedAl... ✕

Suivi par Eya Turki et d'autres

Suivre



**Yassine Brahim** @Yassine\_Brahim ✕

Suivre

Parcourir les catégories · Trouver vos amis

Tendances · Modifier

#galaartis

#PCPointer

#ThingsNotToDoOnAFirstDate

#bcdebate2013

#cdnpoli

Oakville

CBC

Bernabeu

Madrid

## Tweets



**Oueld Ifriqiya** @Khamousss

1 min

Le FLN, un État dans l'État. #Algérie

Ouvrir



**Tunisia News** @tunisnews

1 min

#Tunisie #Tunisia Tunisie - Politique : Soumaya Ghannouchi s'attaque à la France [bit.ly/14Q4tKu](http://bit.ly/14Q4tKu)

Ouvrir



**Tunisia News** @tunisnews

1 min

Urgent : Report de la grève des transporteurs de carburants : Tunisie : Actualités : Société : ... [goo.gl/8wS0F](http://goo.gl/8wS0F) #Tunisie #Tunisia

Ouvrir



**SAT-7 Arabic** @sat7ar

1 min

[youtu.be/qEzWEEUC07w?a](http://youtu.be/qEzWEEUC07w?a) via @YouTube : برنامج كرمالك عراق - القدس

Afficher le média



**koora** @KooraTunisie

2 min

Mi-temps : RM-DOR 0-0 (1-4)

Ouvrir



**radioexpressfm** @radioexpressfm

3 min

L'interview de Hedi Ben Abbes dans Tac au Tac [fb.me/IHx5mvzb](http://fb.me/IHx5mvzb)

Ouvrir



@khaw/amagreb1 الاخوان عدو الاتمان

4 min

علماء: فيروس التهاب الكبد الفيروسي عمره أكثر من 82 مليون سنة [fb.me/PRHluygX](http://fb.me/PRHluygX)

#أوووووووووف

Ouvrir



Housseem dridi

Voir ma page de profil

33

TWEETS

164

ABONNEMENTS

20

ABONNÉS

Écrire un nouveau Tweet...

Suggestions · Actualiser · Tout afficher



Mohamed Ali Chebaane @MedAl... x

Suivi par Eya Turki et d'autres

Suivre



Yassine Brahim @Yassine\_Brahim x

Suivre

Parcourir les catégories · Trouver vos amis

Tendances · Modifier

#galaartis

#PCPointer

#ThingsNotToDoOnAFirstDate

#bcdebate2013

#cdnpoli

Oakville

CBC

Bernabeu

Madrid

### Tweets



Oueld Ifriqiya @Khamousss

1 min

Le FLN, un État dans l'État. #Algérie

Ouvrir



Tunisia News @tunisnews

1 min

#Tunisie #Tunisia Tunisie - Politique : Soumaya Ghannouchi s'attaque à la France bit.ly/14Q4tKu

Ouvrir



Tunisia News @tunisnews

1 min

Urgent : Report de la grève des transporteurs de carburants : Tunisie : Actualités : Société : ... goo.gl/8wS0F #Tunisie #Tunisia

Ouvrir



SAT-7 Arabic @sat7ar

1 min

youtu.be/qEzWEEUC07w?a via @YouTube برنامج كرمالك عراق - القدس

Afficher le média



koora @KooraTunisie

2 min

Mi-temps : RM-DOR 0-0 (1-4)

Ouvrir



radioexpressfm @radioexpressfm

3 min

L'interview de Hedi Ben Abbes dans Tac au Tac fb.me/IHx5mvzb

Ouvrir



@khaw/amagreb1 الاخوان عدو الاتمان

4 min

fb.me/PRHluygX علماء: فيروس التهاب الكبد الفيروسي عمره أكثر من 82 مليون سنة

#أوووووووف

Ouvrir



# Notre Objectif

- Détection de l'ensemble des événements *EV*.
- Trouver les dates saillantes.
- Traiter des *tweets* qui portent sur la Tunisie.
- Corpus :
  - Nous avons extrait 276 505 tweets, entre le 08/02/2012 et 15/04/2012 (67 jours).



# Dialecte tunisien

j'ai voté, ta7ya tounes #TnElec #Vote

تصويرة بن علي رجعت في حلق الوادي

Retour de Ben Ali à La Goulette <http://t.co/RqVXr5Hu> #tunisie #tnelec



# Détection des évènements :

## *Twitter*

- *Twitter* constitue un excellent moyen pour diffuser des informations, pour discuter des évènements et pour donner des avis.
- Données accessibles via APIs.
- Plusieurs recherches ont montré que le contenu de *Twitter* reflète étroitement l'intérêt et les préoccupations des utilisateurs en temps réel (Becker2011, Ozdikis2012).



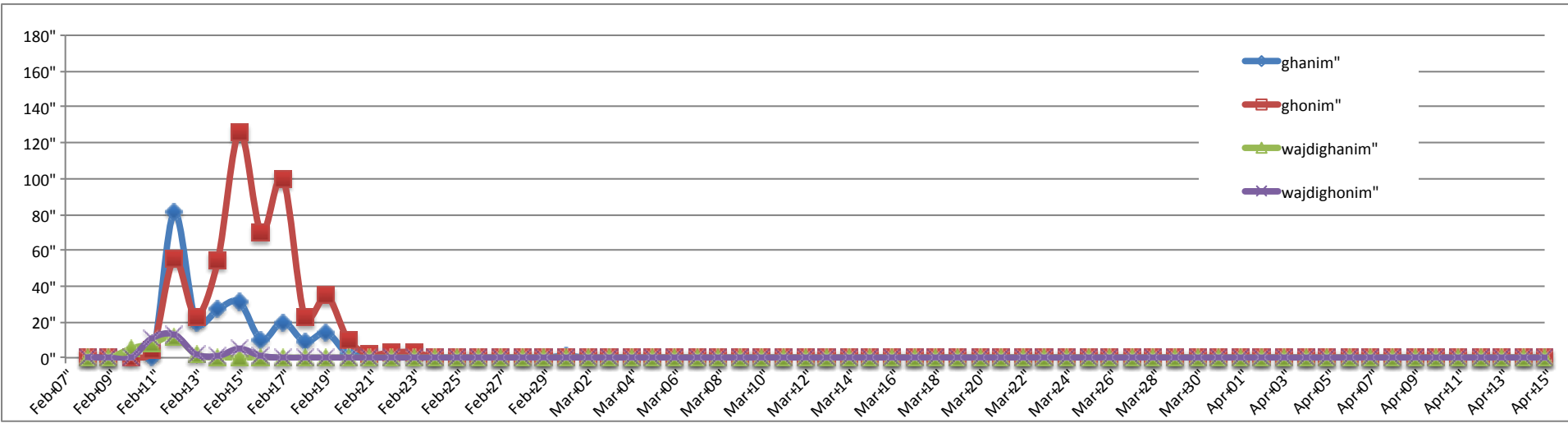


# Détection des évènements

- Un évènement est représenté par un ensemble de termes.
  - Disparition de l'avion *Malaysia Airlines*: {*#PrayForMH370*, *#MH370*, *#MH370Flight*, *#MalaysiaAirlines*, etc.}.
- Objectif :
  - Regrouper automatiquement les termes représentant un même sujet.
  - Trouver l'ensemble d'évènements *EV* (fréquence des termes, etc.).
- Défi :
  - Supporter les conventions d'écriture, les fautes, la taille réduite d'un message, etc.

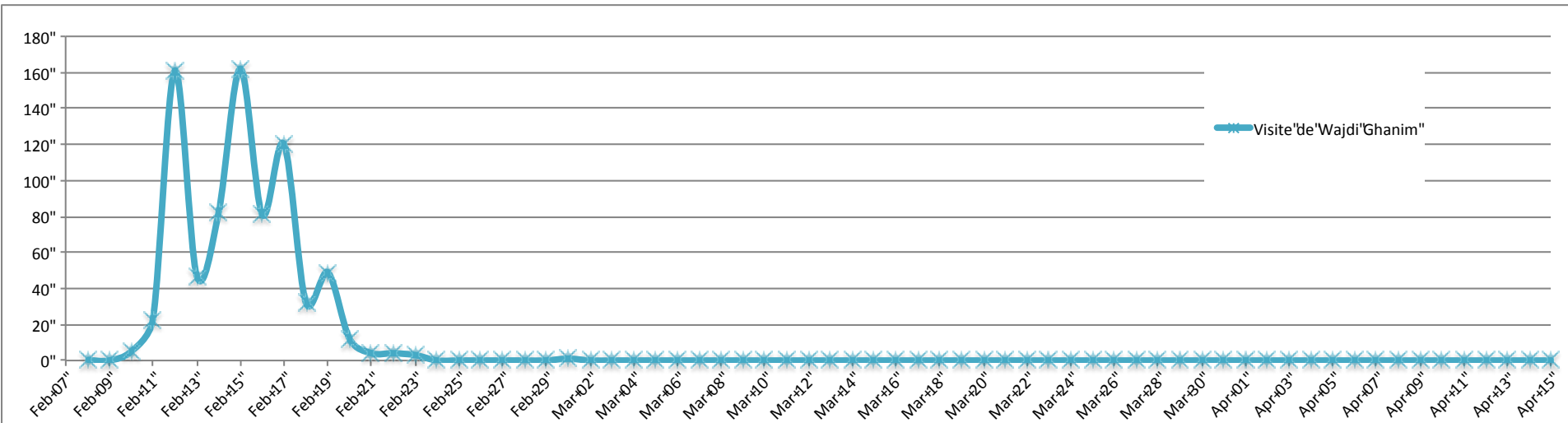


# Regroupement : pourquoi ?



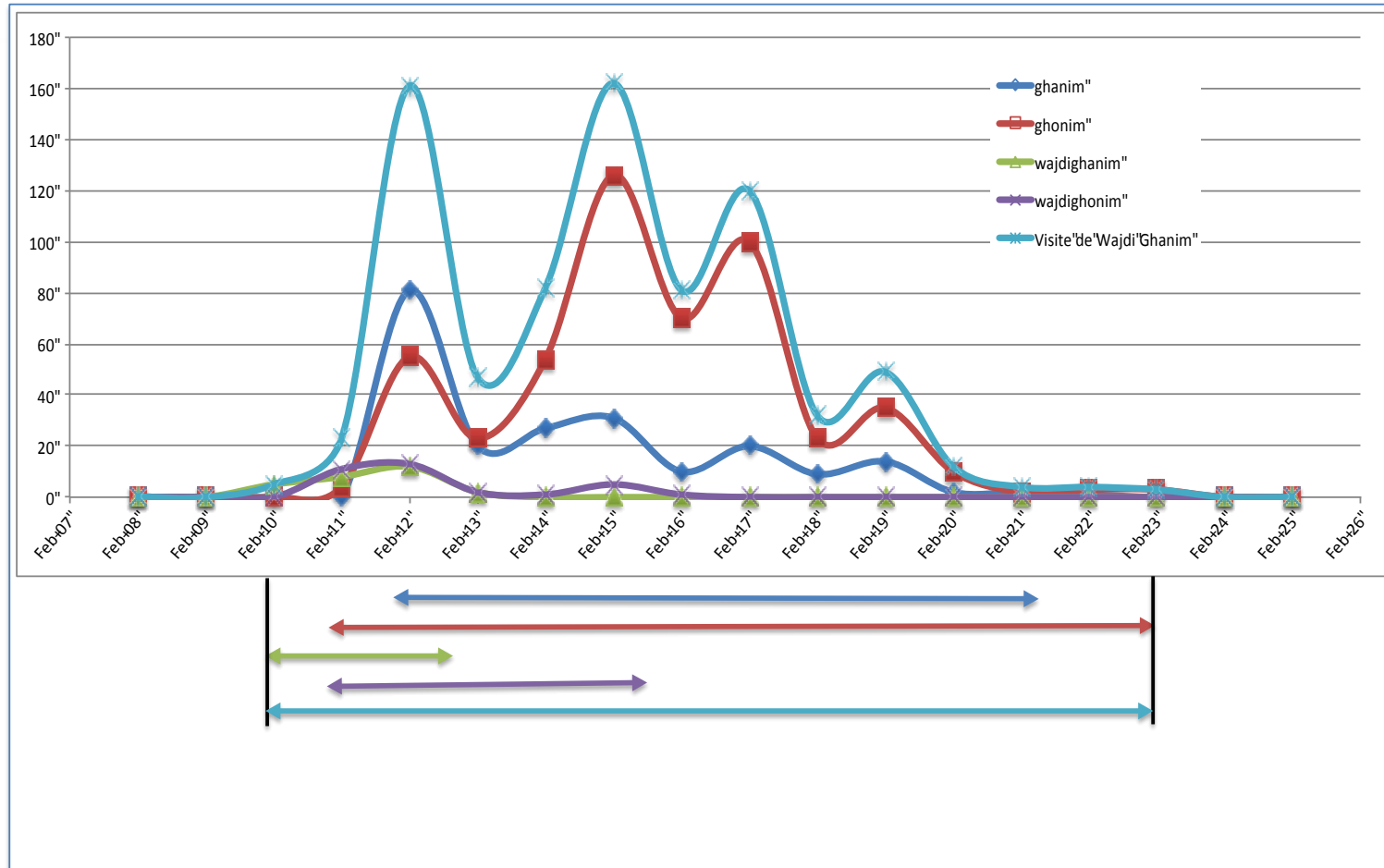


# Regroupement : pourquoi ?



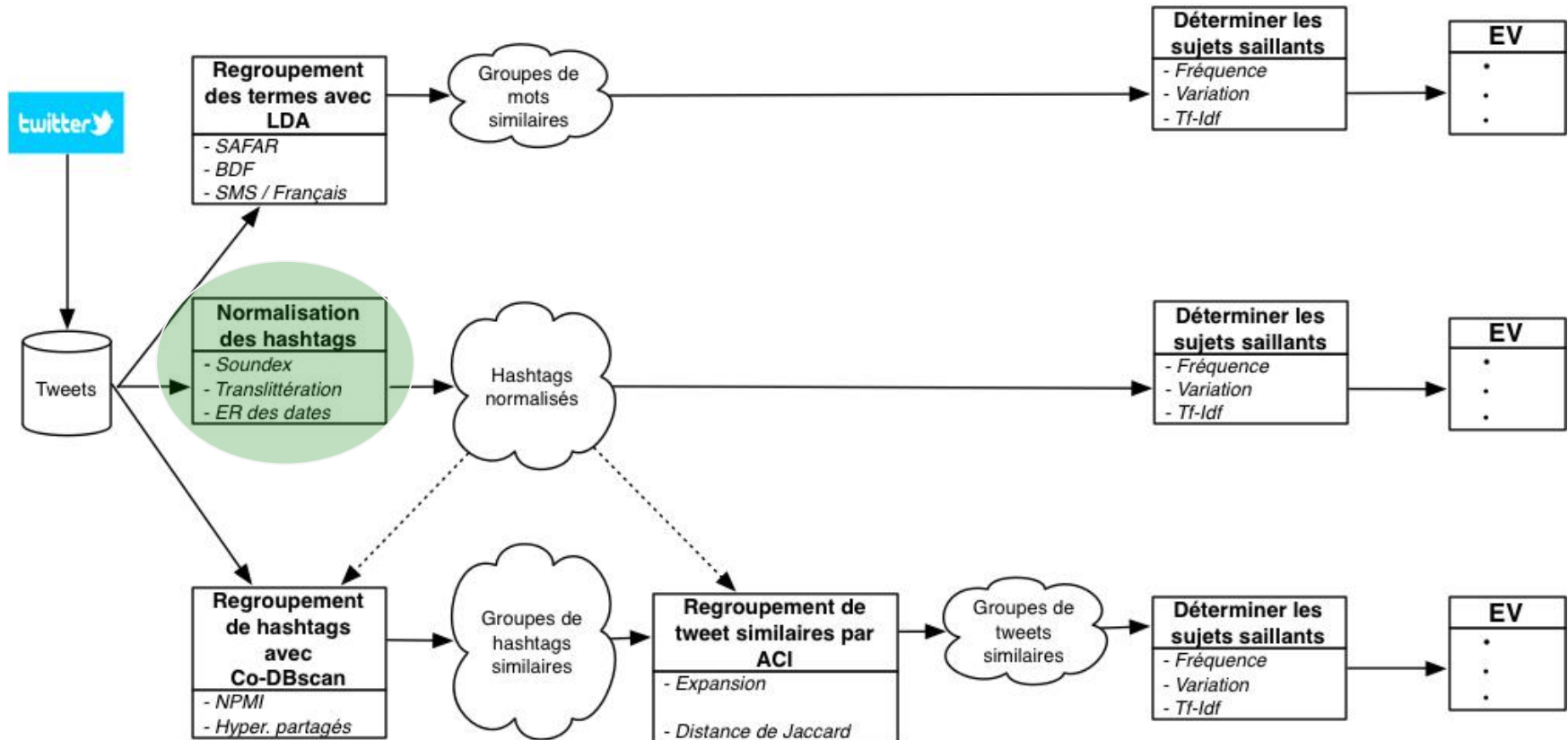


# Regroupement : pourquoi ?





# Démarche : vue globale





# Normalisation des hashtags

- **Soundex** [*Russel et Odell, 1922*]:

- Normaliser les termes qui ont une même prononciation.
- Code les mots qui ont la même prononciation par la même chaîne de caractères.

GHNMO = {#ghanim, #ghenim, #ghnaim, #ghnim, #ghoneim, #ghonem, ...}.

S0000 = {#sousse, #suisse, #ouais, ...} (Noisy Soundex)

- **Normalisation des dates :**

- 9avril = {#9avril, #9april, #9avil, #9avirl, #9أفريل}

- **Translittération :**

غنيمة → GHNMO



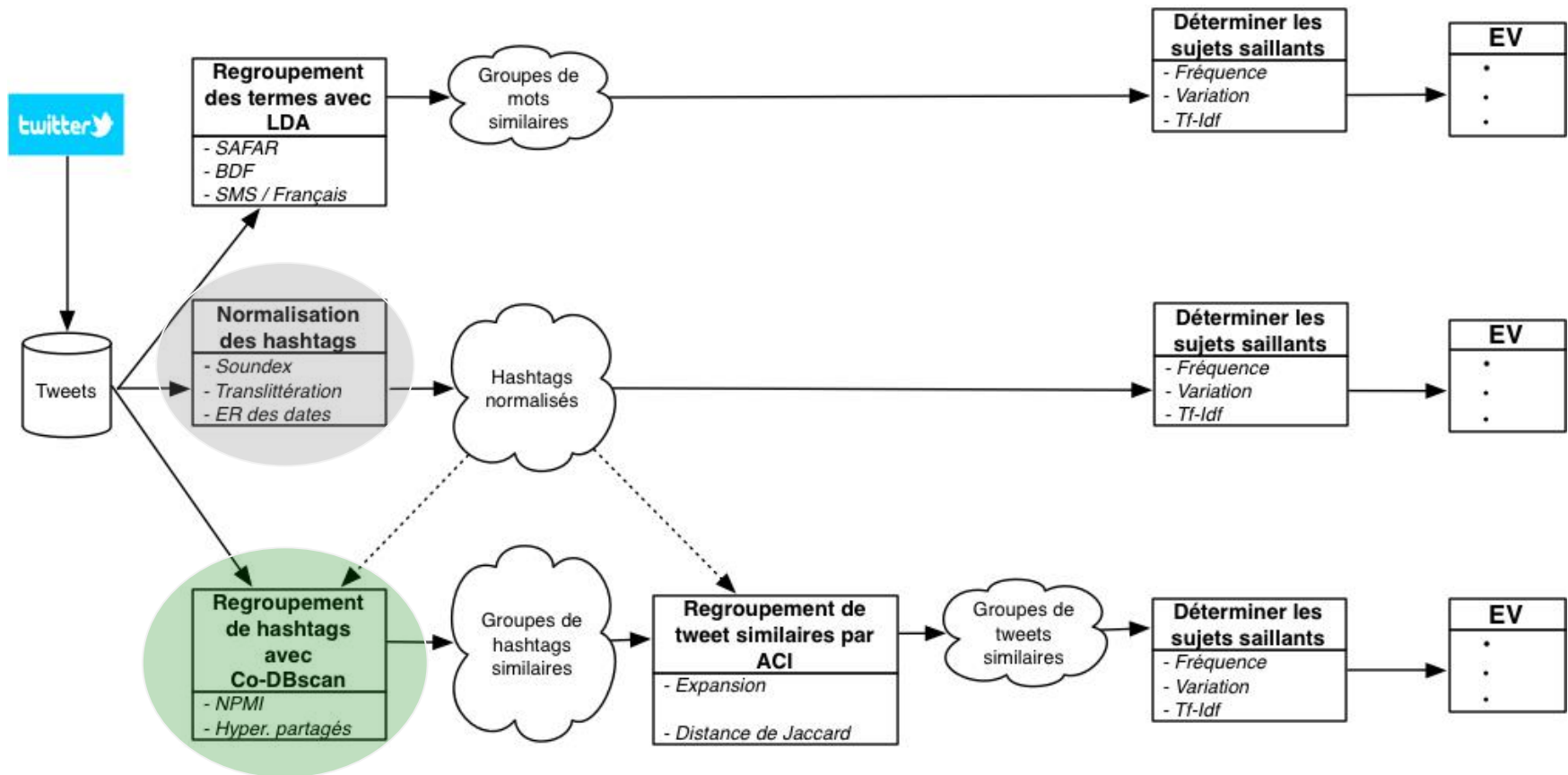
# Normalisation des hashtags : Résultats

<b>Nb. hashtags</b>	12 218
<b>Nb. hashtags écrits en latin</b>	11 693
<b>Nb. groupes Soundex initiaux</b>	7 810
<b>Nb. groupes de Soundex après la normalisation des dates</b>	7 781
<b>Nb. groupes de Soundex en disjoignant les <i>Noisy Soundex</i></b>	8 750
<b>Nb. groupes de Soundex après la translittération</b>	9 033

**Précision** :  $\approx 96\%$ . ( 15 plus importants clusters)



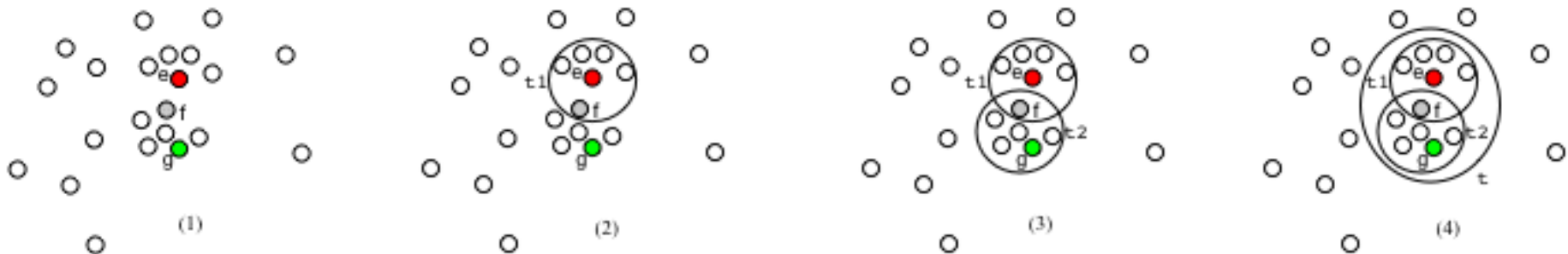
# Démarche : vue globale





# Regroupement hashtags : *DBScan* [Ester, 1996]

- Nombre de *clusters* est déterminé par l'algorithme.
- Paramètres : *MinPts*, *epsilon*.



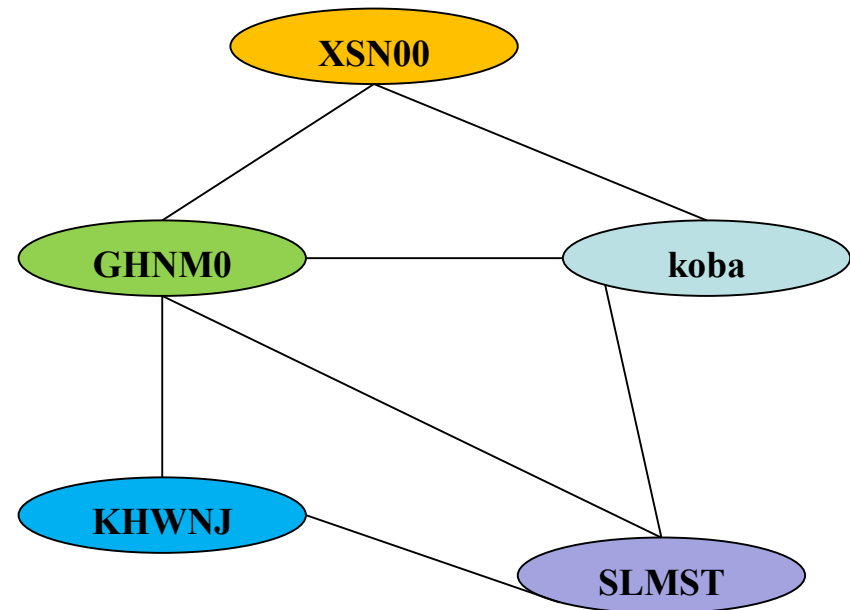
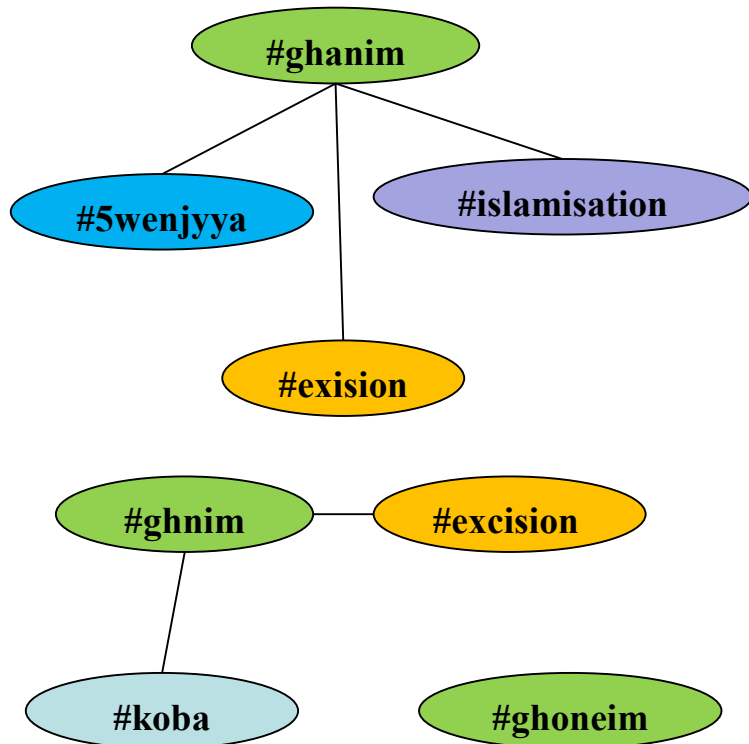


# Regroupement hashtags : *DBScan* [Ester, 1996]

- **Mesures de similarité :**
  1. *NPMI* (*Normalized Pointwise Mutual Information (NPMI)*) : Les hashtags qui apparaissent ensemble, sont sémantiquement similaires.
  2. *Hyper* : Deux hashtags apparaissant avec les mêmes hyperliens suggère que ces hashtags portent sur la même sujet.
- **Hashtags considérés :**
  - *NPMI* : 9 793 /12 218 hashtags.
  - *Hyper* : 6 008 /12 218 hashtags.
- **Problème** : création d'un immense *cluster*. **P.ex** :  $Dbscan_{npmi}$  : 6639 hashtags dans un même cluster.
- **Solution** : vérification de similarité avant l'ajout d'un nouvel élément.  
(*CoDbscan*)



# CoDBScan : Amélioration avec Soundex



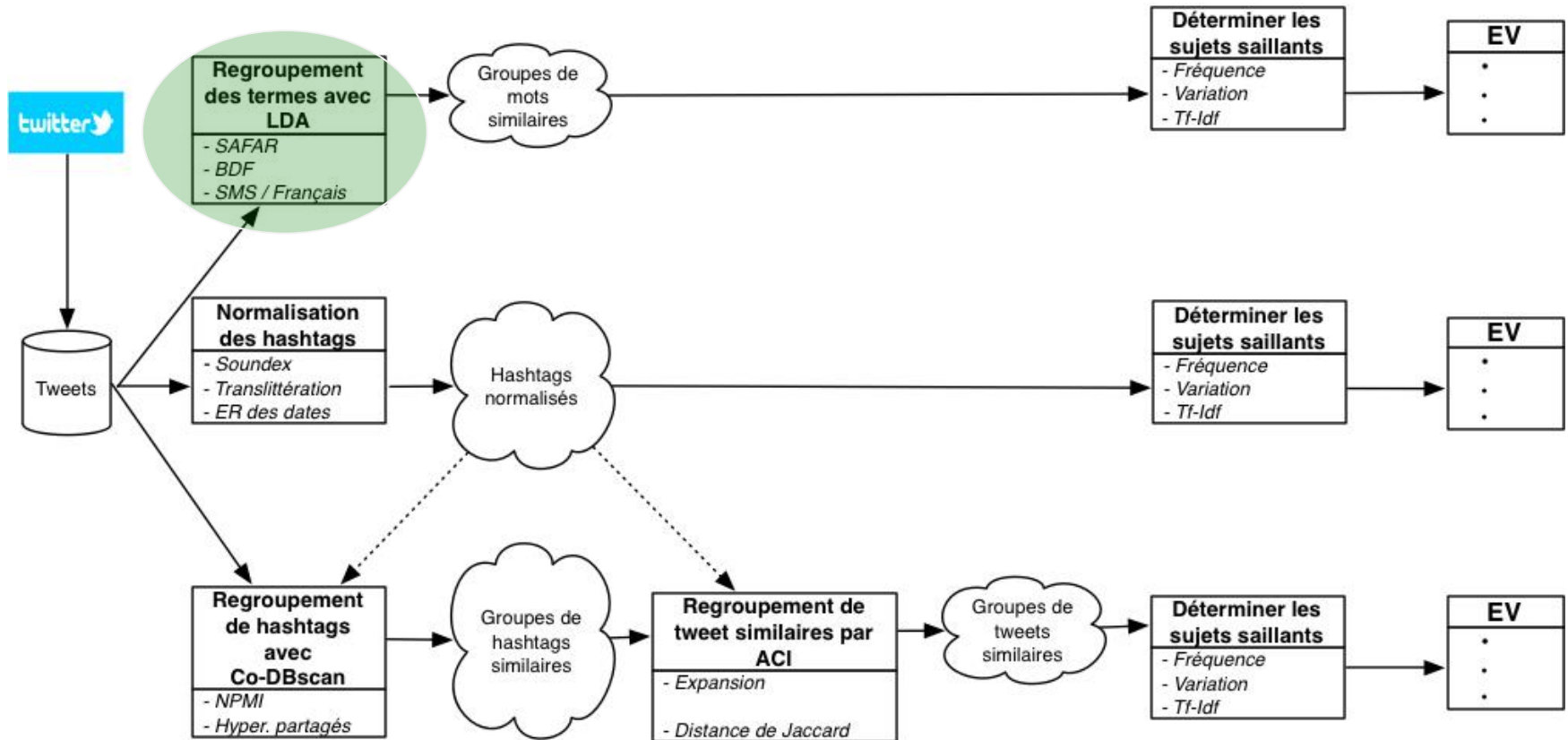


# CoDBScan : Validation manuelle

	NH considérés	Nb clusters	NH regroupés	NH 20 clusters	Précision (20 clusters)
<i>CoDBscan<sub>npmi</sub></i> ( $\epsilon=0,7$ )	9 973	928	5 011	485	94 %
<i>CoDBscan<sub>npmiWithSndx</sub></i> ( $\epsilon=0,6$ )	10 929	908	6 633	509	90 %
<i>CoDBscan<sub>Hyper</sub></i> ( $\epsilon=0,5$ )	6 008	1 213	3 556	457	92 %
<i>CoDBscan<sub>HyperWithSndx</sub></i> ( $\epsilon=0,4$ )	7 431	925	4 660	1 048	91 %



# Démarche : vue globale





# Regroupement : termes

- Utiliser des techniques de *Topic Model* (p.ex. **LDA** [**Blei et al. , 2003**]) pour regrouper les termes liés à un même sujet.
- Chaque document (tweet) peut être représenté comme un mélange de sujets latents, ou un sujet est lui-même représenté comme une distribution des mots qui ont tendance à co-occurrencer.
- Les mots fortement liés à un sujet donné ont les valeurs de probabilité plus grandes.

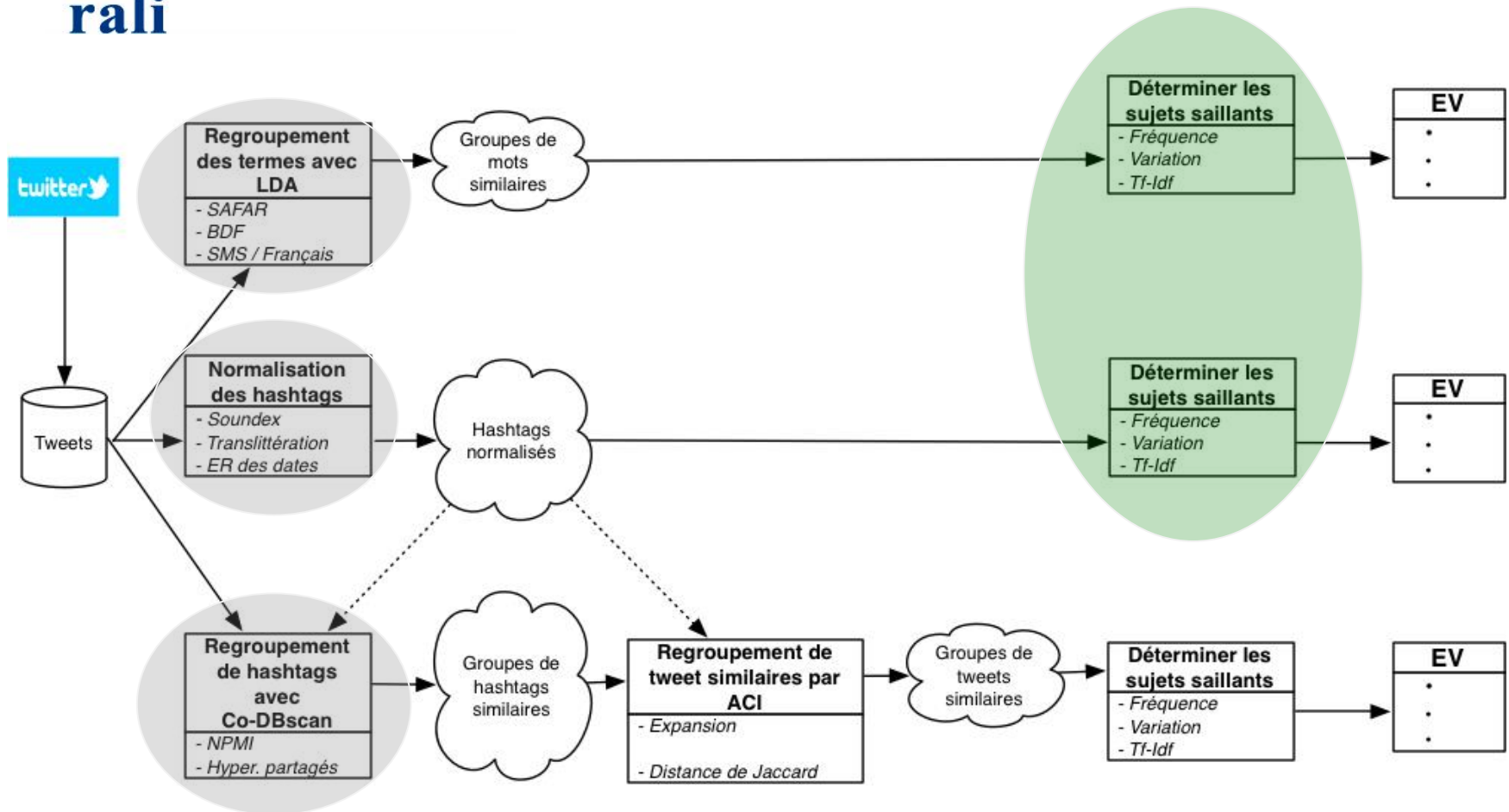


# Amélioration

- Regrouper dans un même document les tweets partageant le (s) même (s) hashtag (s).
- Normaliser certains termes :
  - Lemmatiseur **BDF (Français)** : ‘répond’ et ‘répondu’ sont remplacés par ‘répondre’
  - Corpus parallèle **Français / SMS** : ‘parce que’ -> ‘pcq’.
  - Librairie **SAFAR (arabe)**.



# Démarche : vue globale







# Fréquence de sujets

- Nous avons supposé que chaque groupe réfère à un sujet.
- Un tweet porte sur  $S$  (représenté par un groupe  $g$ ), s'il contient au moins un terme de  $g$ .
- Calculer la fréquence quotidienne de chaque sujet.



# Exemple de résultats

<b>Code Soundex</b>	<b>Hashtags</b>	<b>Dates saillantes</b>	<b>DF</b>	<b>Fréq.</b>	<b>Écart-Type</b>	<b>TF-IDF</b>
9avril	#9avril #9april #9avil #9avirl	09/04	14	6070	516.73	4127,27
MPL00	#empl #emplo #emploi	16/03	67	3025	22.12	0
RCRTM	#recrutement	16/03	67	2854	22.63	0
WTHR0	#weather	15/02, 16/02	65	2043	10.25	68,81



## Exemple de résultats

Code Soundex	Hashtags	Dates saillantes	DF	Fréq.	Écart-Type	TF-IDF
9avril	#9avril #9april #9avil #9avirl	09/04	14	6070	516,73	4127,27
20mars	#20mars, #20مارس	20/03	15	1422	124,59	924,27
ugtt	#ugtt	25/02	50	1797	97,76	182,93
NHD00	#enahda #enanhda #enhada	21/02	66	1872	66,41	12,23

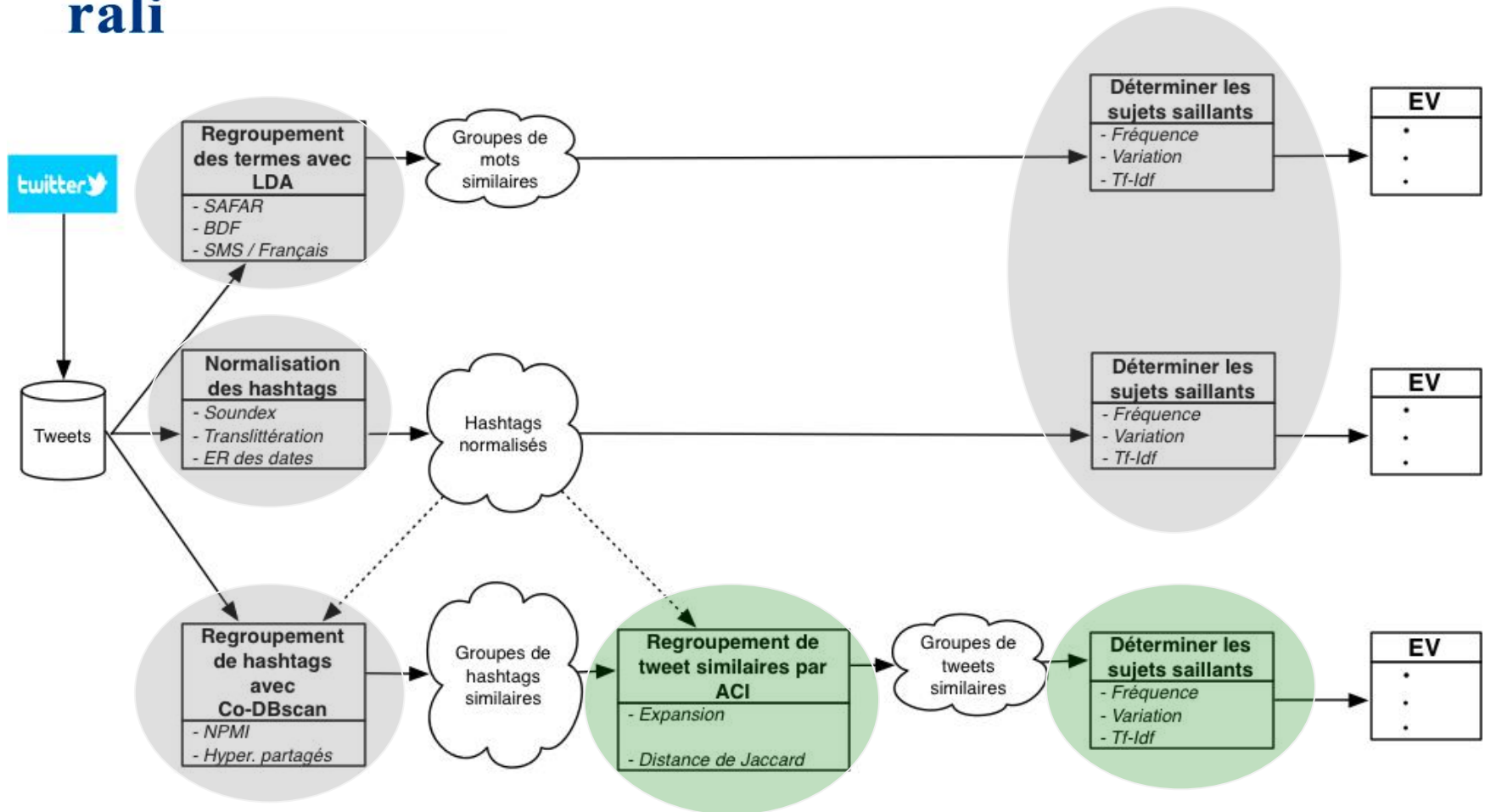


## Exemple de résultats

<b>Code Soundex</b>	<b>Hashtags</b>	<b>Dates saillantes</b>	<b>DF</b>	<b>Fréq.</b>	<b>Écart-Type</b>	<b>TF-IDF</b>
9avril	#9avril #9april #9avil #9avirl	09/04	14	6070	516.73	4127,27
TNPHR	#tunpharma	04/04 04/05	15	1710	52,79	1111,47
20mars	#20mars, #20مارس	20/03	15	1422	124,59	924,27
3PLKT	#application	04/04 04/05	15	1355	41,89	880,72



# Démarche : vue globale





# Regroupement de tweets : *expansion*

- Étendre les tweets par d'autres hashtags reliés afin d'améliorer le regroupement.
- Impact :

	Nb. Tweets non regroupés
Sans expansion	4 296
Expansion <i>CoDBscan</i> <sub><i>npmiWithSoundex</i></sub>	3 077
Expansion <i>CoDBscan</i> <sub><i>npmiWithSoundex</i></sub>	3 305



# Déterminer les événements

: *dates saillantes*

- Méthode de **Palshikar (2009)** :

08-02	09-02	10-02	11-02	12-02	13-02	14-02	15-02	16-02	17-02	18-02	19-02	20-02	21-02	22-02	23-02	24-02
0	0	1	33	190	74	102	176	118	155	40	55	14	6	4	4	1

- k=2

$$S_i = \frac{\max_{1 \leq j \leq k} (x_i - x_{j-1}) + \max_{1 \leq j \leq k} (x_i - x_{i+j})}{2}$$

08-02	09-02	10-02	11-02	12-02	13-02	14-02	15-02	16-02	17-02	18-02	19-02	20-02	21-02	22-02	23-02	24-02
-	-	-15,5	-4	152,5	6,5	6	80	47	76	-26	32	-8	-3	0,5	-	-

$x_i$  est un pic si :  $S_i > \text{moyenne} (S > 0)$  ET  $(S_i - \text{moyenne}) > \text{écart-type} (S > 0)$



# Déterminer les évènements

: *EV*

- Sélectionner que les sujets saillants : *EV* !
- Adapter la méthode **Palshikar (2009)**.
  - $k$  = nombre de sujets,
  - $x_i$  : valeur du critère (*fréquence, écart-type, Tf-Idf*) utilisé.





# Validation

- L'évaluation de l'exactitude de la méthode est une tâche difficile : **pas de données de référence.**
- Deux méthodes utilisées :
  1. Vérification de l'importance et de la date de chaque évènement auprès d'un ensemble des medias traditionnels (journaux numériques) fiables.
  2. Recours à des experts, dans notre cas des Tunisiens au courant des évènements qui se sont déroulés en Tunisie, afin de distinguer les évènements importants parmi ceux détectés par nos méthodes.



# Validation

<http://rali.iro.umontreal.ca:8080/dridihou/>

## analyse des tweets

home   expériences   **annotation**   contact

### Annotation des événements

Num	Termes	Date (2012)	Jugement		
1	9avril (9avril 9avril 9avril 9avril 9avr )	( 04/09 )	Évènement	pas-Évènement	pas clair
2	20mars (20mars 20mars1956 20mars 20mars1956 20mars2012 )	( 03/20 )	Évènement	pas-Évènement	pas clair
3	ugtt (ugtt )	( 02/25 )	Évènement	pas-Évènement	pas clair
4	NHD00 (enahda enanhda enhada ennahada ennahd )	( 02/21 )	Évènement	pas-Évènement	pas clair
5	3KB00 (acab )	( 04/09 )	Évènement	pas-Évènement	pas clair
6	MNB00 (mannoub mannouba manouba mnaouba منوبة )	( 03/07 )	Évènement	pas-Évènement	pas clair
7	TNPHR (tunpharma )	( 04/04 04/05 )	Évènement	pas-Évènement	pas clair
8	PHRMC (pharmacie pharmaciecentrale pharmacies )	( 04/04 04/05 )	Évènement	pas-Évènement	pas clair
9	3LGRY (algeria algerian algria )	( 02/16 02/21 )	Évènement	pas-Évènement	pas clair
10	GYPT0 (egypt egypt egypt egypte égypte )	( 03/02 04/09 )	Évènement	pas-Évènement	pas clair
11	BHRN0 (baharin bahrain bahraïn bahrein bahreïn )	( 02/14 02/15 03/02 )	Évènement	pas-Évènement	pas clair



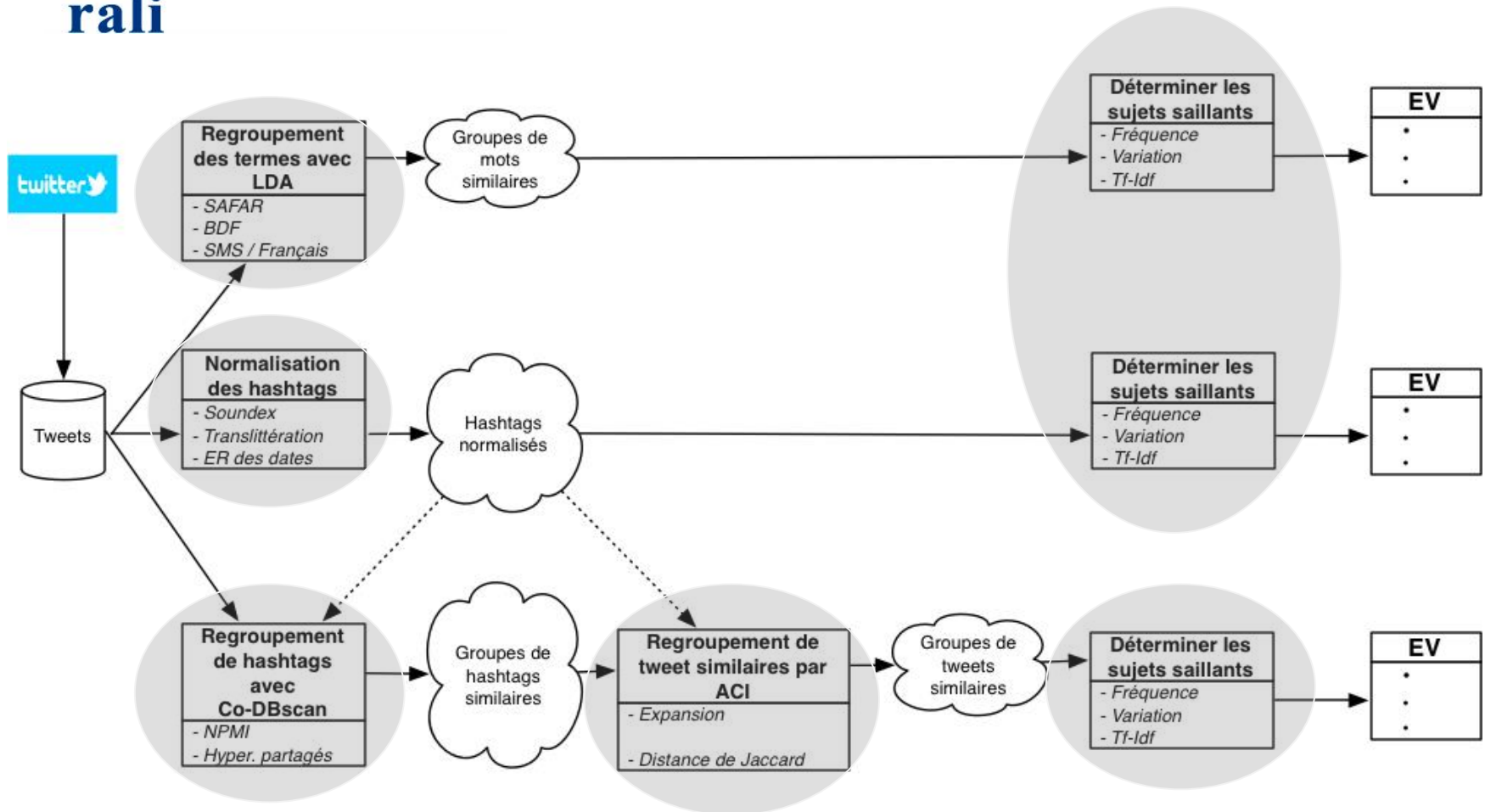
# Évaluation

- Annotation par 10 experts.

	<i>Fréquence</i>		<i>variation</i>		<i>Tf-Idf</i>	
	<i> EV </i>	<b>Pr.</b>	<i> EV </i>	<b>Pr.</b>	<i> EV </i>	<b>Pr.</b>
<b>Groupe de hashtags (normalisation)</b>	123	<b>0,64</b>	88	0,82	81	<b>0,95</b>
<b>Groupe termes (LDA)</b>	53	<b>0,73</b>	67	0,87	74	<b>0,92</b>
<b>Groupe de tweets + expansion</b> <i>CoDBscan<sub>npmiWithSoundex</sub></i>	98	<b>0,61</b>	52	0,8	46	<b>0,93</b>
<b>Groupe de tweets + expansion</b> <i>CoDBscan<sub>HyperWithSoundex</sub></i>	104	<b>0,54</b>	51	0,69	46	<b>0,94</b>



# Démarche : vue globale





# Travaux futurs (1)

- Utiliser d'autres éléments pour identifier *EV*. P.ex : hyperliens.
- Tester nos méthodes sur d'autres corpus :
  - Environ 600 000 tweets portant sur la Tunisie (octobre 2012 – janvier 2013).
  - Plus que 5 millions de tweets envoyés à partir du Québec. (février 2014 – août 2014).
- Déterminer l'opinion publique pour un évènement :
  - Déterminer la proportion de chaque polarité (positive, négative, neutre) à partir des tweets assignés à l'évènement.
  - Pas de ressources disponibles!



# Travaux futurs (2)

- Environ 4 000 tweets annotés.

## analyse des tweets

home	expériences	annotation	contact
<b>Construction d'ensemble d'apprentissage</b>			
Num	Tweet	Langue	Sentiment
1	En Tunisie, il est difficile d'assumer ses missions d'universitaires. Aidons-les en diffusant leurs messages! https://t.co/y0tkxPmO	<input checked="" type="checkbox"/> Français <input type="checkbox"/> Anglais <input type="checkbox"/> Arabe <input type="checkbox"/> Tunisien <input type="checkbox"/> Mixte <input type="checkbox"/> Autre langue	positive négative neutre Information insuffisante
2	حاجب العيون : أسعار العلف من نار ... والسماسة على الخط http://t.co/NeKaTxGj #Tunisie #Tunisia		positive négative neutre Information insuffisante
3	العلمانيين لا يتحملون ديمقراطيتهم العفنة .. شكوى ضد الداعية وجدي غنيم في تونس http://t.co/rPC2xNvp #tunisie		positive négative neutre Information insuffisante
4	Souvenez-vous des guignols qui venaient nous réciter leurs programmes à la #TTN. C'était ça leur campagne électorale! #LesInconnus #tnElec		positive négative neutre Information insuffisante
5	Tunisie: L'INRIC qualifie "d'injustes" les accusations contre les journalistes http://t.co/YvEissHM		positive négative neutre Information insuffisante
6	Tunisie: Le reflet de l'ombre de Naceur Belhaj Bettaïeb: [La Presse] S'appuyant... http://t.co/7RCdzTZn #artisanat		positive négative neutre Information insuffisante
7	Titre alternatif: 'La blague de l'année' <a href="http://t.co/SPvWEnCr">http://t.co/SPvWEnCr</a>		positive négative neutre Information insuffisante
8	#immobilier TUNISIE LEASING : Les caractéristiques de l'opération d'absorption par la société de sa filiale la... http://t.co/gWANAGFQ		positive négative neutre Information insuffisante



# Conclusion

- Les hashtags sont des indices pour déterminer les évènements.
- La fréquence ne reflète pas l'importance d'un sujet.
- Tâches de validation est coûteuse.
- Notre système : haute précision.



[@toutLeMonde](#): Merci pr vtr attention. [#soutenance](#) [#thèse](#) [#Houssem](#)