

**ANR CORPUS GEOMEDIA**

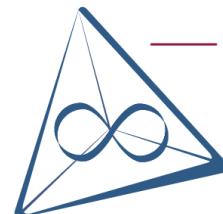
Paris, 24-25 novembre 2014

# Agrégation multiéchelle de l'information médiatique

Robin Lamarche-Perrin

Jean-Marc Vincent

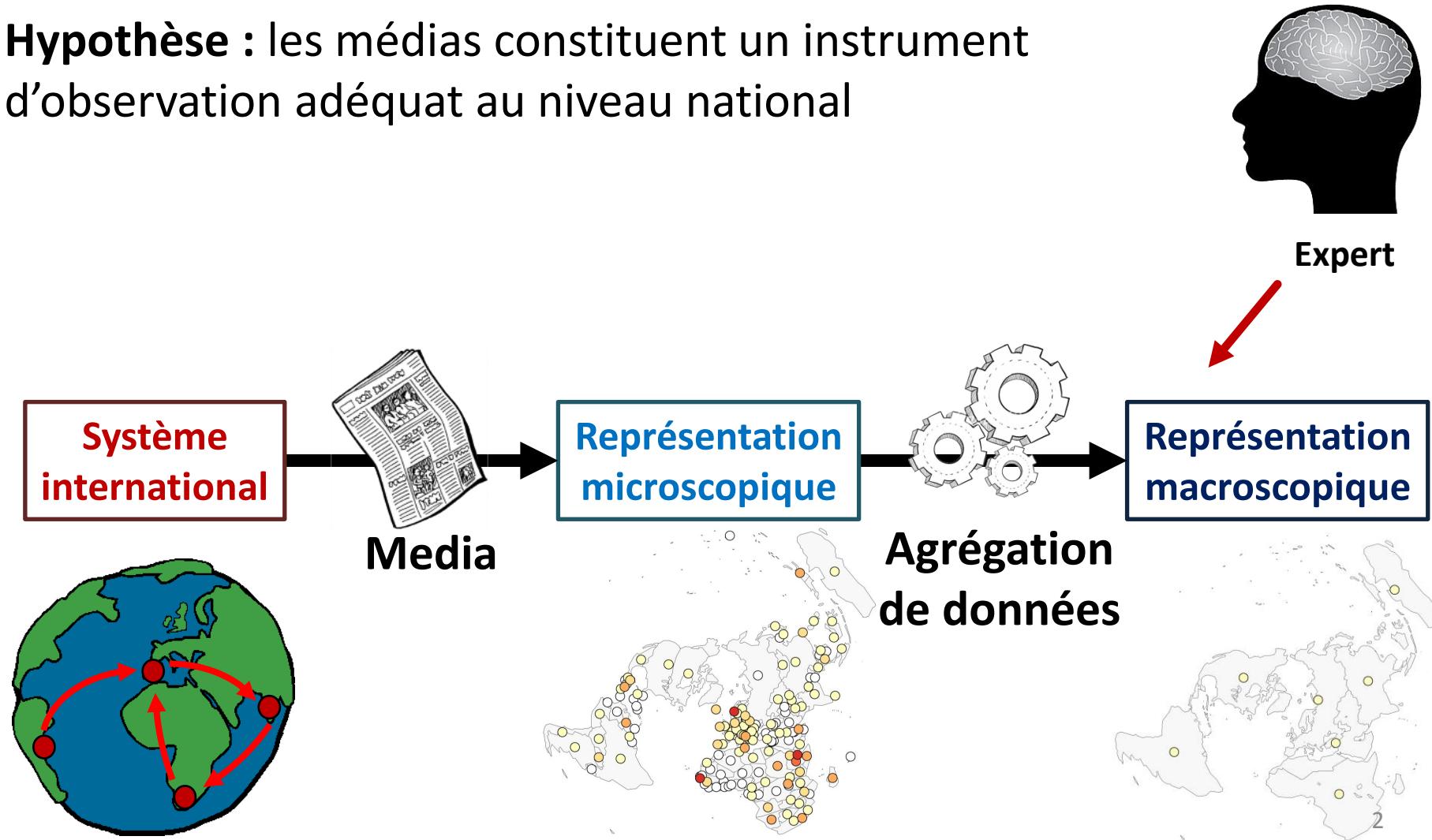
Yves Demazeau



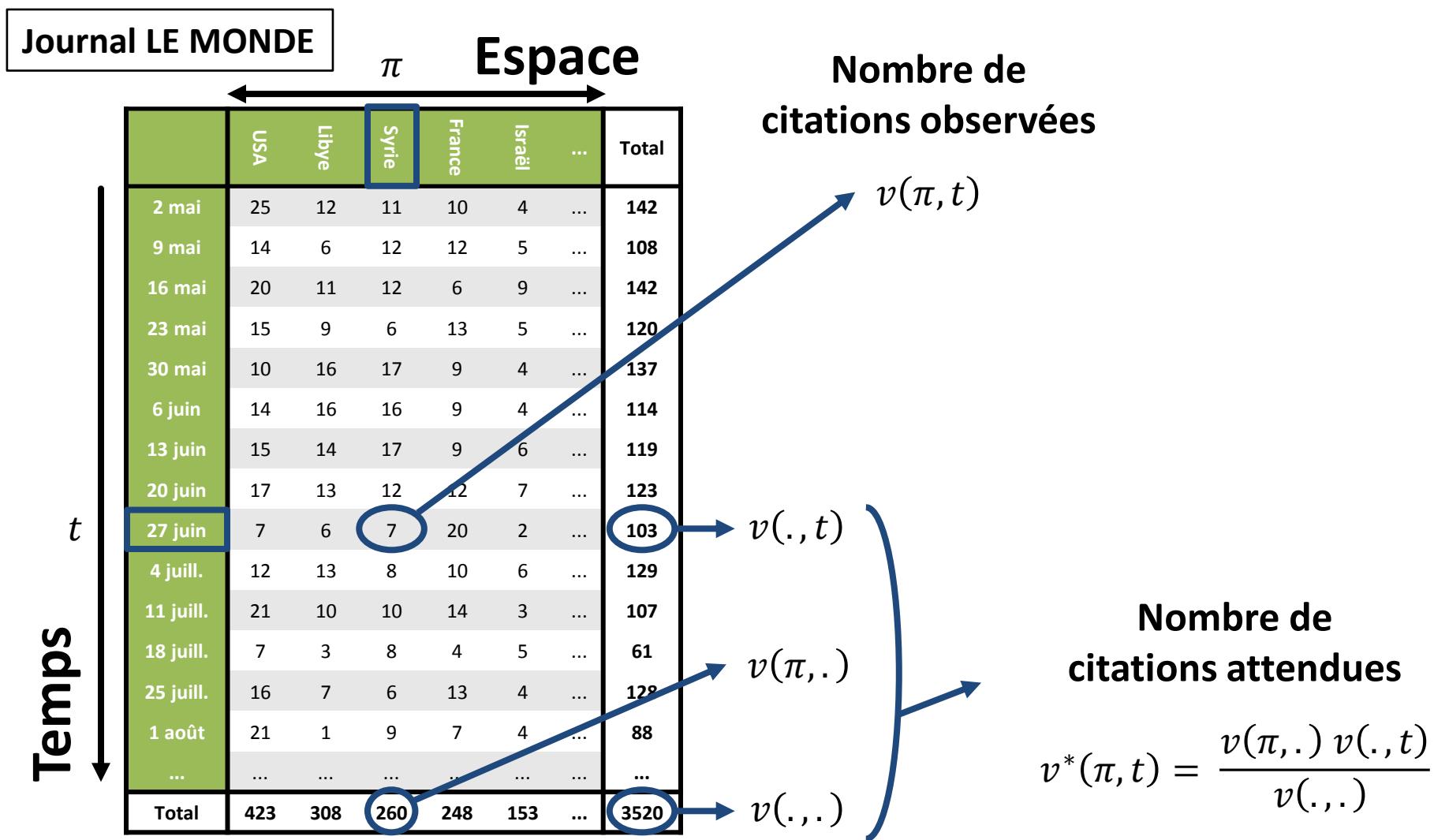
Max Planck Institute for  
**Mathematics**  
in the Sciences

# Analyse médiatique des relations internationales

**Hypothèse :** les médias constituent un instrument d'observation adéquat au niveau national



# Représentation microscopique du système international



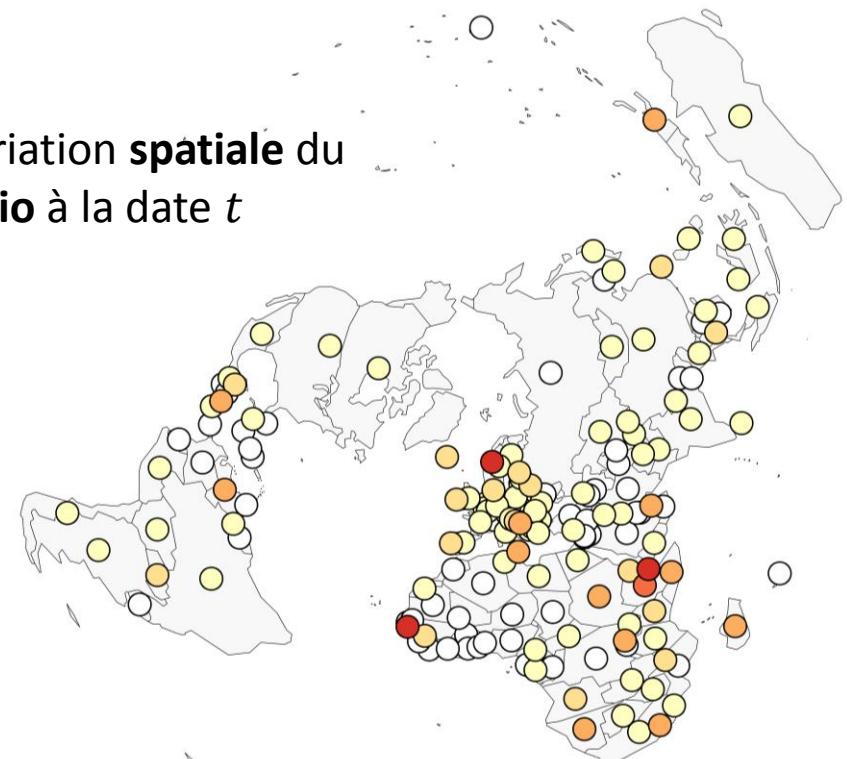
# Représentation géographique

**Espace**

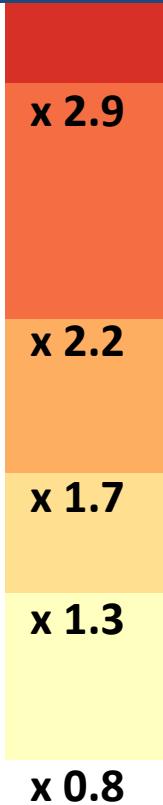
	USA	Libye	Syrie	France	Israël	...	Total
2 mai	25	12	11	10	4	...	142
9 mai	14	6	12	12	5	...	108
16 mai	20	11	12	6	9	...	142
23 mai	15	9	6	13	5	...	120
30 mai	10	16	17	9	4	...	137
6 juin	14	16	16	9	4	...	114
13 juin	15	14	17	9	6	...	119
20 juin	17	13	12	12	7	...	123
27 juin	7	6	7	20	2	...	103
4 juill.	12	13	8	10	6	...	129
11 juill.	21	10	10	14	3	...	107
18 juill.	7	3	8	4	5	...	61
25 juill.	16	7	6	13	4	...	128
1 août	21	1	9	7	4	...	88
...	...	...	...	...	...	...	...
Total	423	308	260	248	153	...	3520

**Ratio des citations observées et des citations attendues**

$$\rho(\pi, t) = \frac{v(\pi, t)}{v^*(\pi, t)} = \frac{v(\pi, t) v(., .)}{v(\pi, .) v(., t)}$$

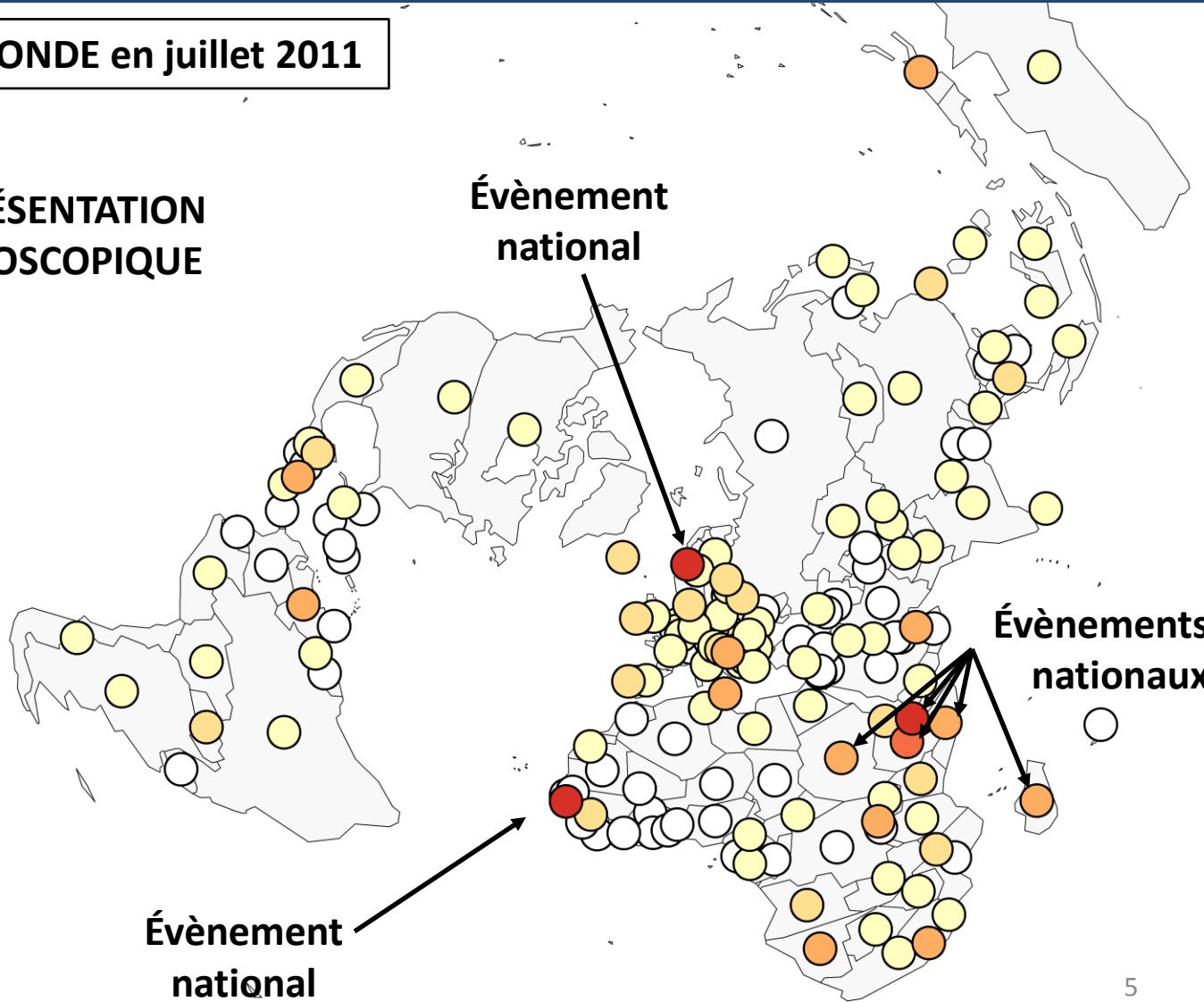


# Détection d'évènements médiatiques

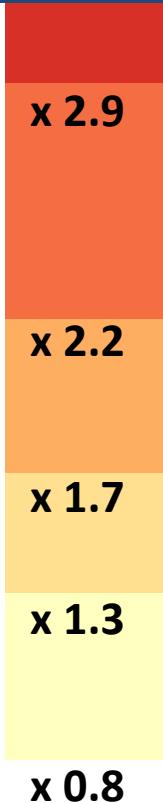


Journal LE MONDE en juillet 2011

REPRÉSENTATION  
MICROSCOPIQUE

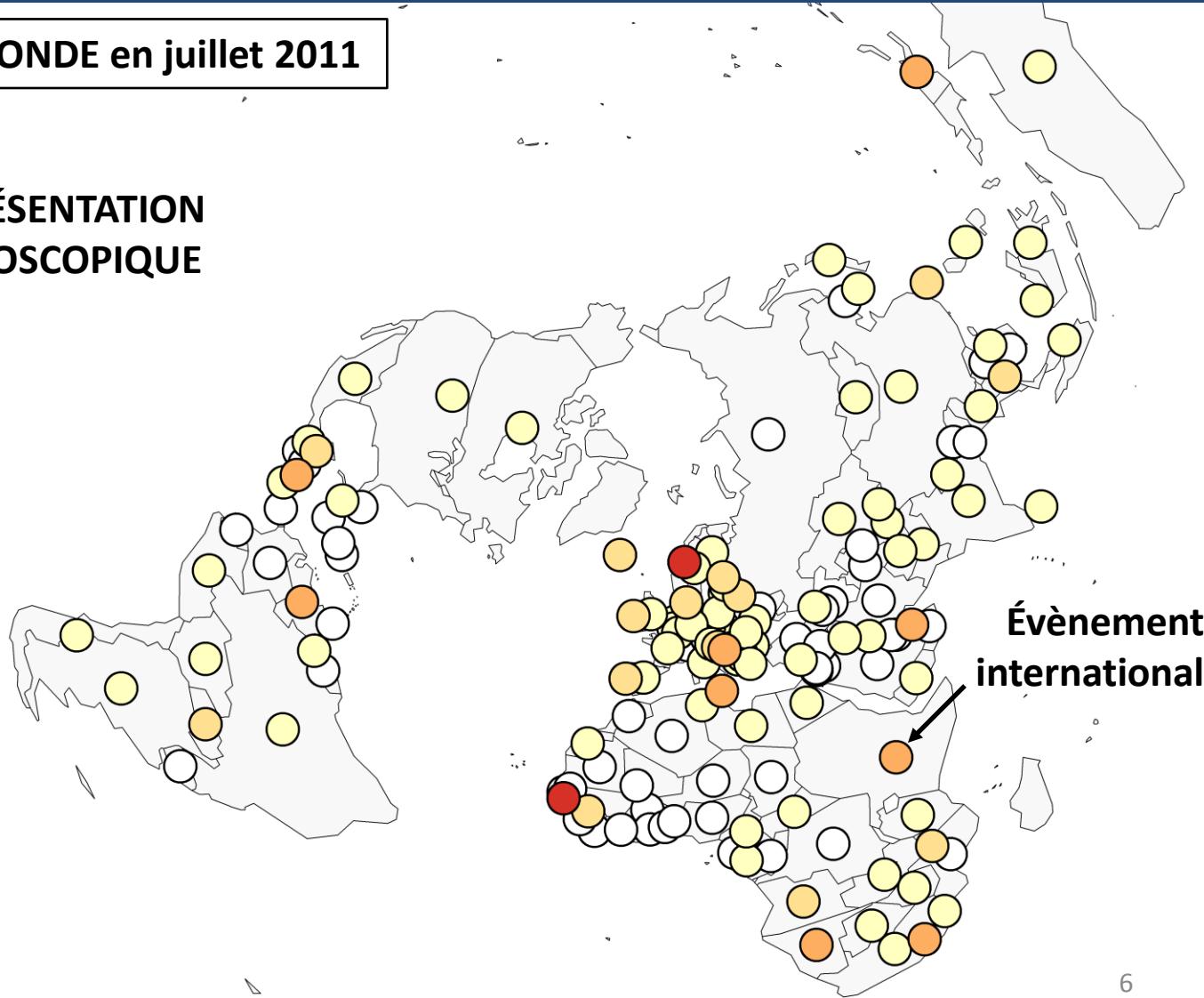


# Détection d'évènements médiatiques



Journal LE MONDE en juillet 2011

REPRÉSENTATION  
MICROSCOPIQUE

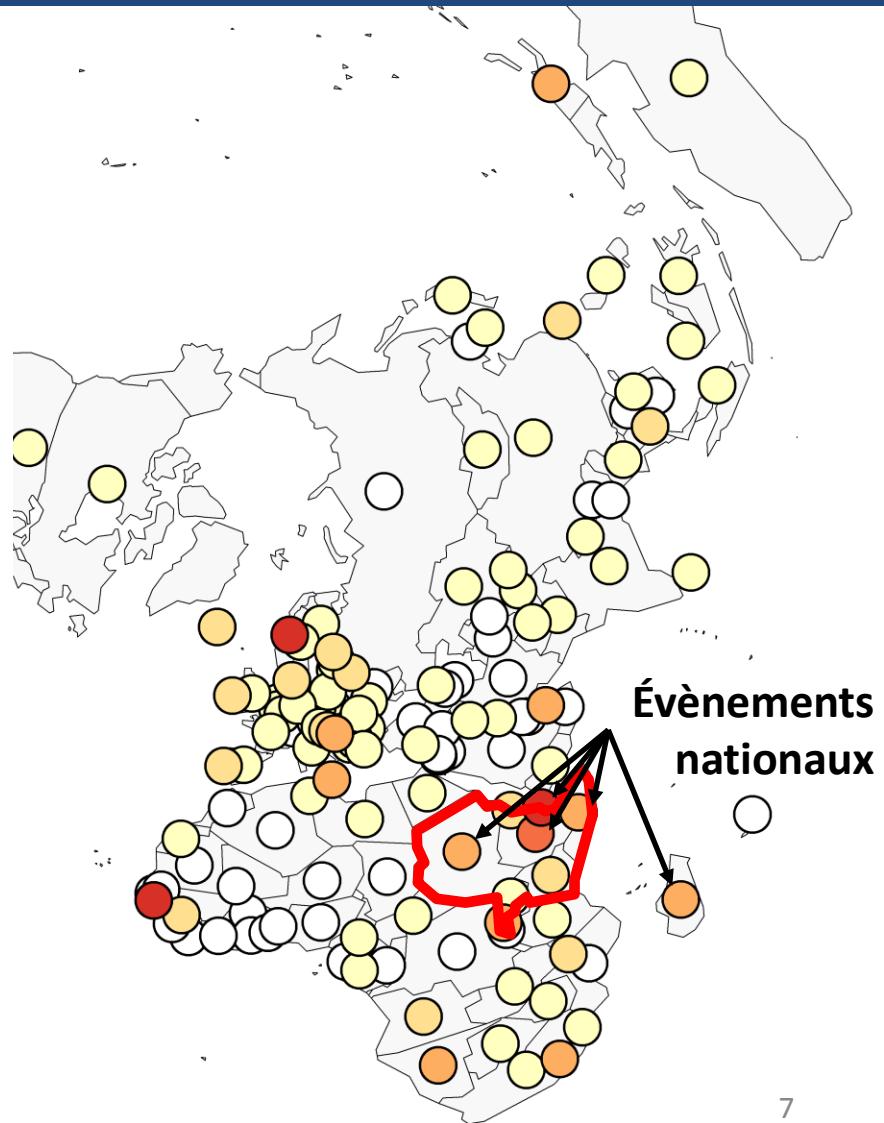


# Agrégation de données

**Espace**

Temps

	USA	Libye	Syrie	France	Israël	...	Total
2 mai	25	12	11	10	4	...	142
9 mai	14	6	12	12	5	...	108
16 mai	20	11	12	6	9	...	142
23 mai	15	9	6	13	5	...	120
30 mai	10	16	17	9	4	...	137
6 juin	14	16	16	9	4	...	114
13 juin	15	14	17	9	6	...	119
20 juin	17	13	12	12	7	...	123
27 juin	7	6	7	20	2	...	103
4 juill.	12	13	8	10	6	...	129
11 juill.	21	10	10	14	3	...	107
18 juill.	7	3	8	4	5	...	61
25 juill.	16	7	6	13	4	...	128
1 août	21	1	9	7	4	...	88
...	...	...	...	...	...	...	...
<b>Total</b>	<b>423</b>	<b>308</b>	<b>260</b>	<b>248</b>	<b>153</b>	...	<b>3520</b>

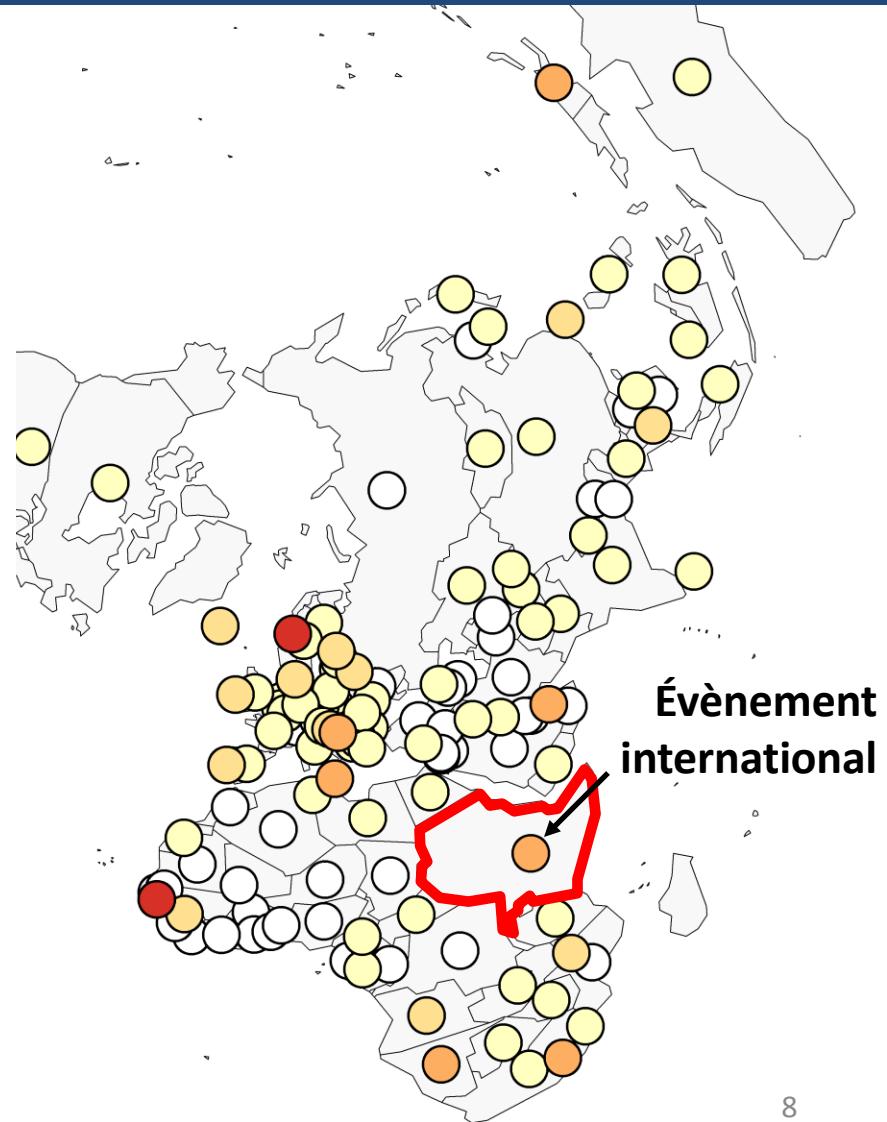


# Agrégation de données

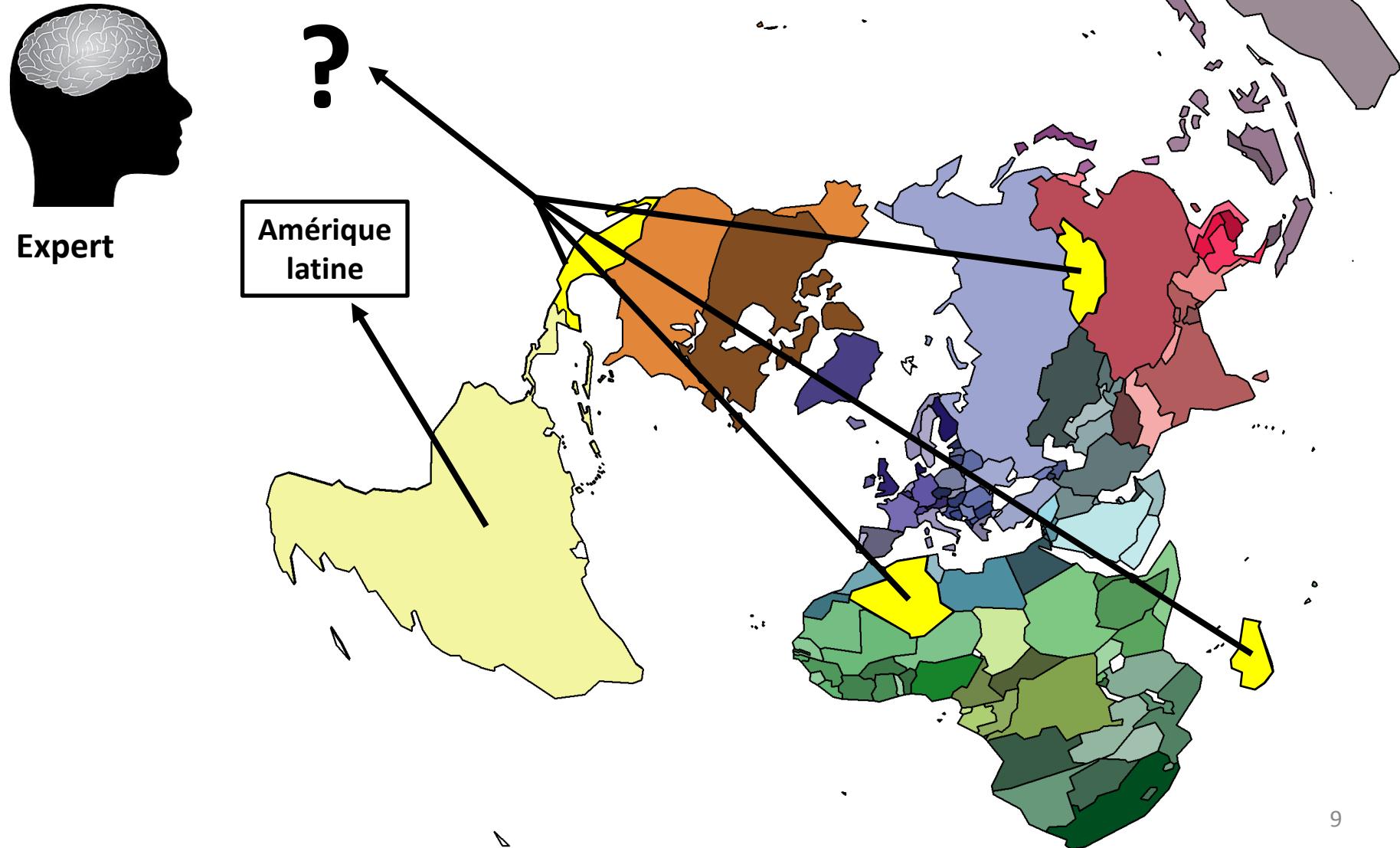
**Espace**

Temps

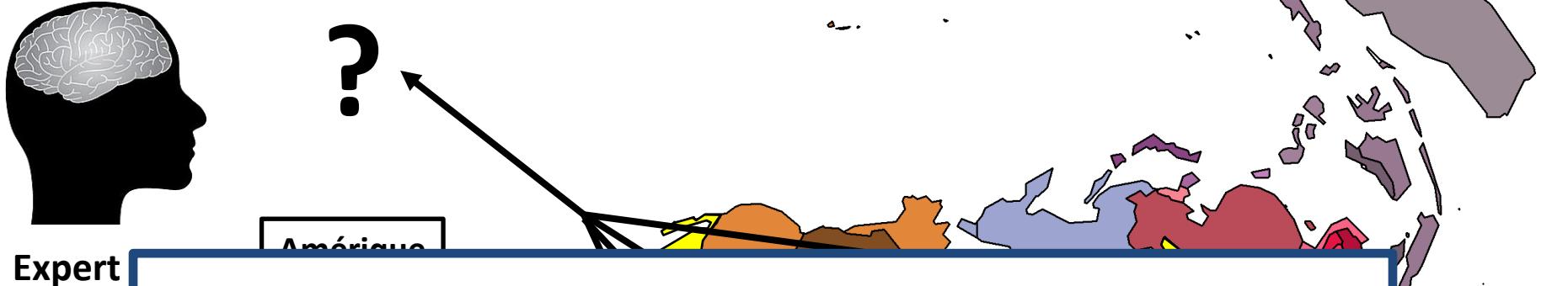
	USA	Agrégat	Israël	...	Total
2 mai	25	13+11+10	4	...	142
9 mai	14	6+12+12	5	...	108
16 mai	20	11+12+6	9	...	142
23 mai	15	9+6+13	5	...	120
30 mai	10	16+17+9	4	...	137
6 juin	14	16+16+9	4	...	114
13 juin	15	14+17+9	6	...	119
20 juin	17	13+12+12	7	...	123
27 juin	7	6+7+20	2	...	103
4 juill.	12	13+8+10	6	...	129
11 juill.	21	10+10+14	3	...	107
18 juill.	7	3+8+4	5	...	61
25 juill.	16	7+6+13	4	...	128
1 août	21	1+9+7	4	...	88
...	...	...	...	...	...
<b>Total</b>	<b>423</b>	<b>308+260+248</b>	<b>153</b>	...	<b>3520</b>



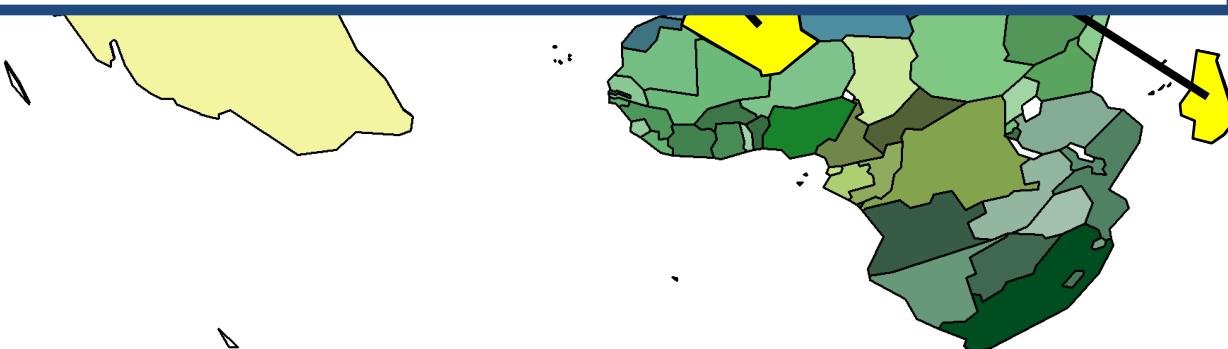
# P1 : Sémantique géographique des agrégats



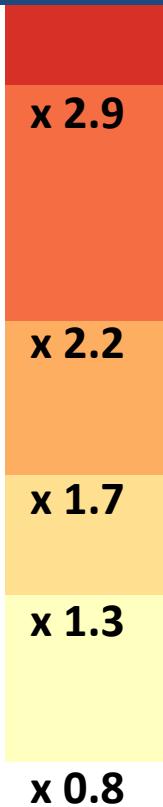
# P1 : Sémantique géographique des agrégats



**Problème 1 :** Comment engendrer des abstractions cohérentes avec l'espace géographique ?

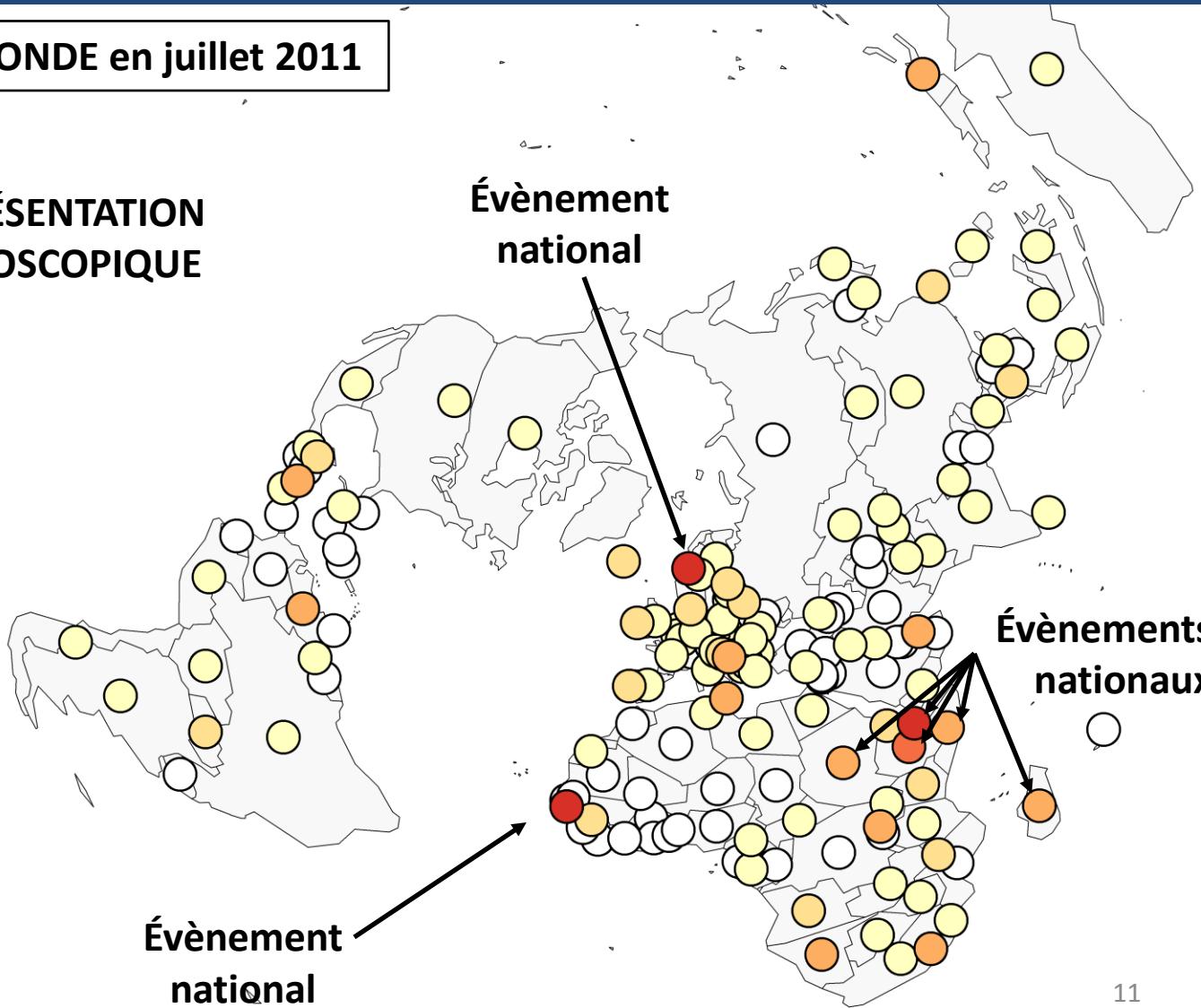


# P2 : Niveaux de représentation

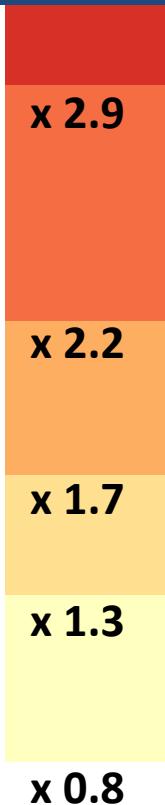


Journal LE MONDE en juillet 2011

REPRÉSENTATION  
MICROSCOPIQUE

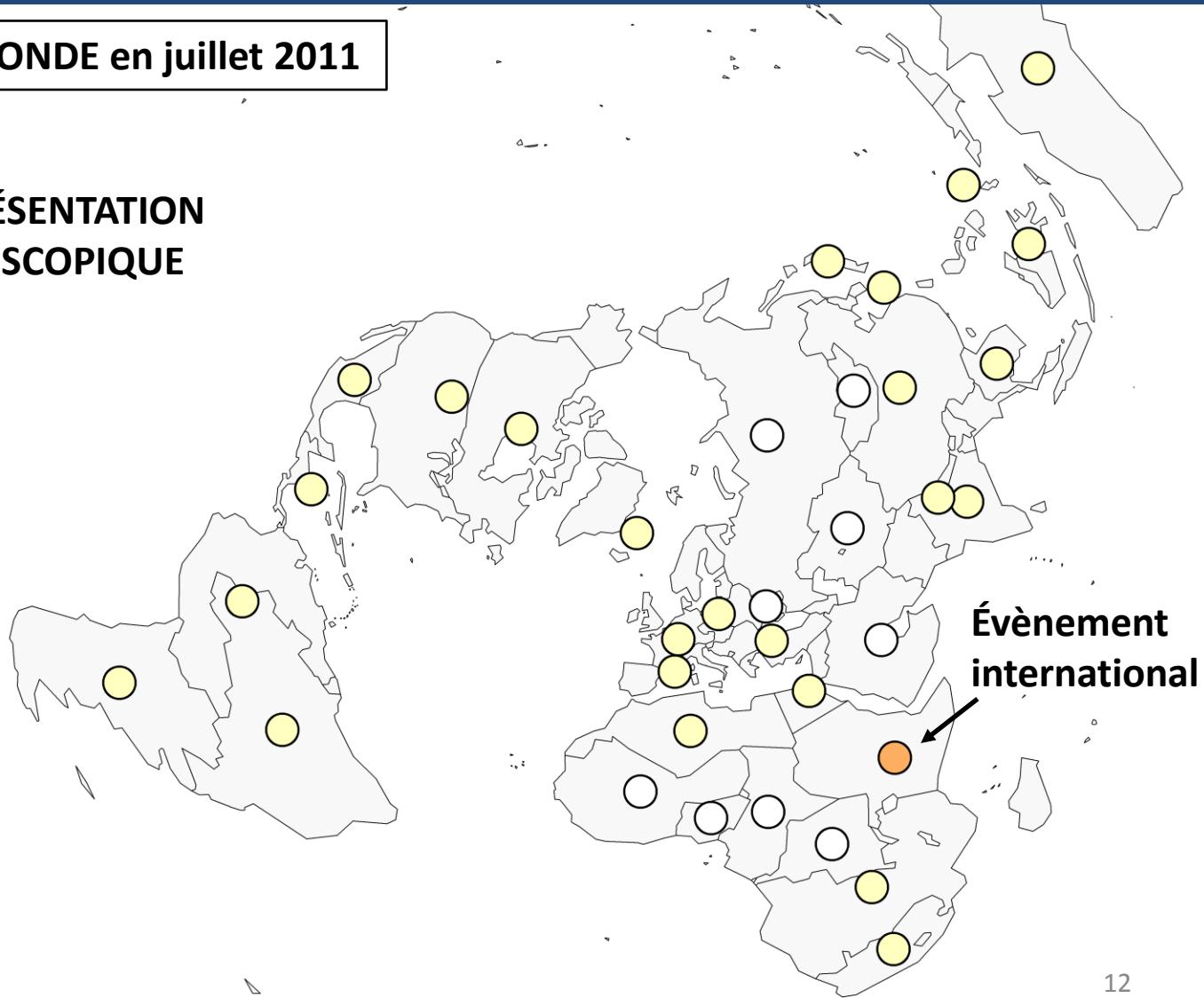


# P2 : Niveaux de représentation

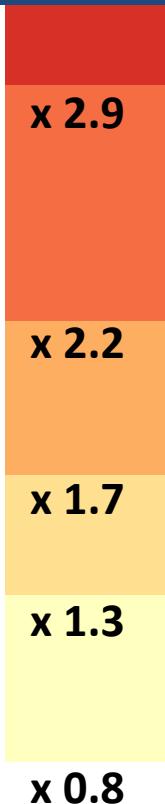


Journal LE MONDE en juillet 2011

REPRÉSENTATION  
MÉSOSCOPIQUE

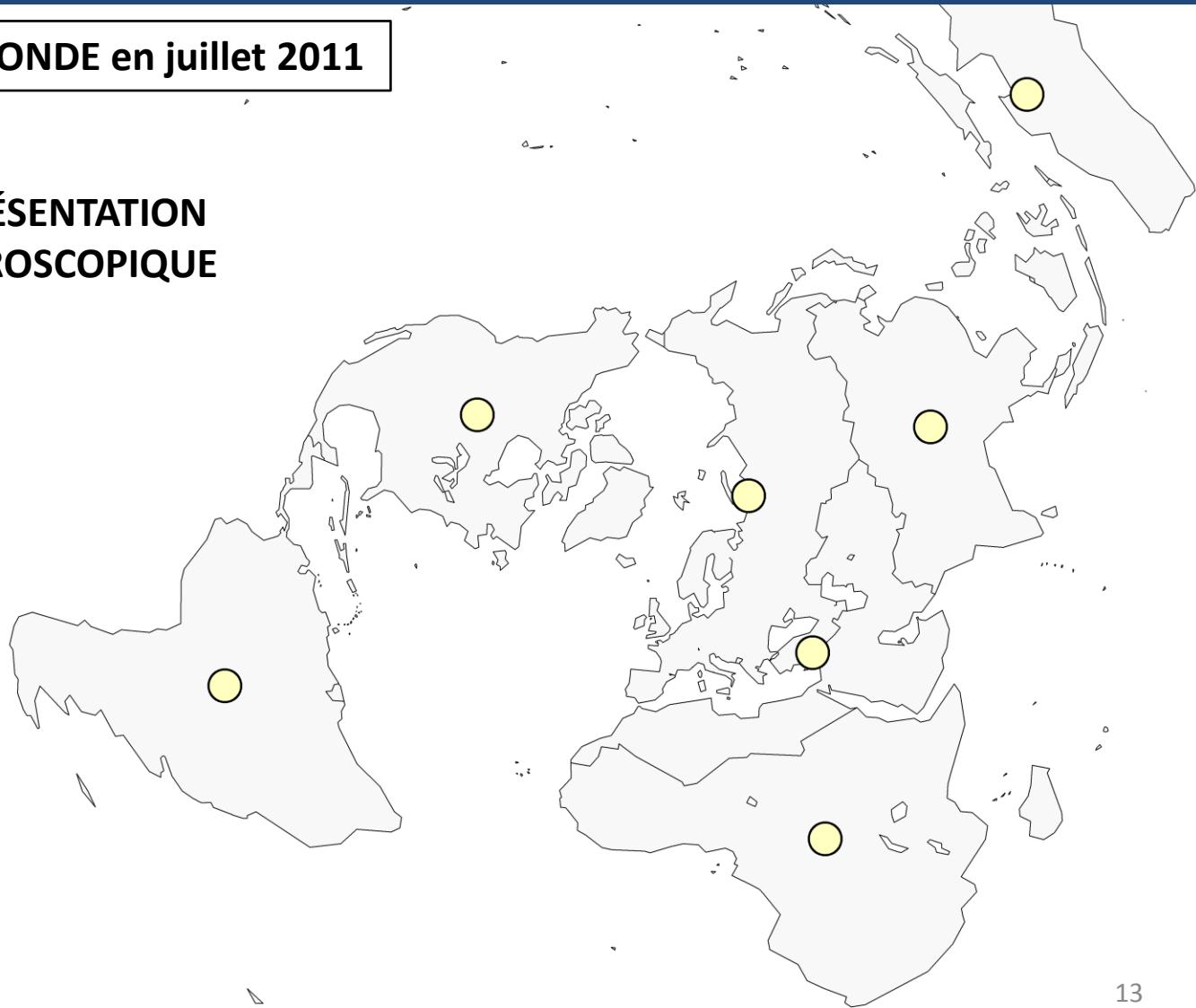


# P2 : Niveaux de représentation

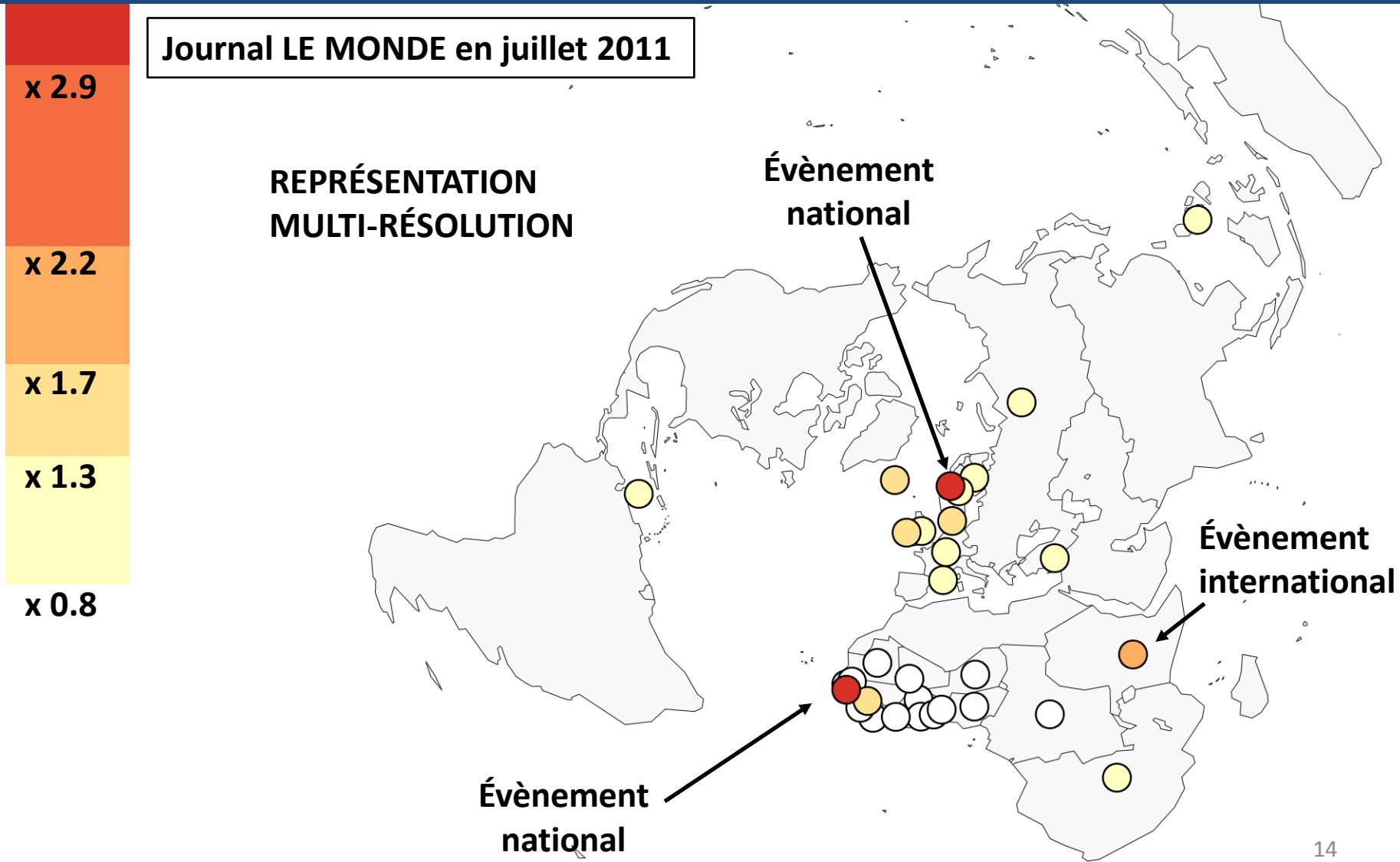


Journal LE MONDE en juillet 2011

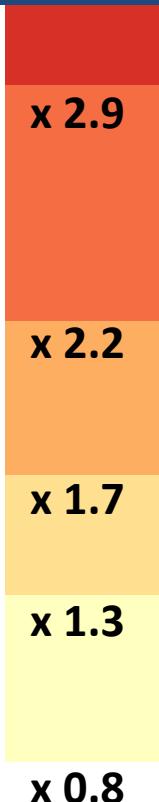
REPRÉSENTATION  
MACROSCOPIQUE



# P2 : Niveaux de représentation



# P2 : Niveaux de représentation



Journal LE MONDE en juillet 2011

REPRÉSENTATION  
MULTI-RÉSOLUTION

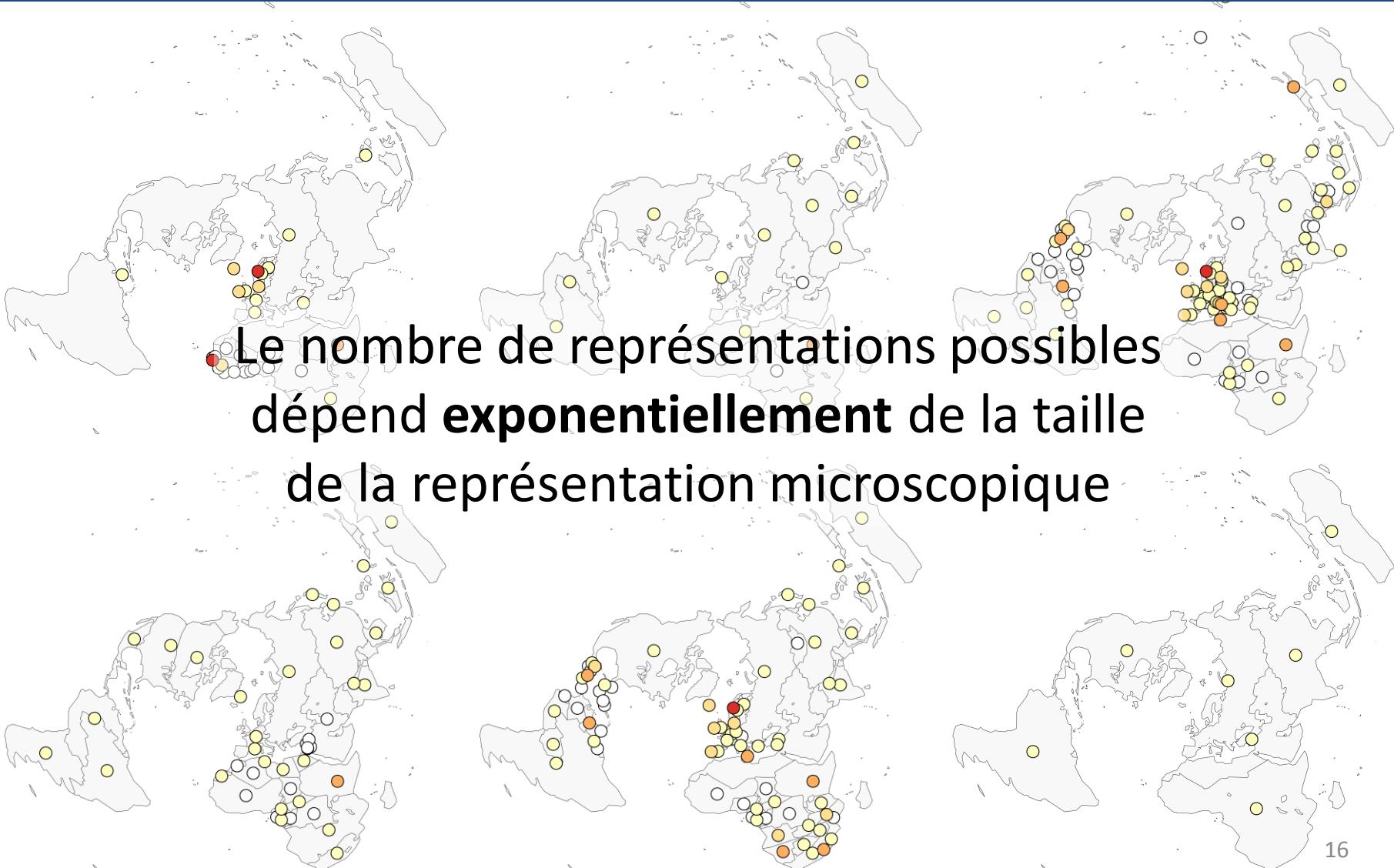
**Problème 2 : Comment trouver  
les niveaux de représentation  
pertinents pour l'analyse ?**

Évènement  
national

Évènement  
international

Évènement  
national

# P3 : Calcul des meilleures représentations

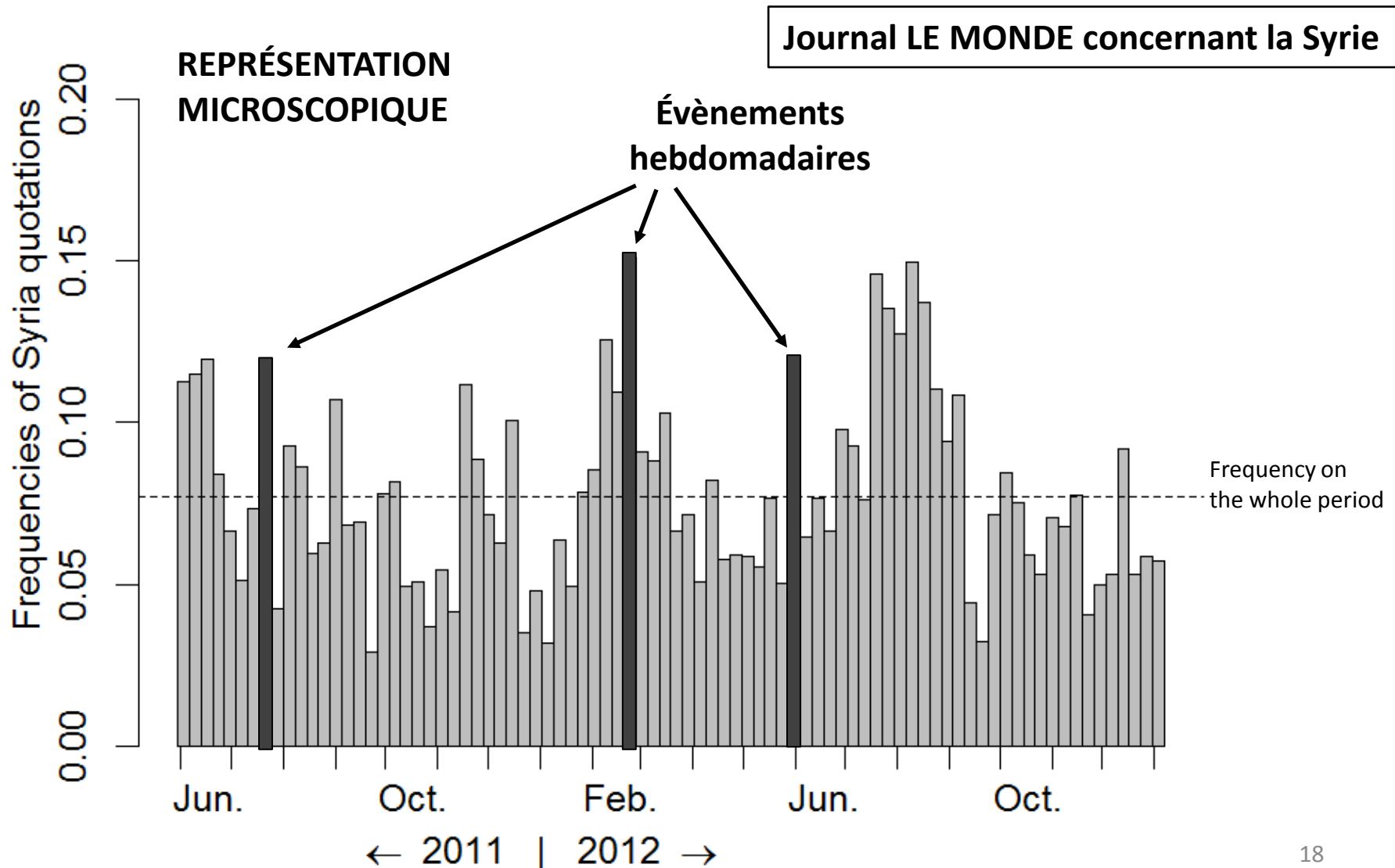


- Le nombre de représentations possibles dépend **exponentiellement** de la taille de la représentation microscopique

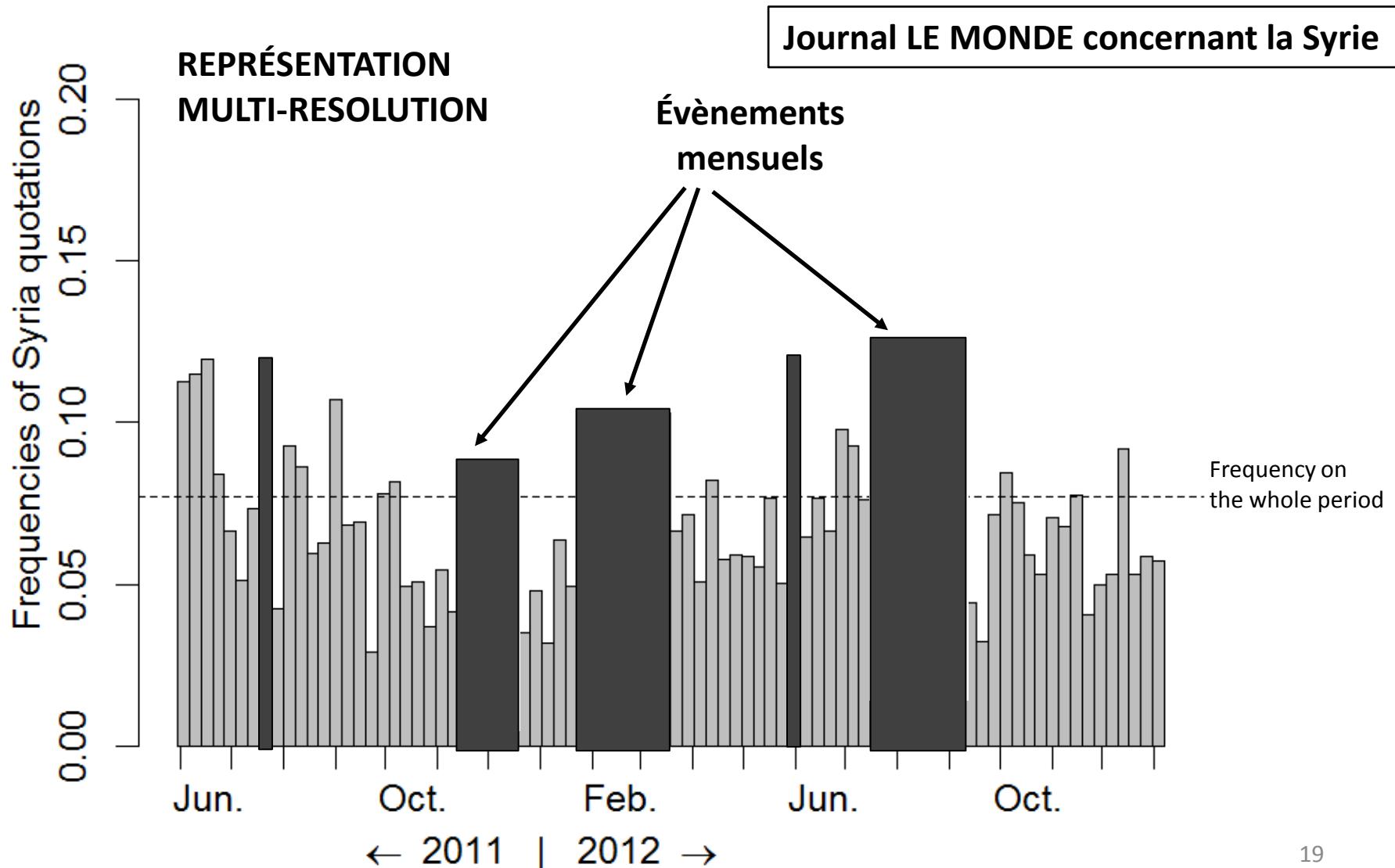
# P3 : Calcul des meilleures représentations

**Problème 3 : Comment calculer les « meilleures » représentations de manière efficace ?**

# Et pour d'autres dimensions du système ?

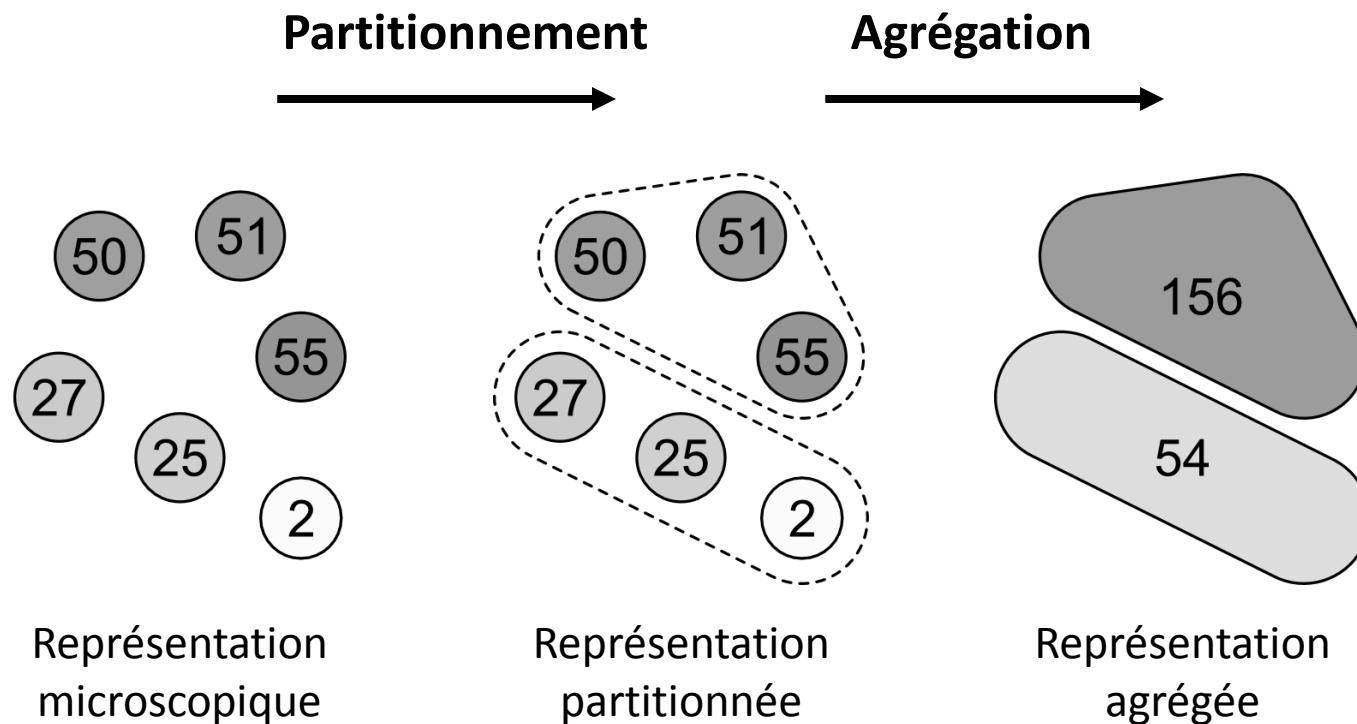


# Et pour d'autres dimensions du système ?

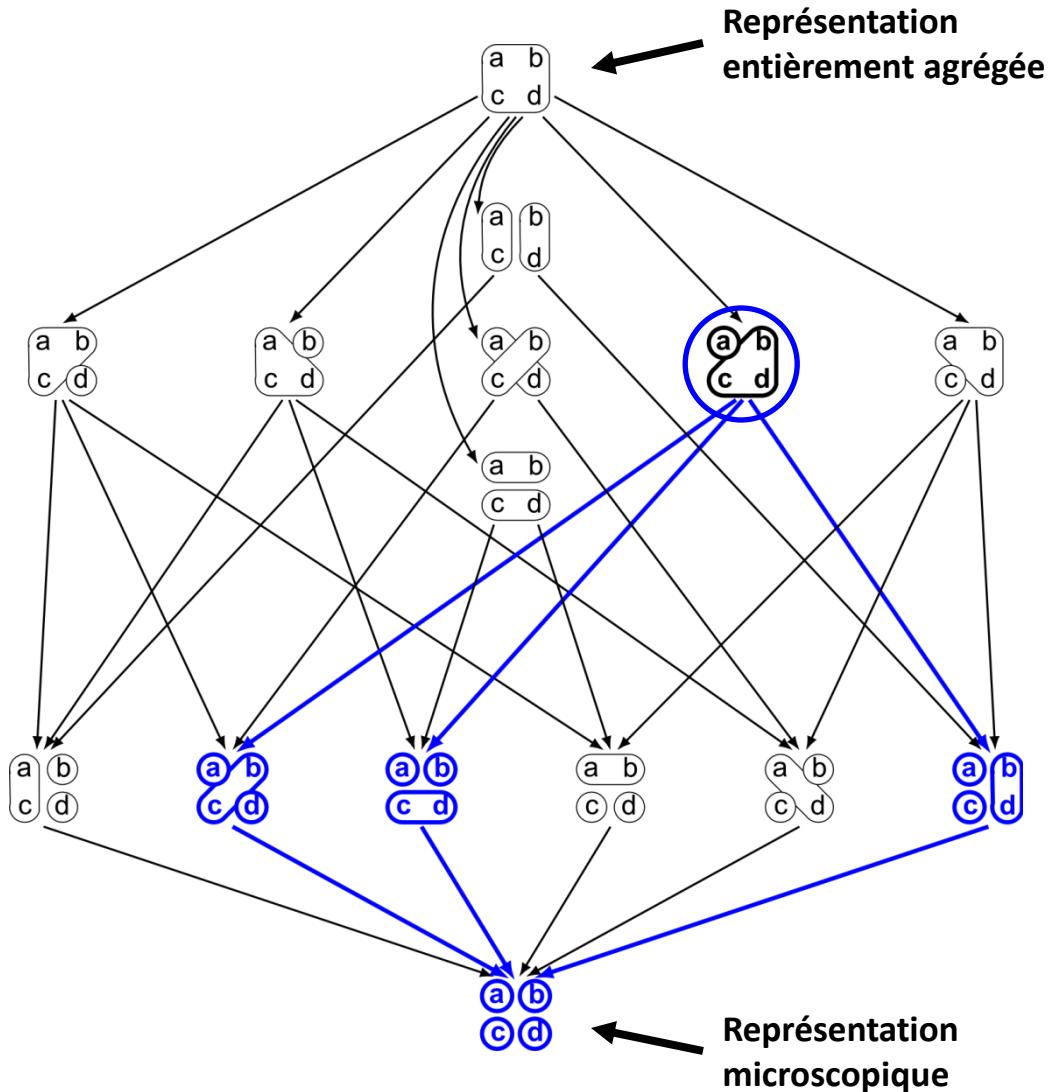


# **PRÉSERVER LA STRUCTURE DE L'ESPACE GÉOGRAPHIQUE**

# Le processus d'agrégation



# Ensemble des partitions possibles



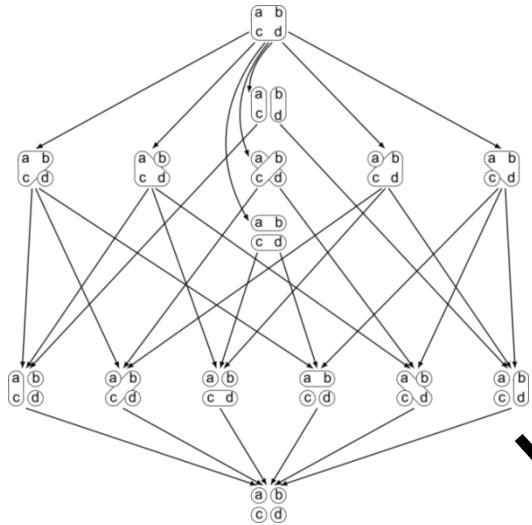
## Structure algébrique

Ordre partiel sur l'ensemble des partitions possibles

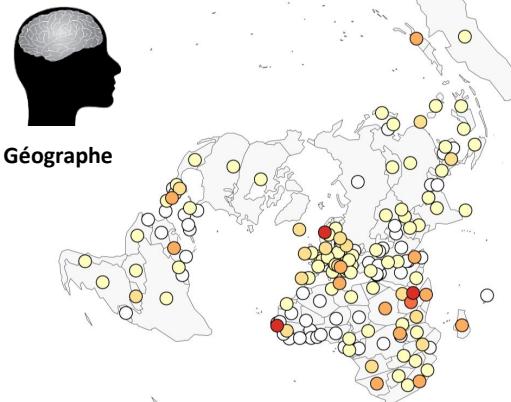
→ relations de raffinement

# Problème et objectif

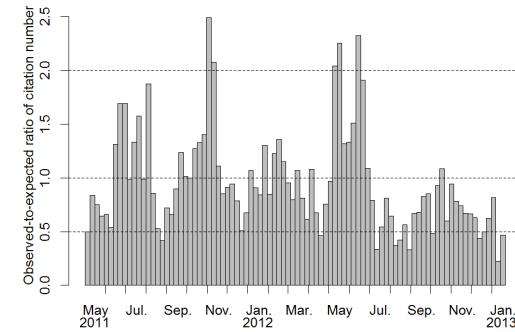
## Ensemble des partitions possibles



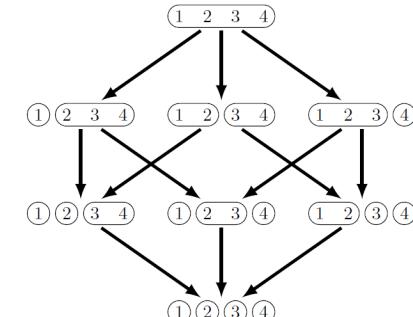
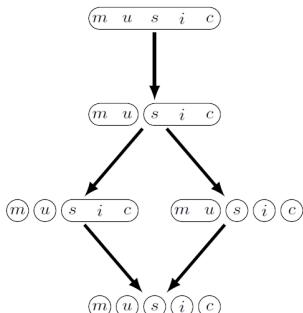
## Sémantique géographique



## Sémantique temporelle

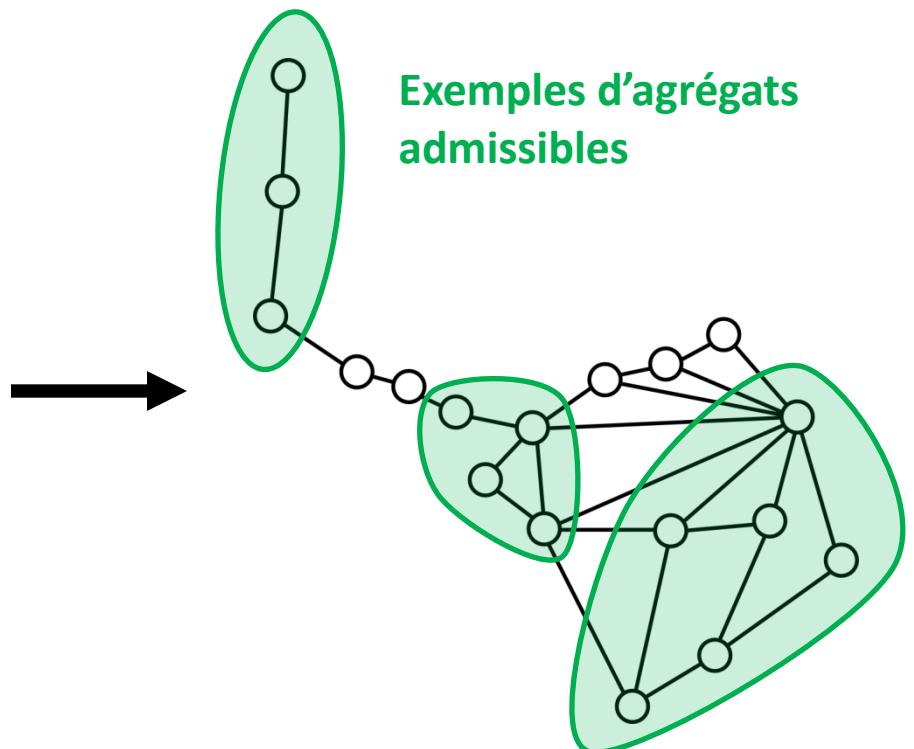


**Construire l'ensemble  
des partitions admissibles**



# Conserver la relation de voisinage

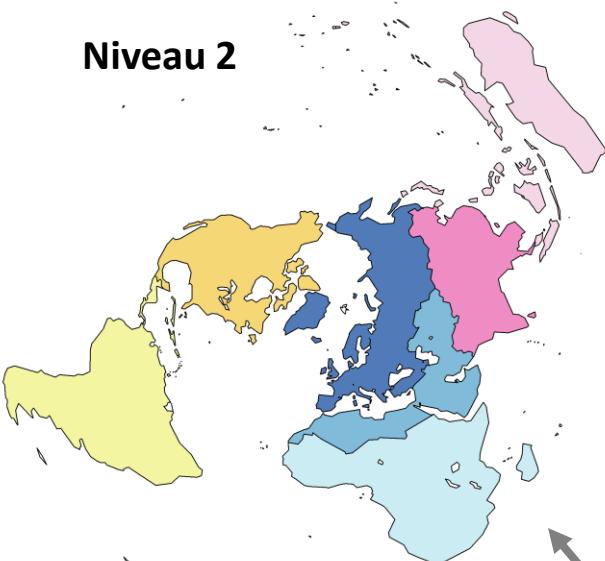
**Agrégats admissibles** : ensembles de pays connexes vis-à-vis du graphe de voisinage



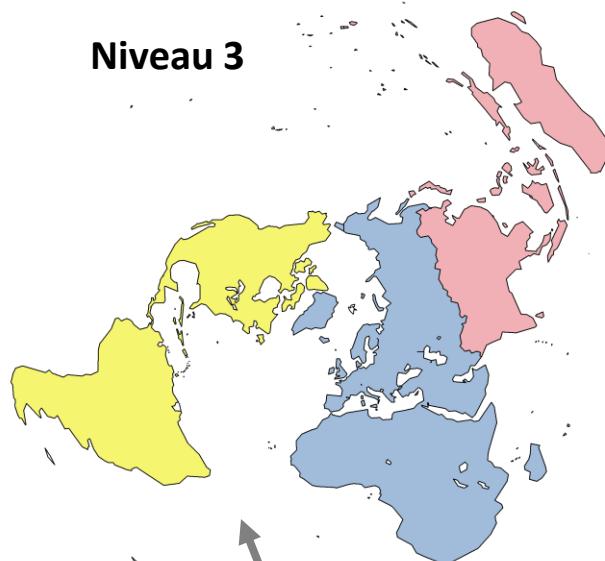
# La hiérarchie WUTS

[Grasland et Didelon, 2007]

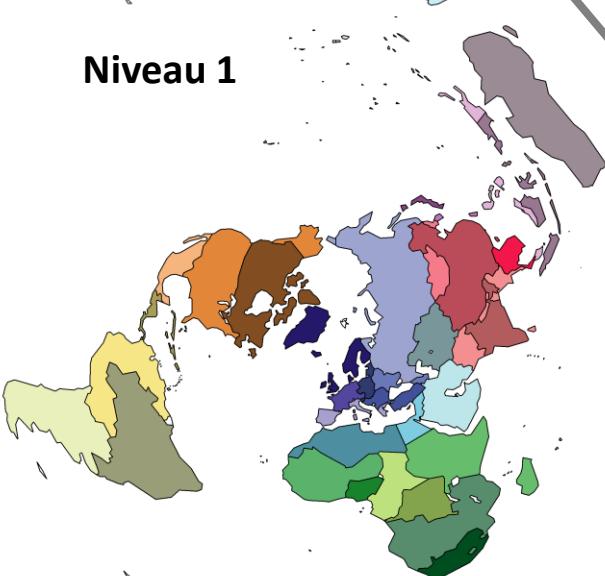
Niveau 2



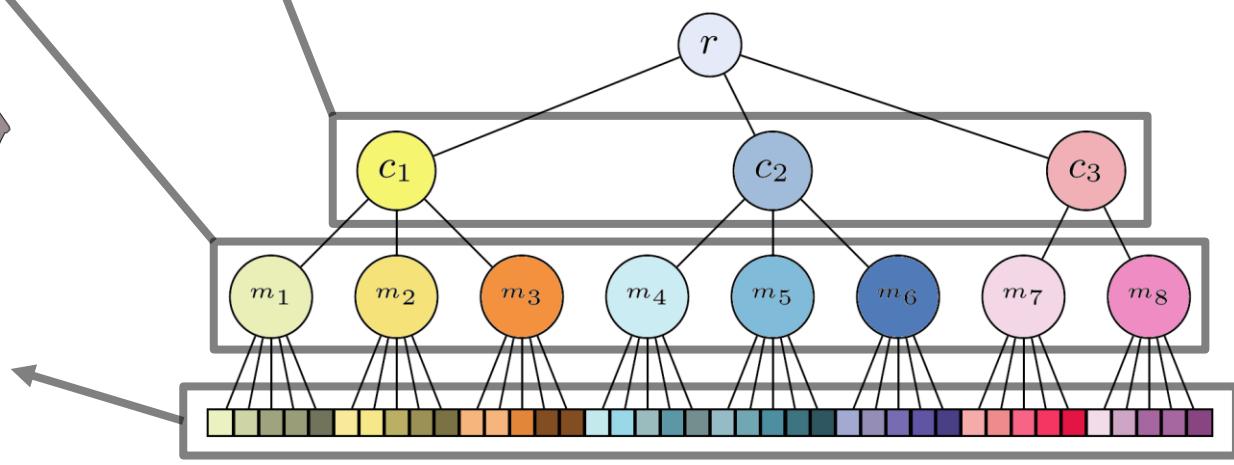
Niveau 3



Niveau 1



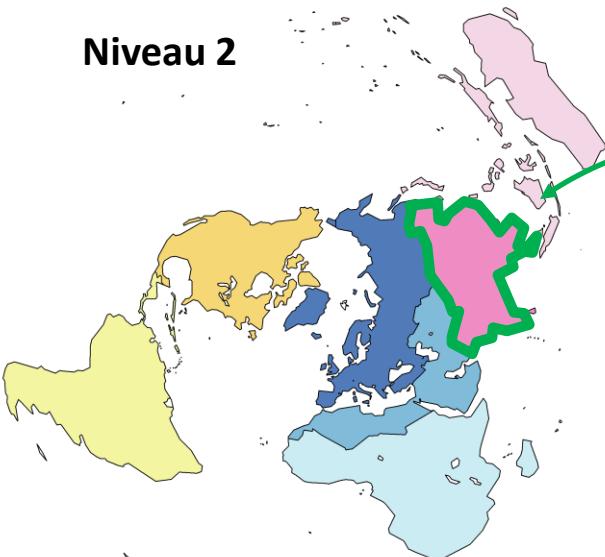
**Agrégats admissibles :**  
ensembles de pays proches  
sur le plan politique,  
culturel, économique, etc.



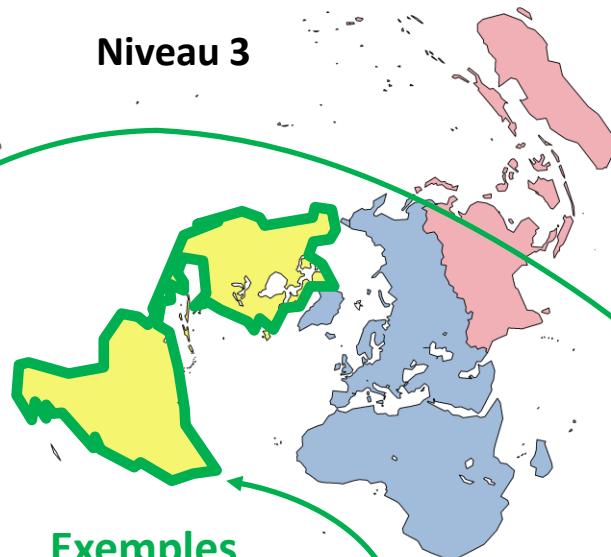
# La hiérarchie WUTS

[Grasland et Didelon, 2007]

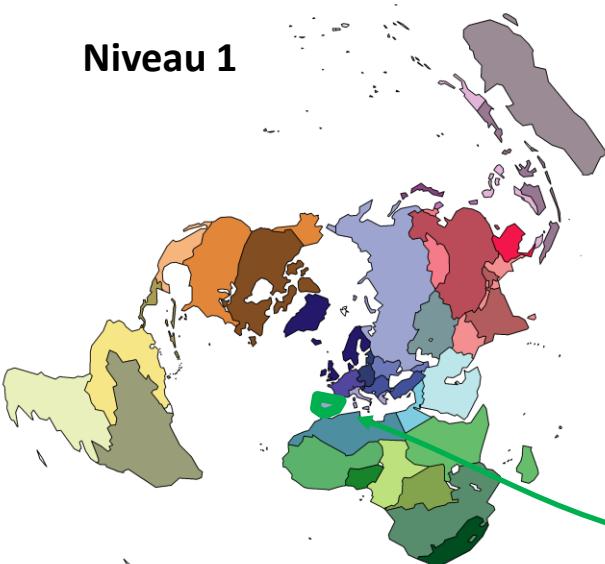
Niveau 2



Niveau 3

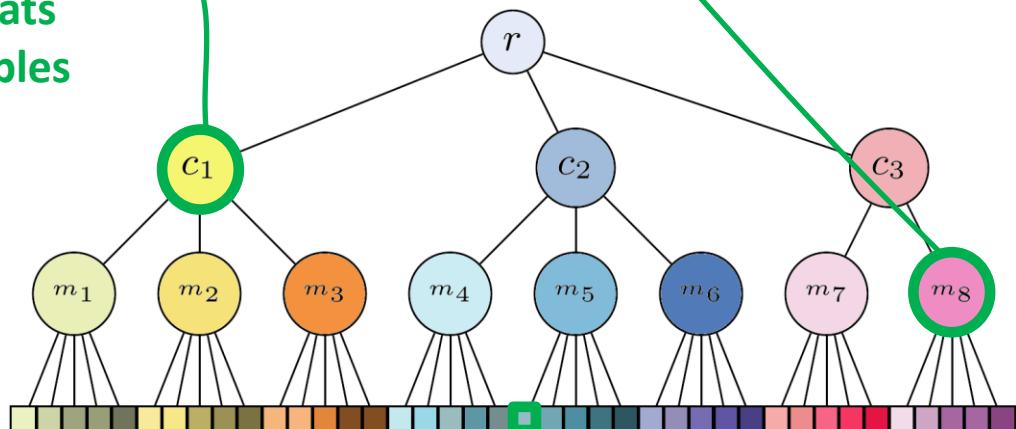


Niveau 1



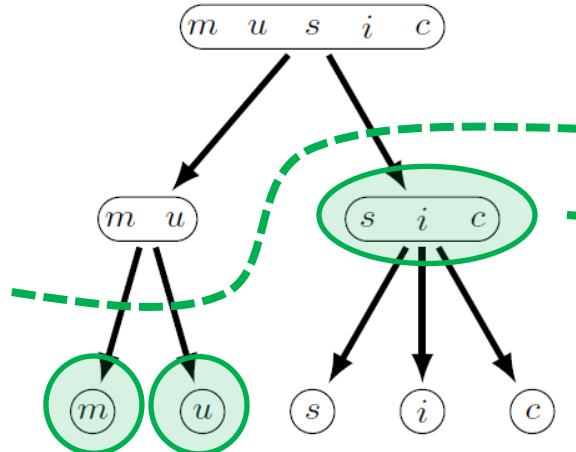
Exemples  
d'agrégats  
admissibles

**Agrégats admissibles :**  
ensembles de pays proches  
sur le plan politique,  
culturel, économique, etc.

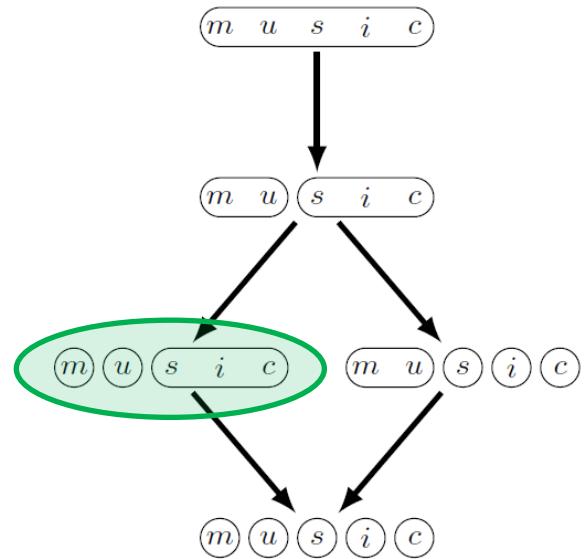


# Agrégation selon une hiérarchie

**Agrégats admissibles**  
(nœuds de la hiérarchie)

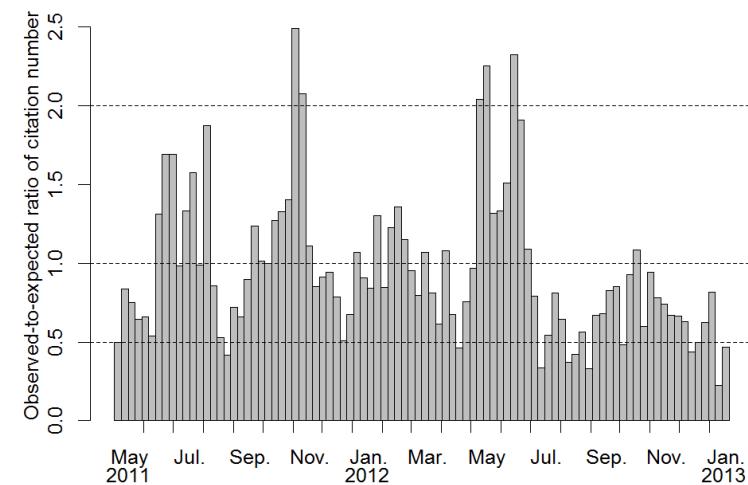


**Partitions admissibles**  
(coupe dans la hiérarchie)

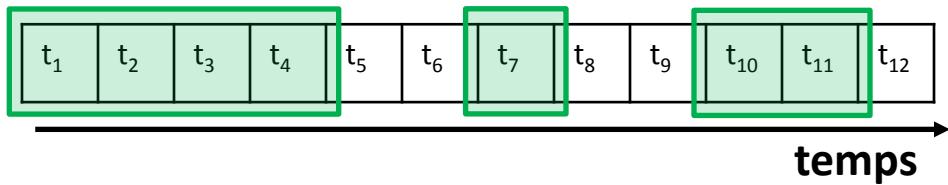


# Conserver l'ordre du temps

Agrégats admissibles :  
intervalles de temps

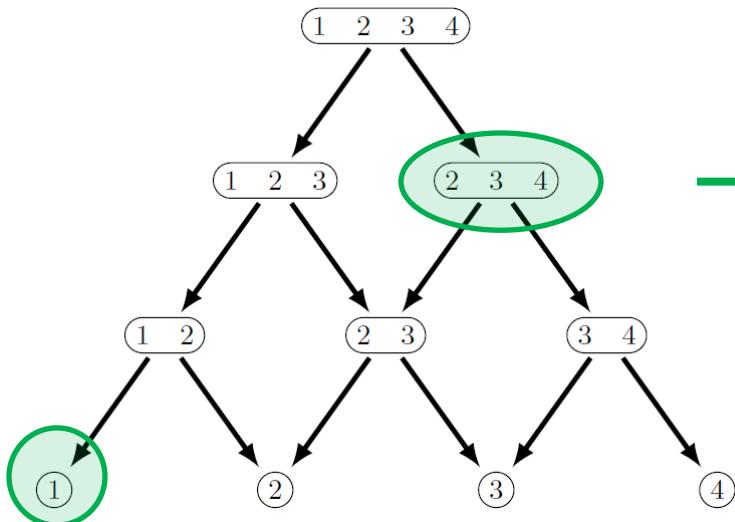


Exemples d'agrégats admissibles

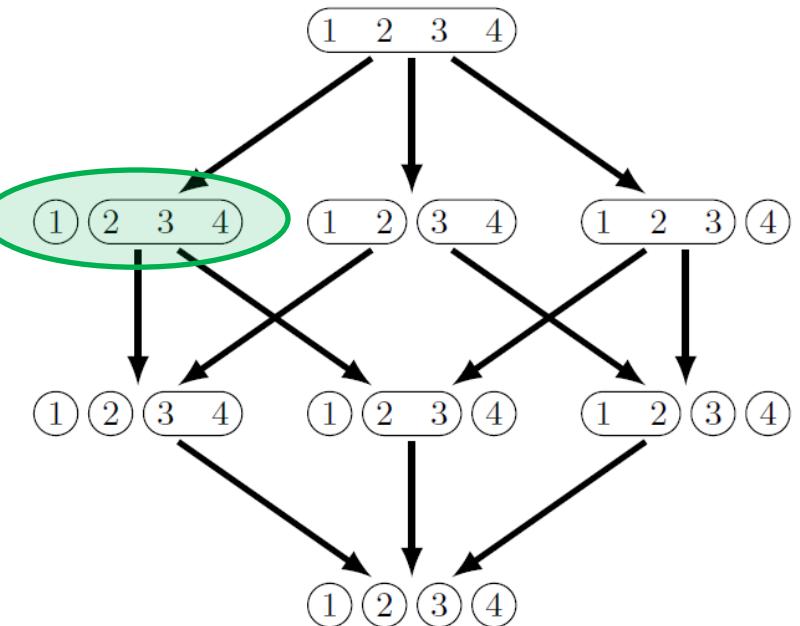


# Agrégation selon un ordre total

**Agrégats admissibles**  
(intervalles de temps)

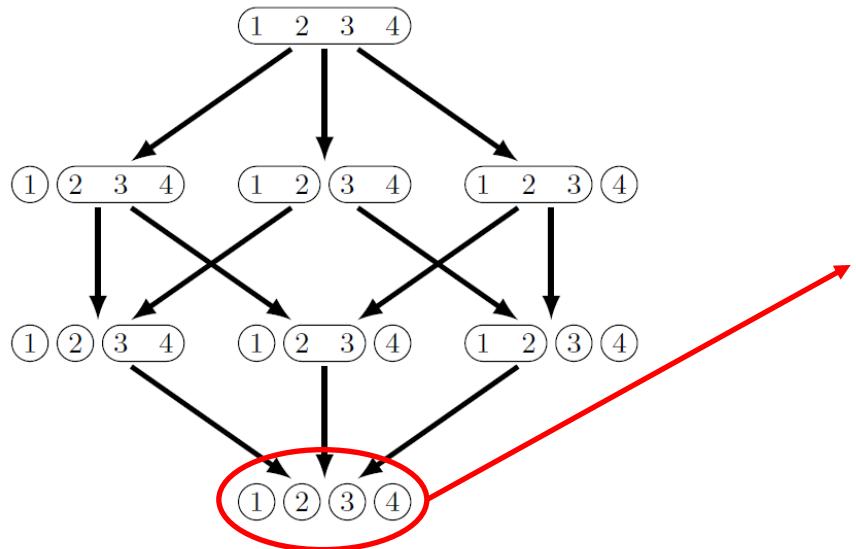


**Partitions admissibles**  
(séquences d'intervalles)



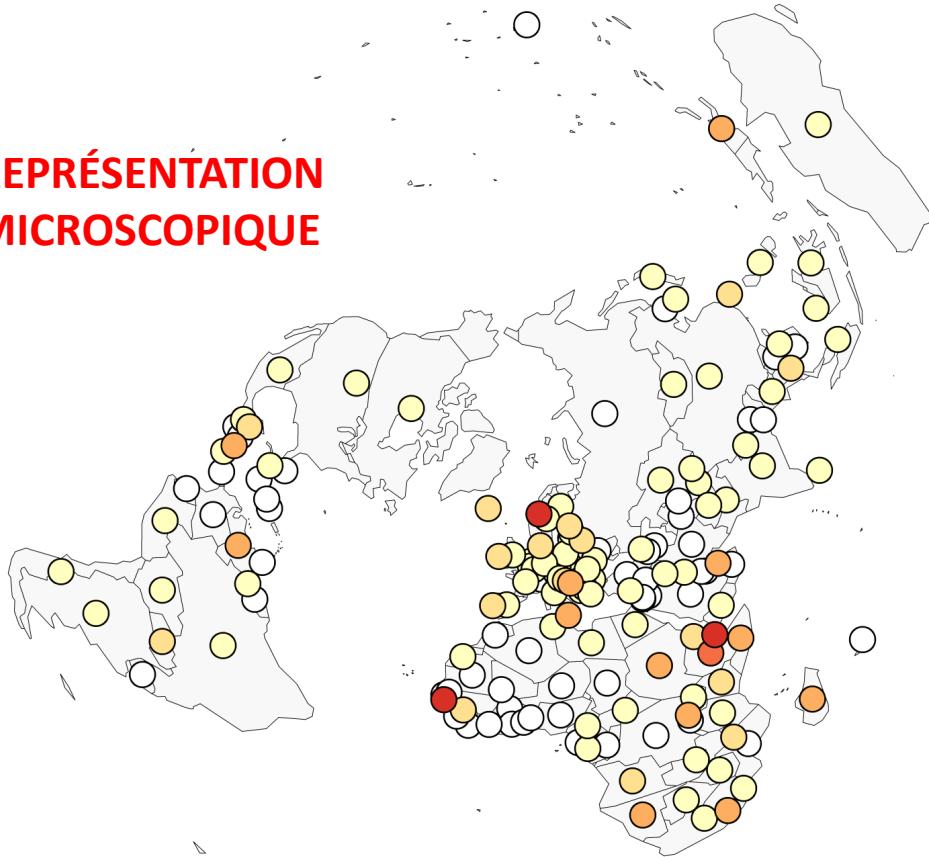
# **AGRÉGATION ET INFORMATION**

# Objectifs et difficultés



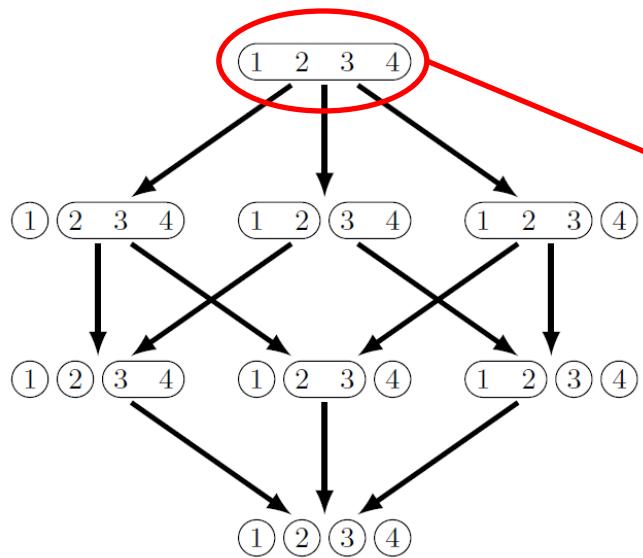
REPRÉSENTATION  
MICROSCOPIQUE

Quelle partition admissible  
est la meilleure pour un jeu  
de données particulier ?



→ TROP COMPLEXE POUR  
PASSER À L'ÉCHELLE

# Objectifs et difficultés



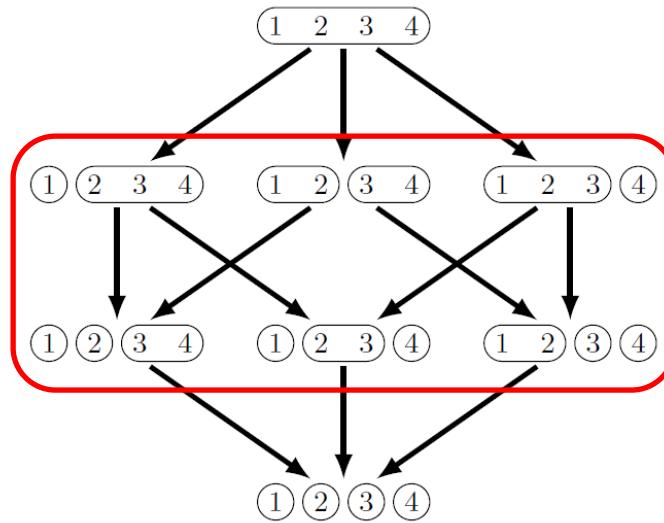
REPRÉSENTATION  
ENTIÈREMENT AGRÉGÉE



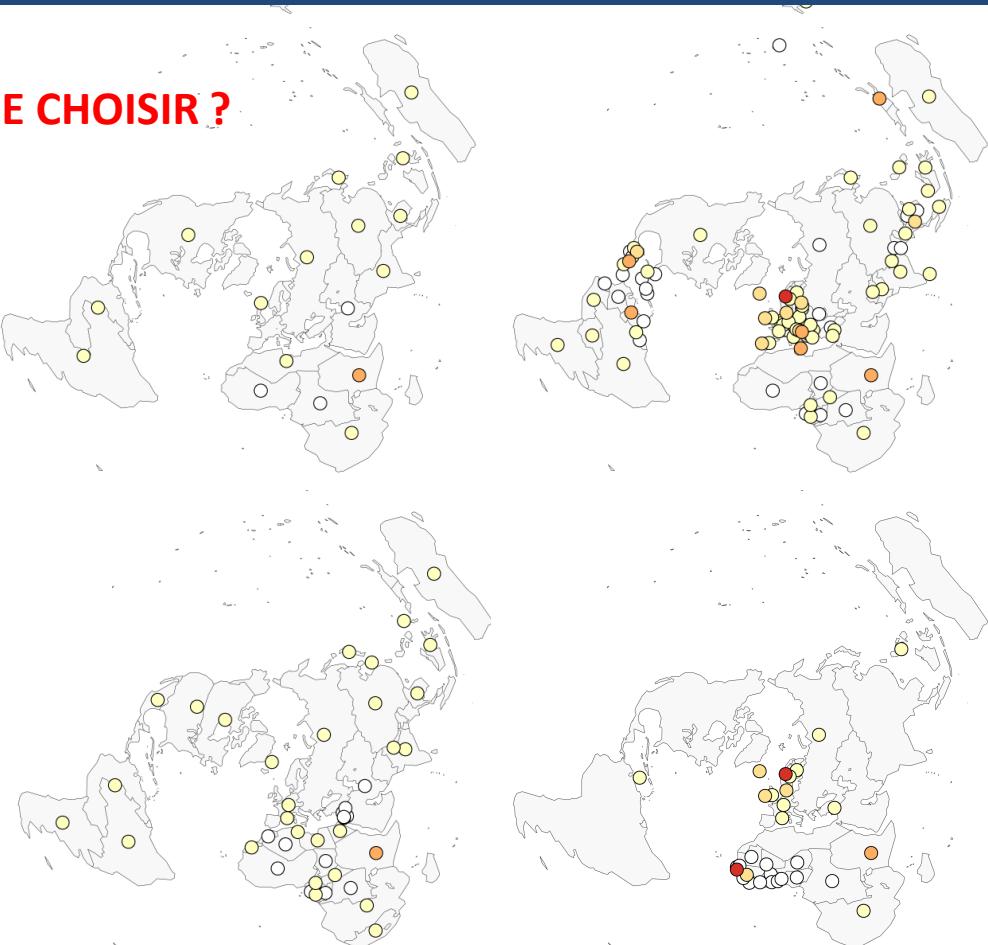
Quelle partition admissible  
est la meilleure pour un jeu  
de données particulier ?

→ NE DONNE QUE TRÈS  
PEU D'INFORMATION

# Objectifs et difficultés

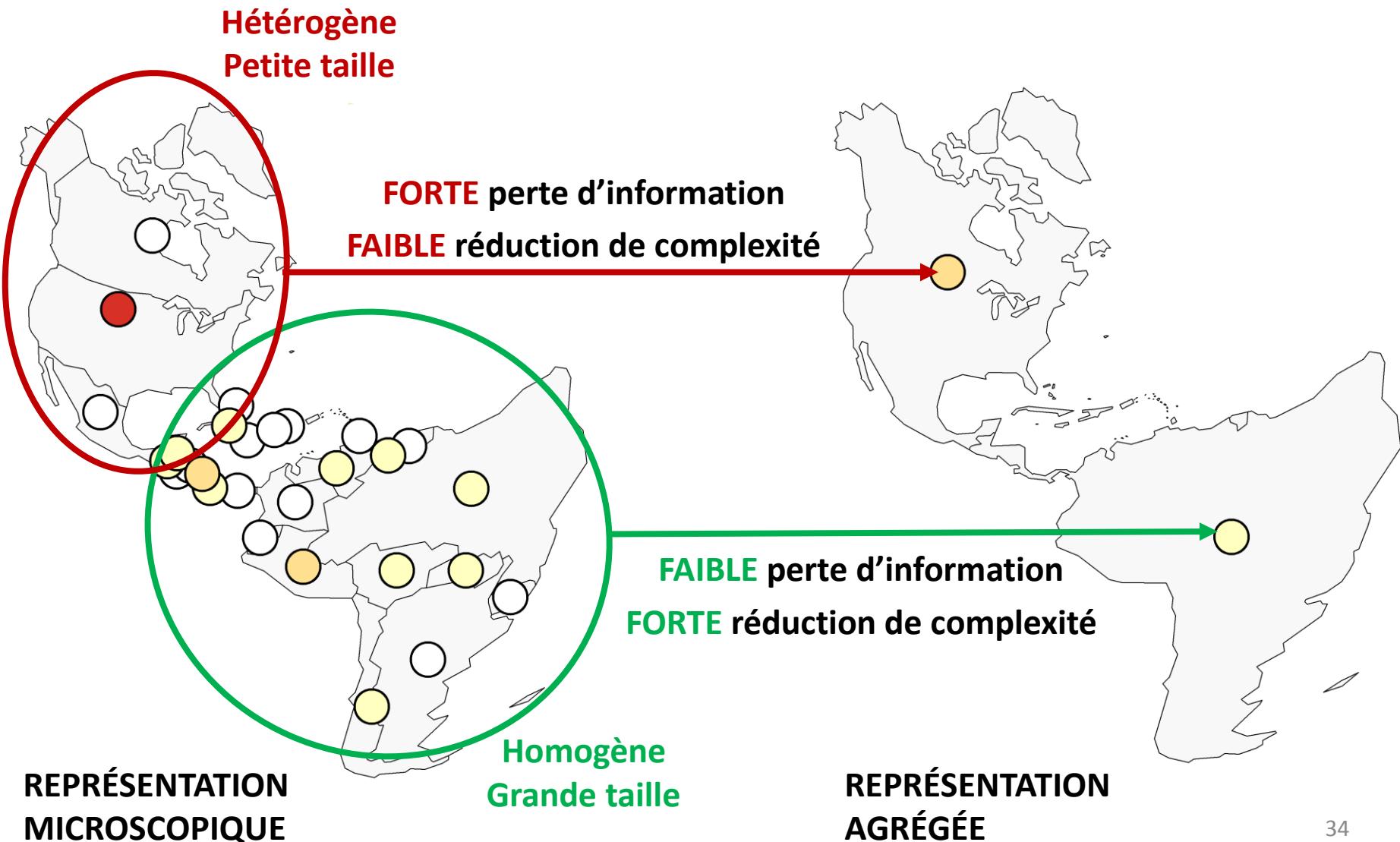


QUE CHOISIR ?

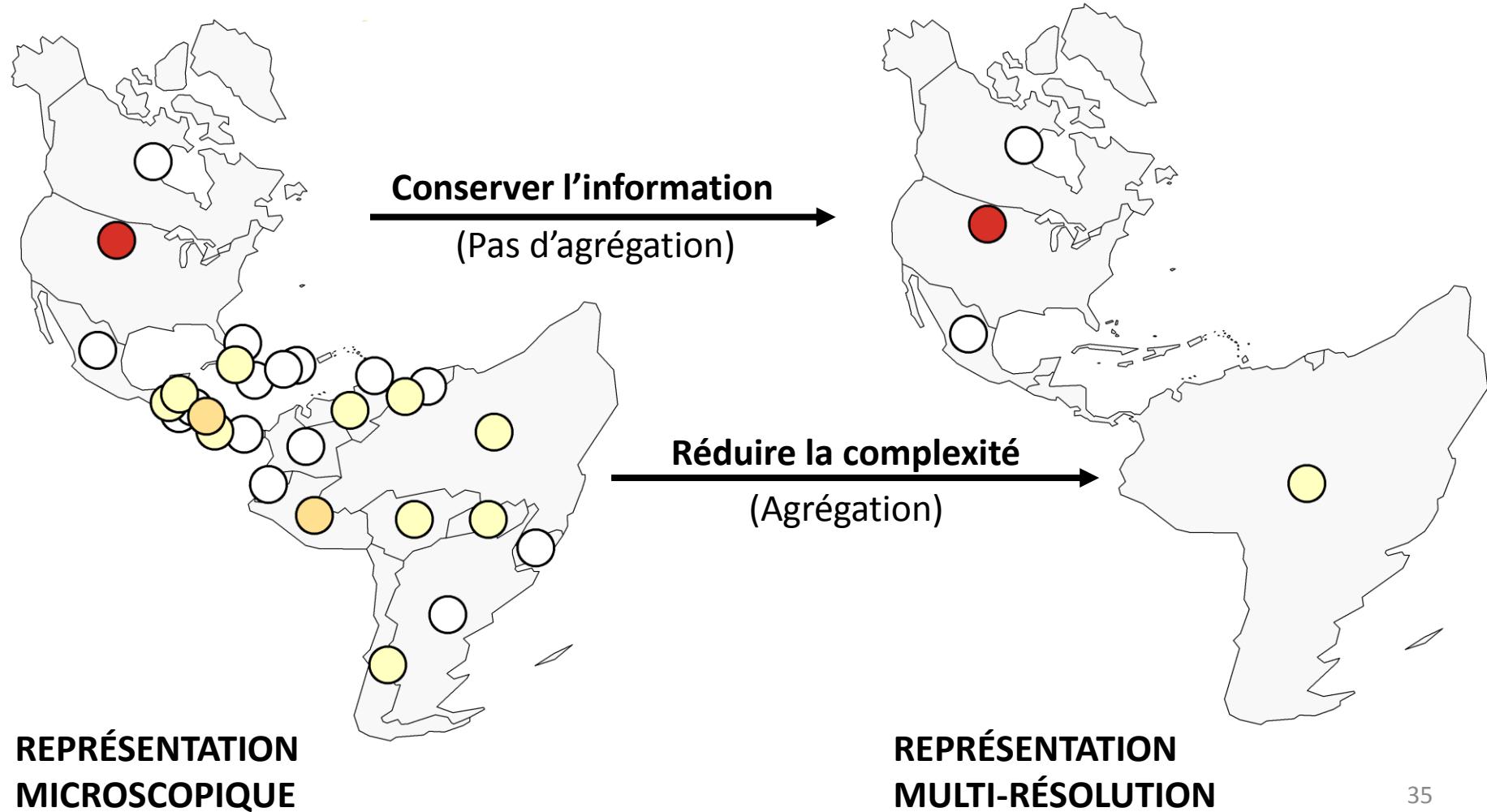


Quelle partition admissible  
est la meilleure pour un jeu  
de données particulier ?

# Complexité et information



# Complexité et information



# Mesures de qualité

[Lamarche-Perrin et al., ECCS 2012]

La **complexité** dépend de la **tâche à accomplir** et des **outils de description** disponibles

[Bonabeau et Dessalles, 1997]

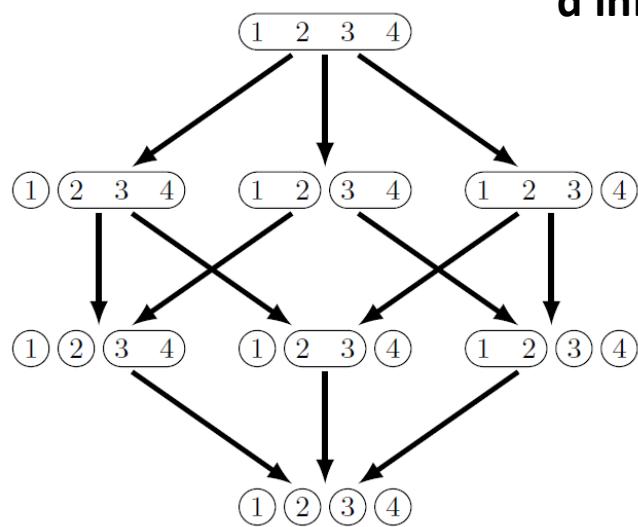
La **perte d'information** est mesurée par la **divergence** entre deux distributions de probabilité

[Kullback et Leibler, 1951]

Nombre d'agrégats représentés :

$$T(\mathcal{X}) = |\mathcal{X}|$$

Complexité



Perte  
d'information

Divergence de Kullback-Leibler :

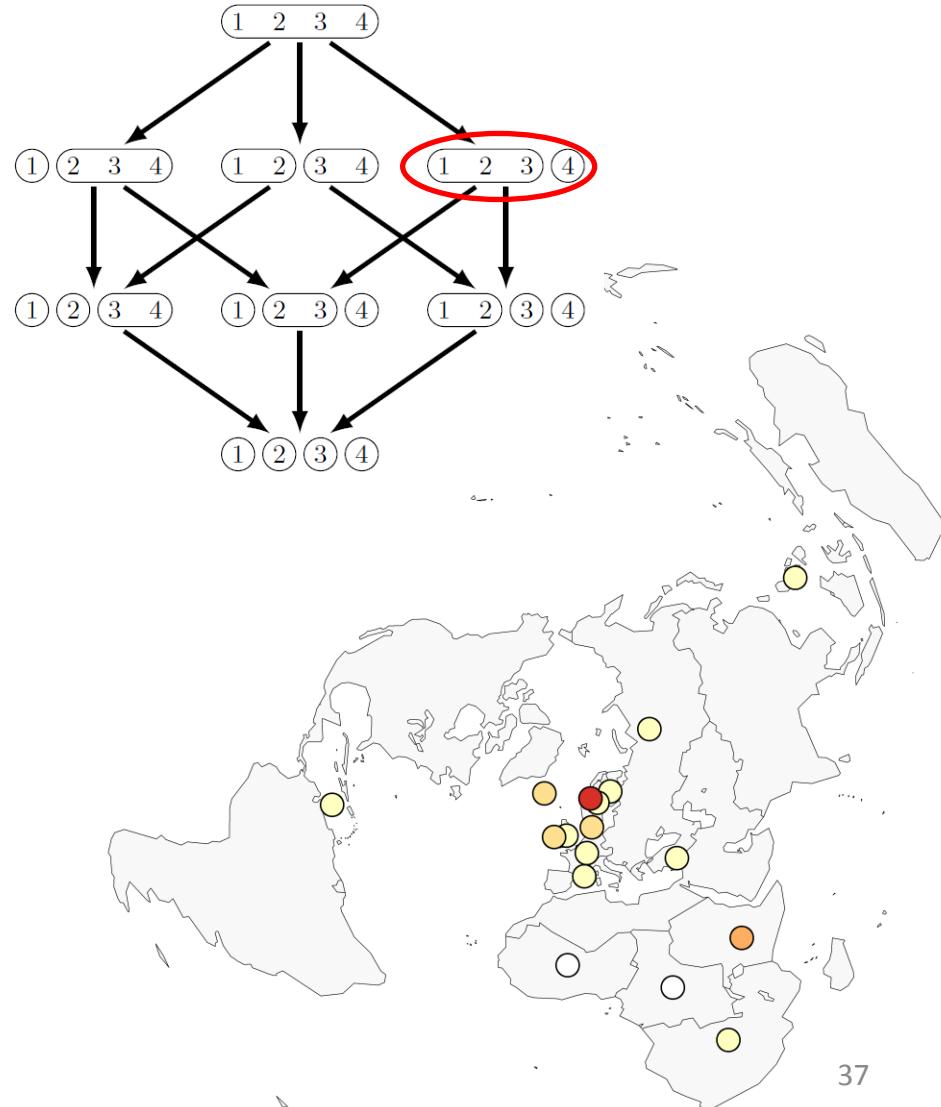
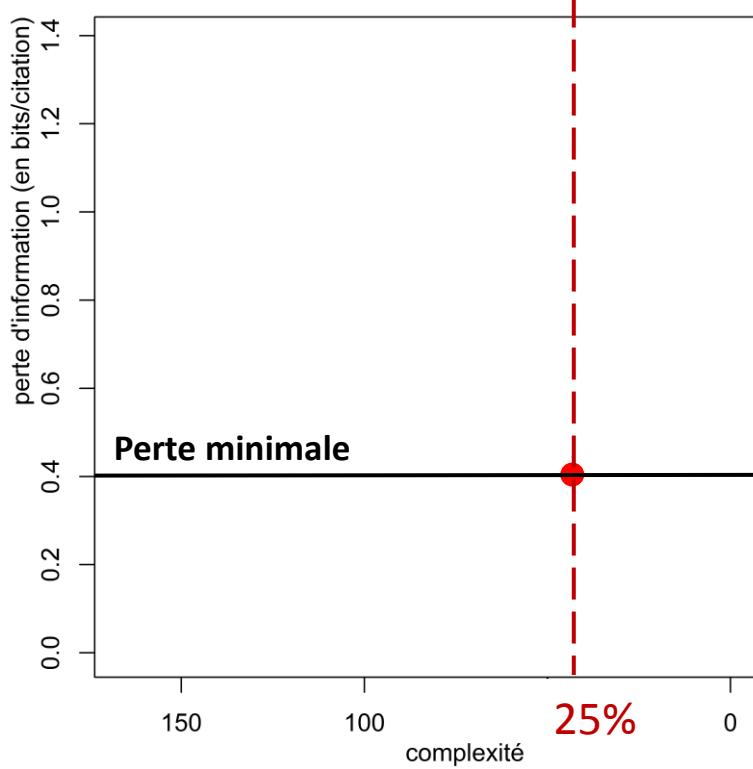
$$D(\mathcal{X}) = \sum_{X \in \mathcal{X}} \sum_{x \in X} \frac{v(x)}{v(\Omega)} \log_2 \left( \frac{v(x) |X|}{v(X)} \right)$$

# Optimisation des mesures de qualité

Deux critères d'évaluation indépendants

Compromis de qualité :

$$CQL_\alpha = \alpha \frac{\Delta T}{\Delta T_{\max}} - (1 - \alpha) \frac{D}{D_{\max}}$$

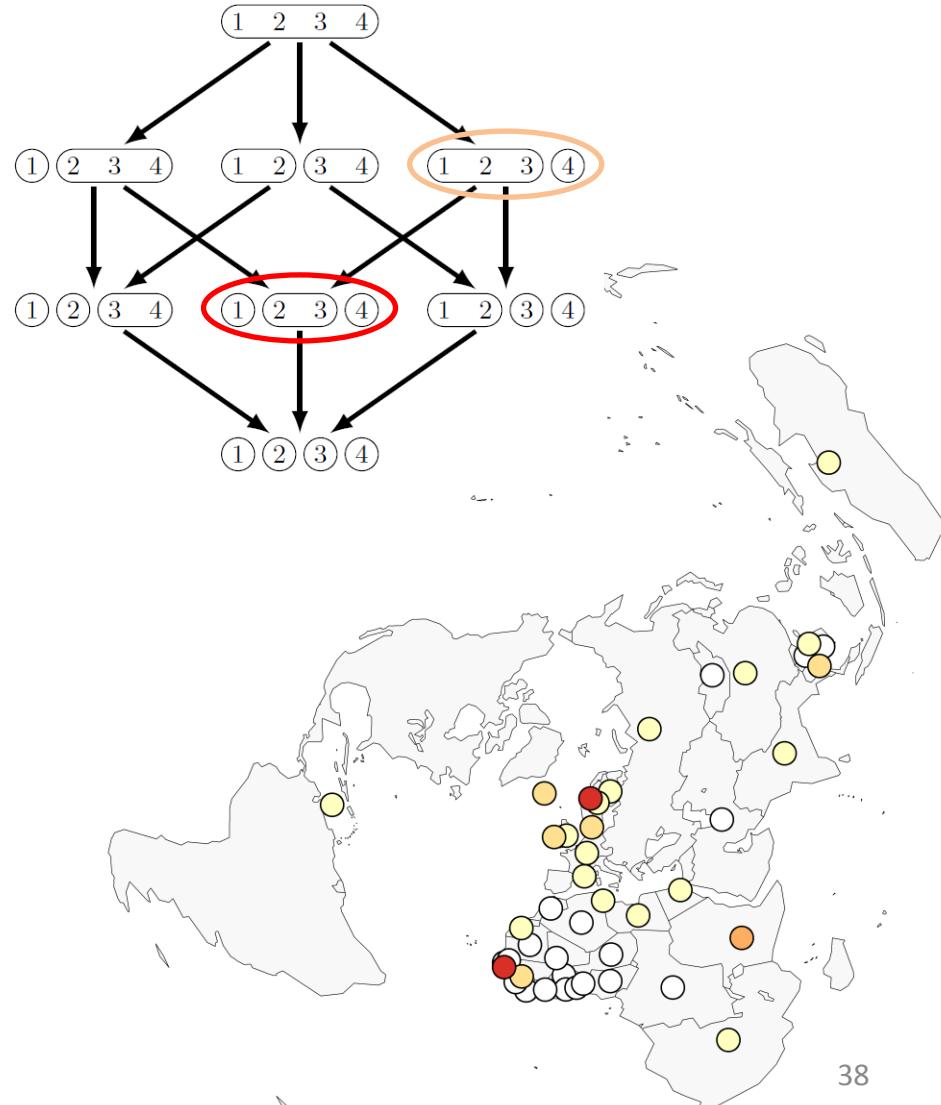
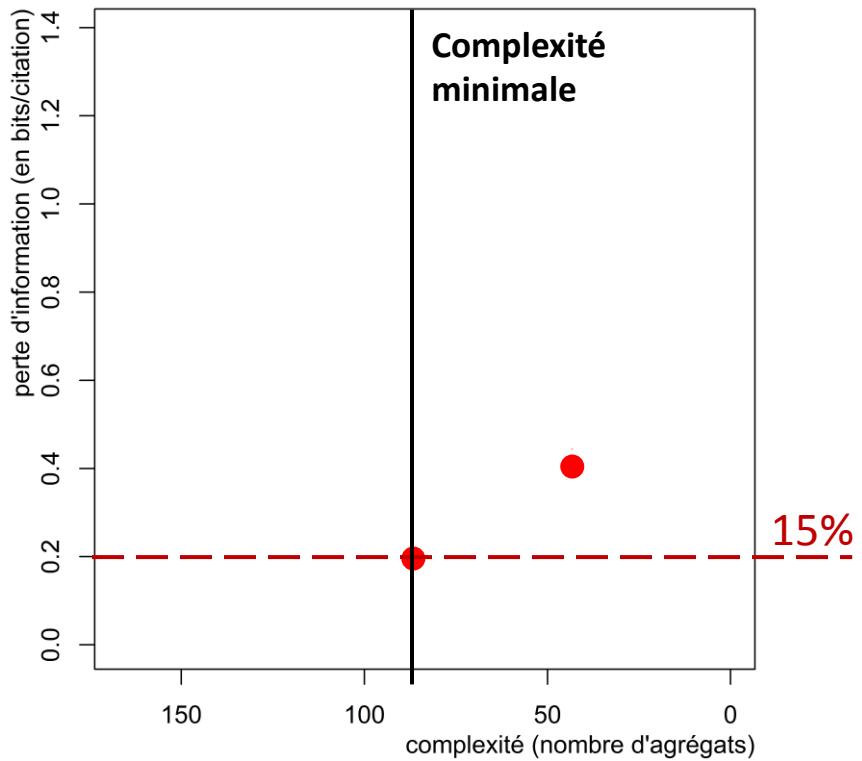


# Optimisation des mesures de qualité

Deux critères d'évaluation indépendants

Compromis de qualité :

$$CQL_\alpha = \alpha \frac{\Delta T}{\Delta T_{\max}} - (1 - \alpha) \frac{D}{D_{\max}}$$

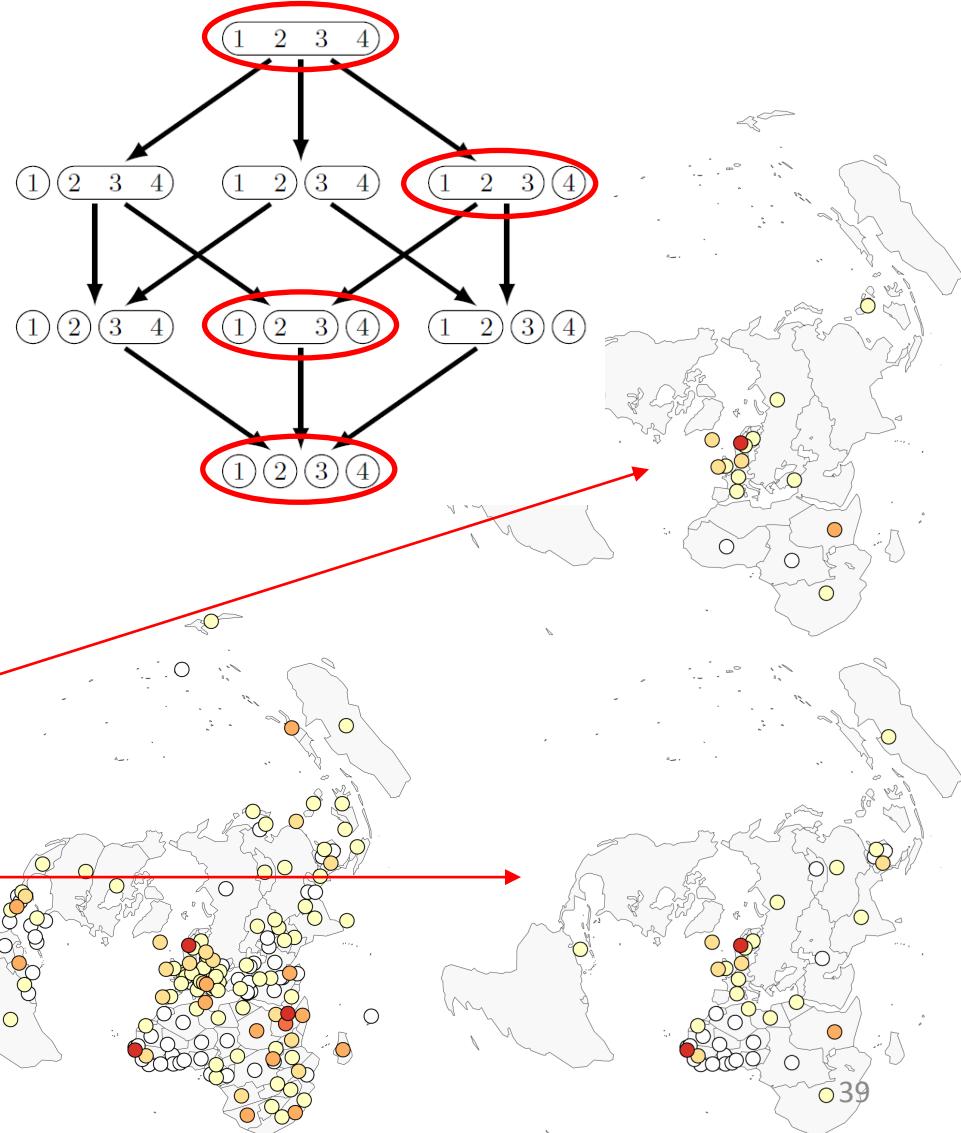
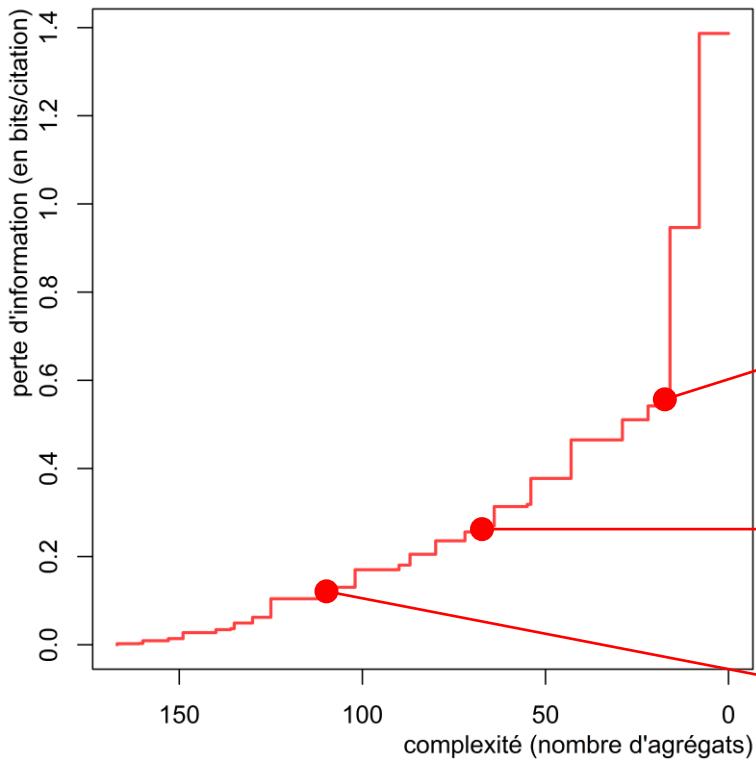


# Optimisation des mesures de qualité

Deux critères d'évaluation indépendants

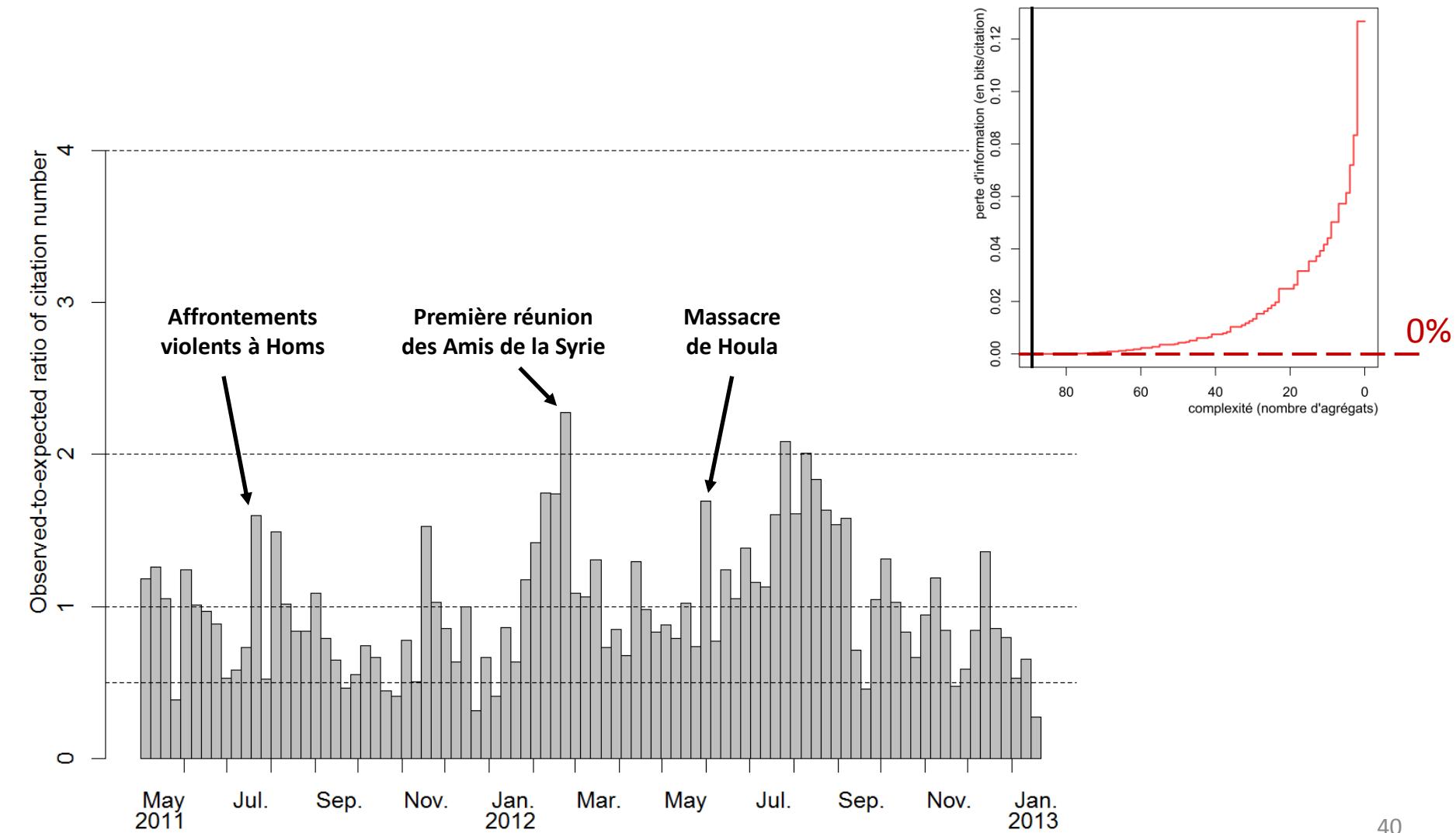
Compromis de qualité :

$$CQL_\alpha = \alpha \frac{\Delta T}{\Delta T_{\max}} - (1 - \alpha) \frac{D}{D_{\max}}$$



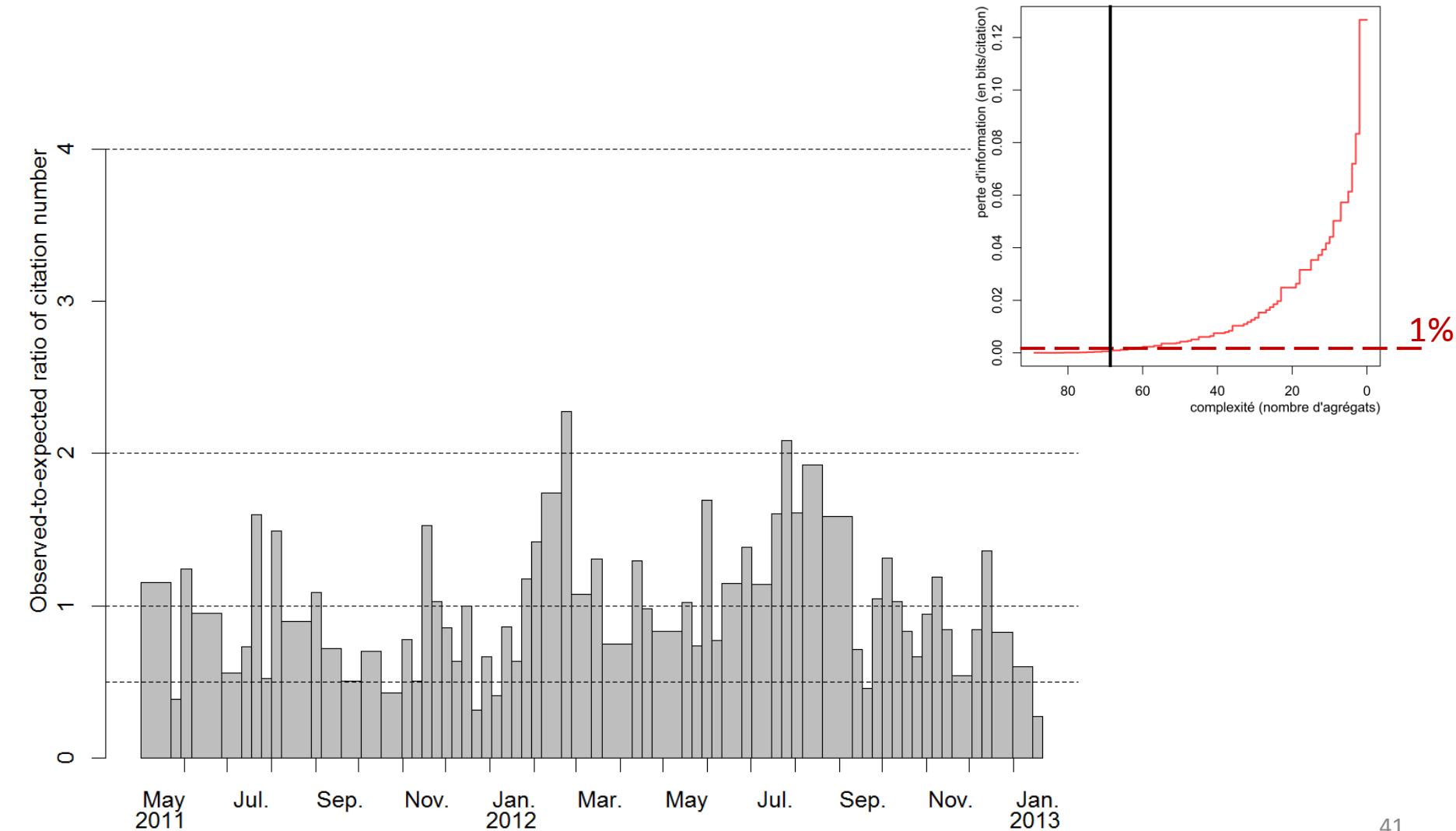
# La Syrie vue par LE MONDE

[Giraud, Grasland, Lamarche-Perrin *et al.*, ECTQG 2013]



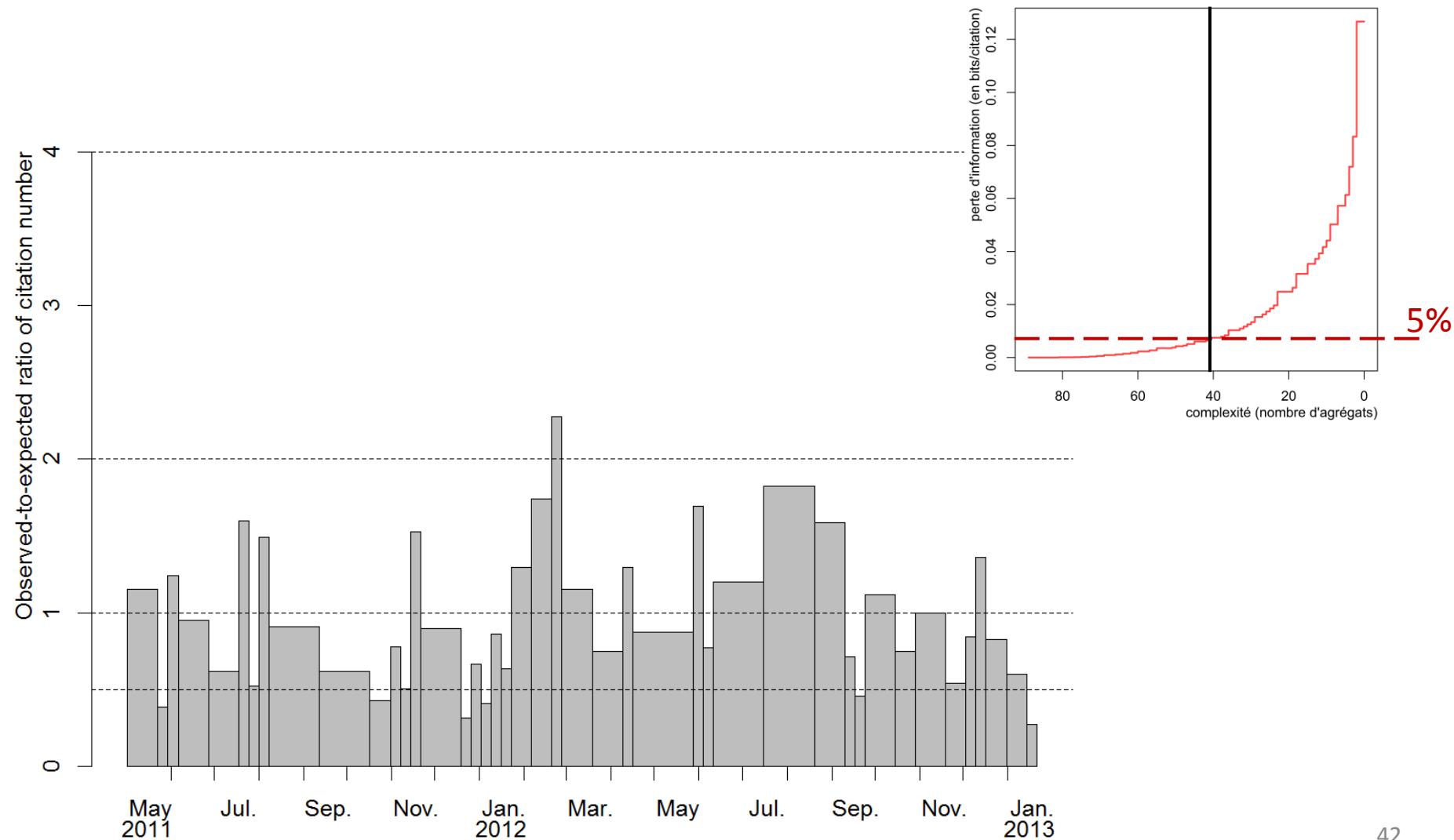
# La Syrie vue par LE MONDE

[Giraud, Grasland, Lamarche-Perrin *et al.*, ECTQG 2013]



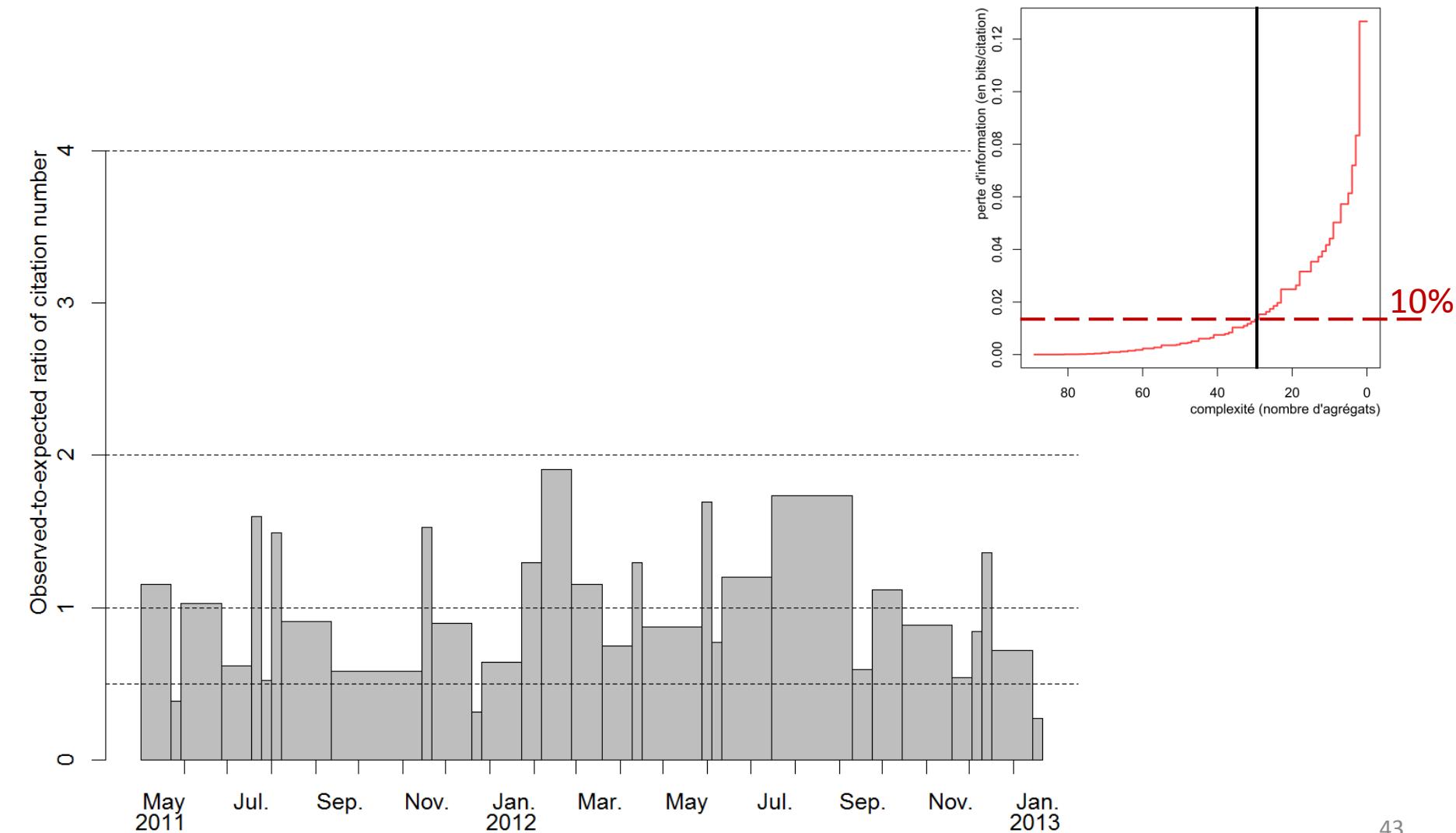
# La Syrie vue par LE MONDE

[Giraud, Grasland, Lamarche-Perrin *et al.*, ECTQG 2013]



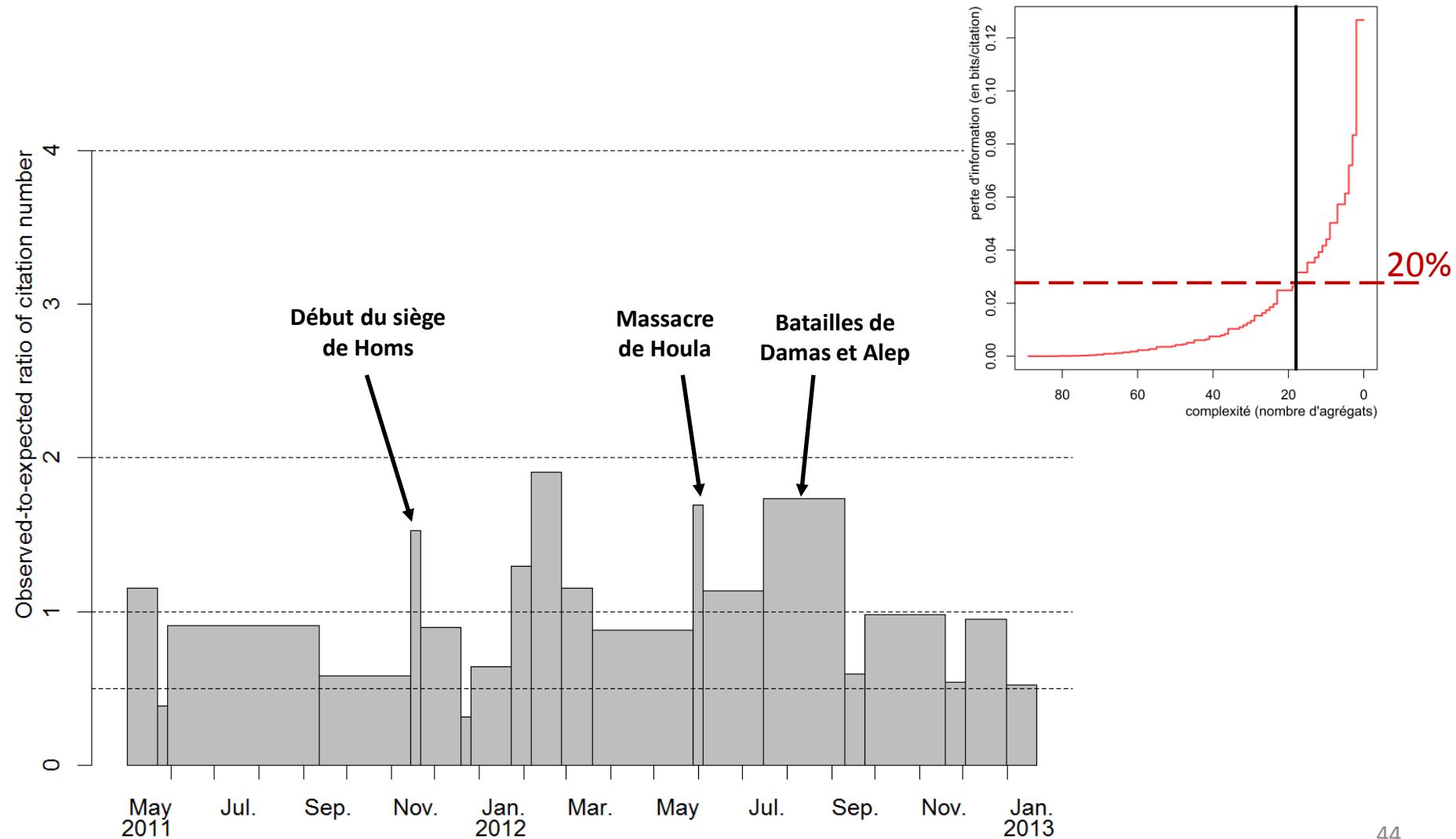
# La Syrie vue par LE MONDE

[Giraud, Grasland, Lamarche-Perrin *et al.*, ECTQG 2013]



# La Syrie vue par LE MONDE

[Giraud, Grasland, Lamarche-Perrin *et al.*, ECTQG 2013]

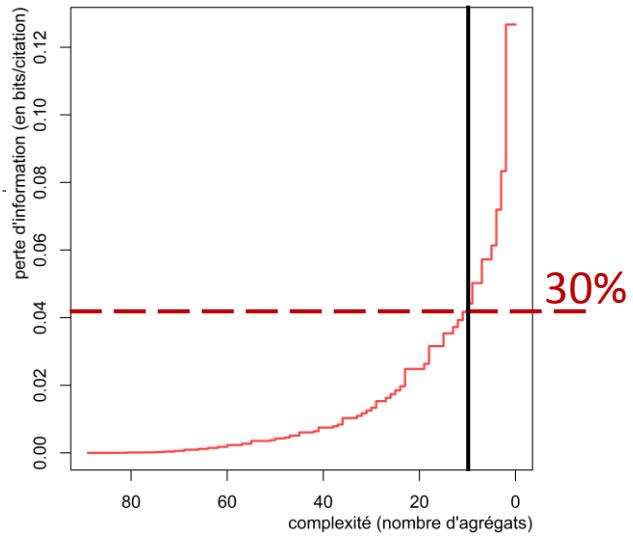
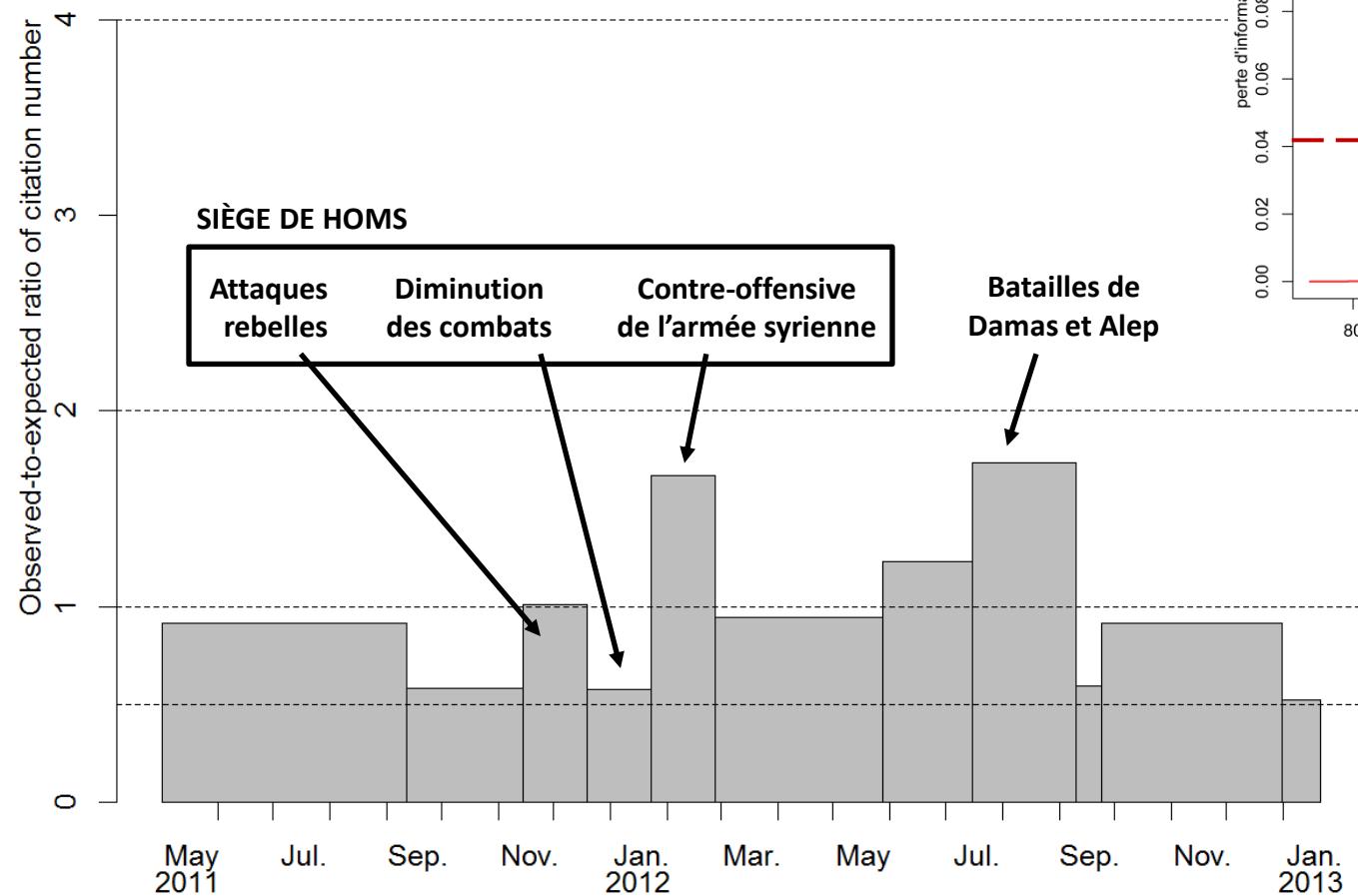


# La Syrie vue par LE MONDE

[Giraud, Grasland, Lamarche-Perrin *et al.*, ECTQG 2013]

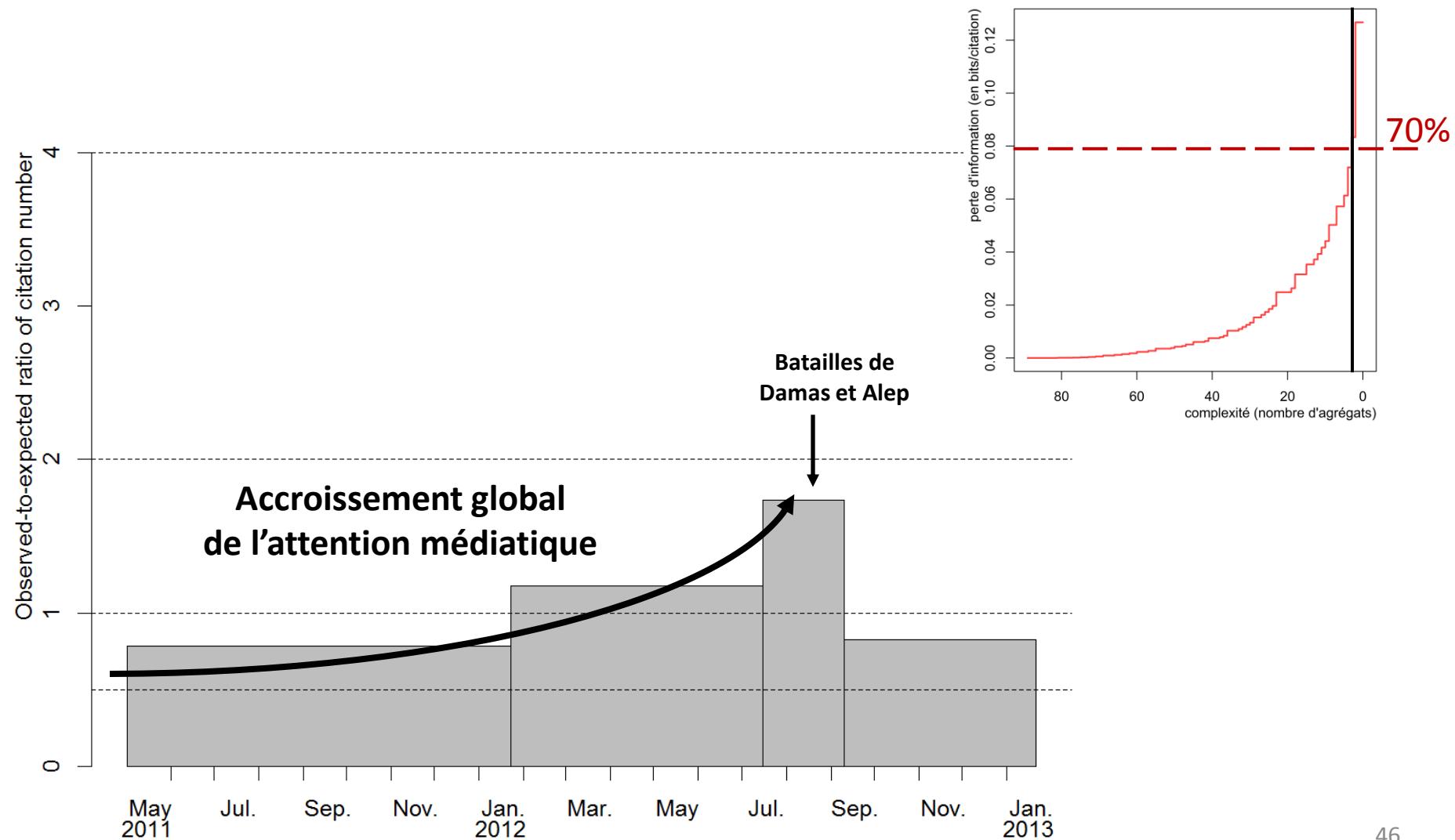
Source : Wikipedia

*Timeline of the Syrian civil war*  
*Siege of Homs*



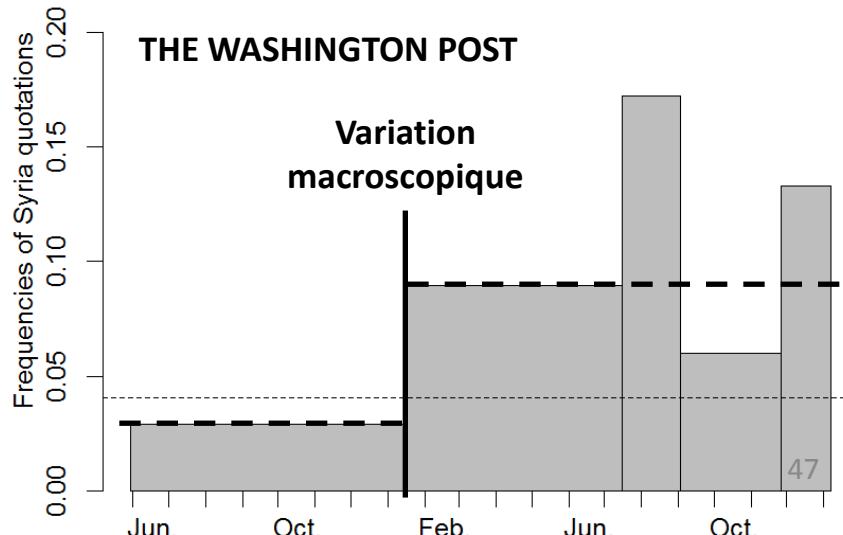
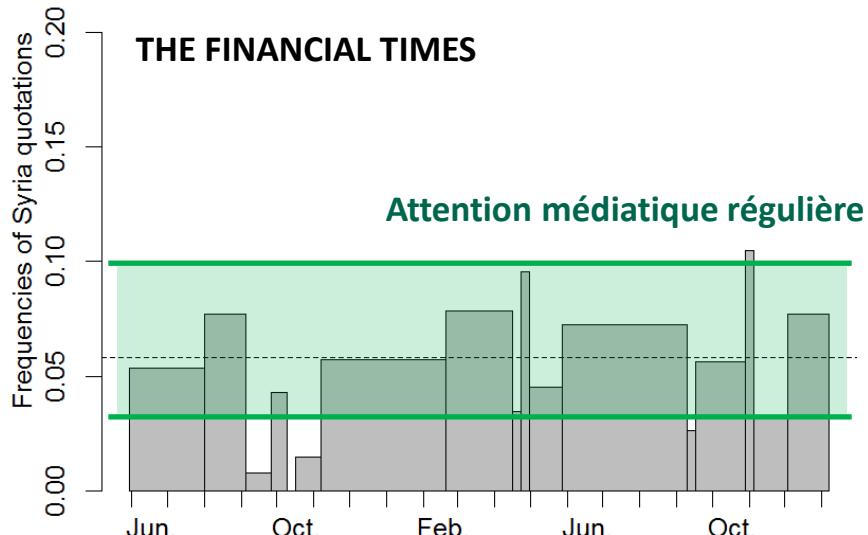
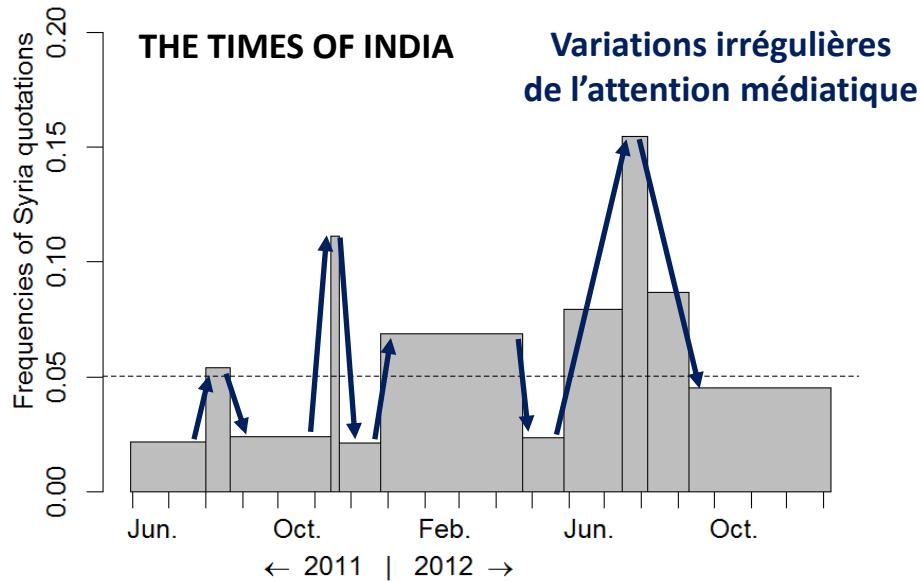
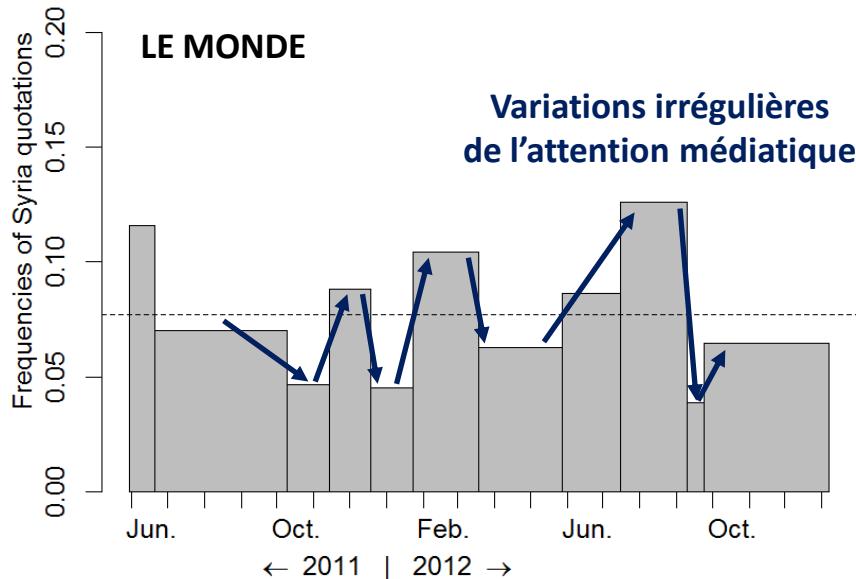
# La Syrie vue par LE MONDE

[Giraud, Grasland, Lamarche-Perrin *et al.*, ECTQG 2013]

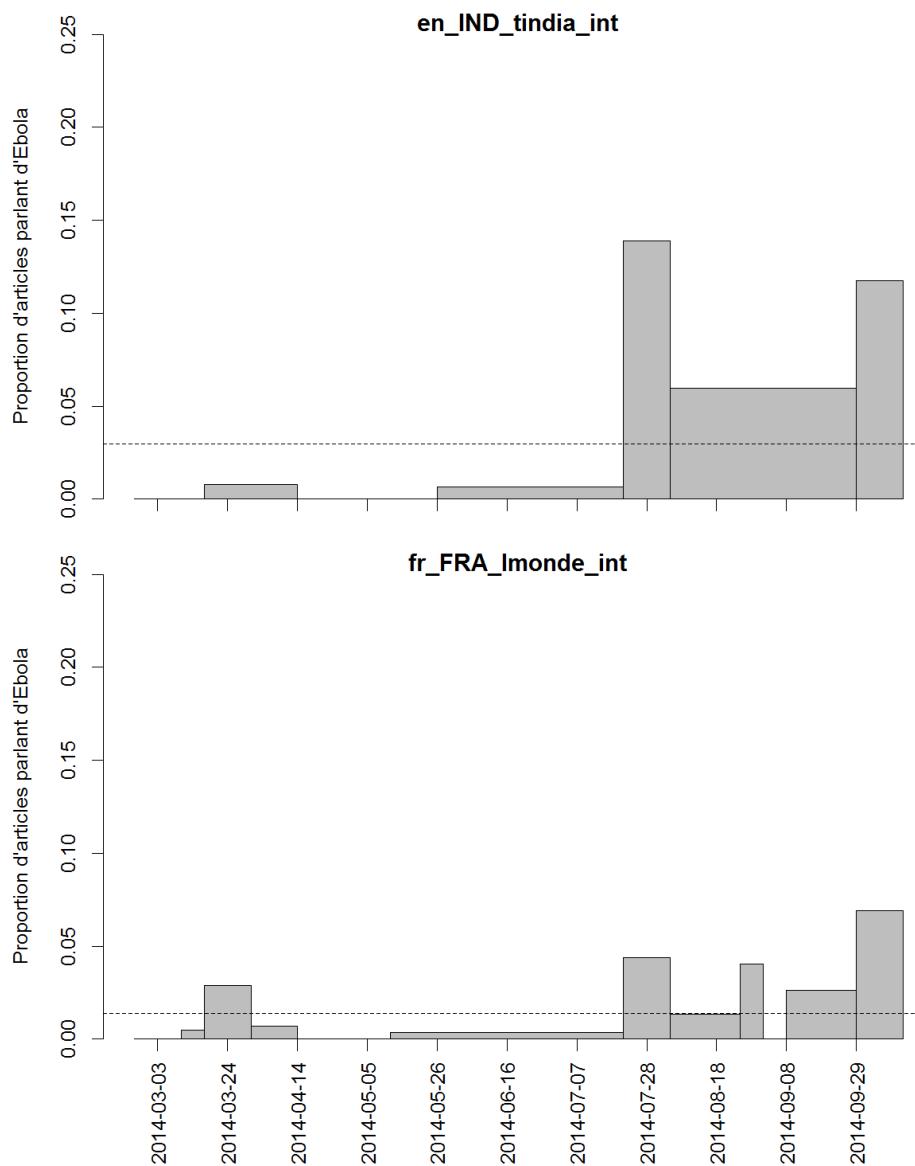
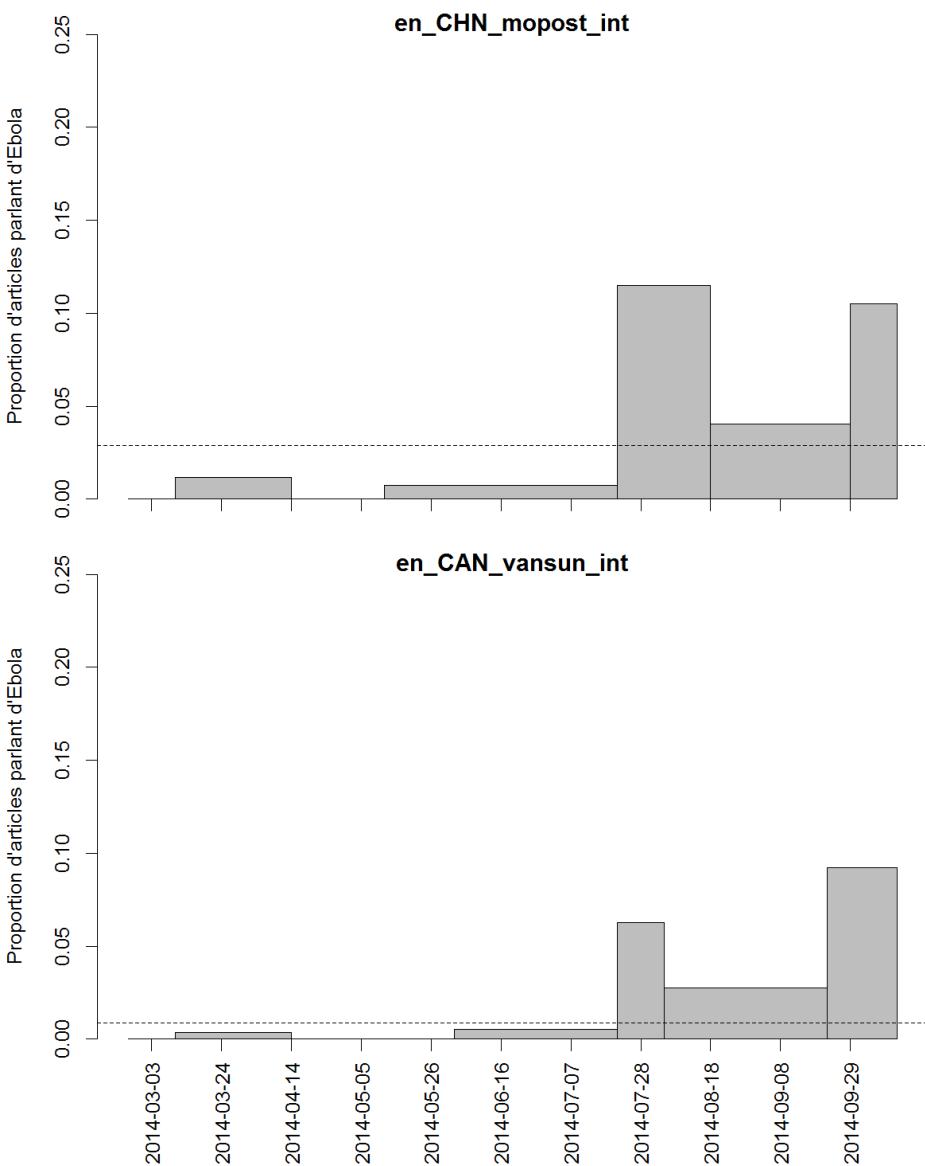


# La Syrie vue par 4 journaux

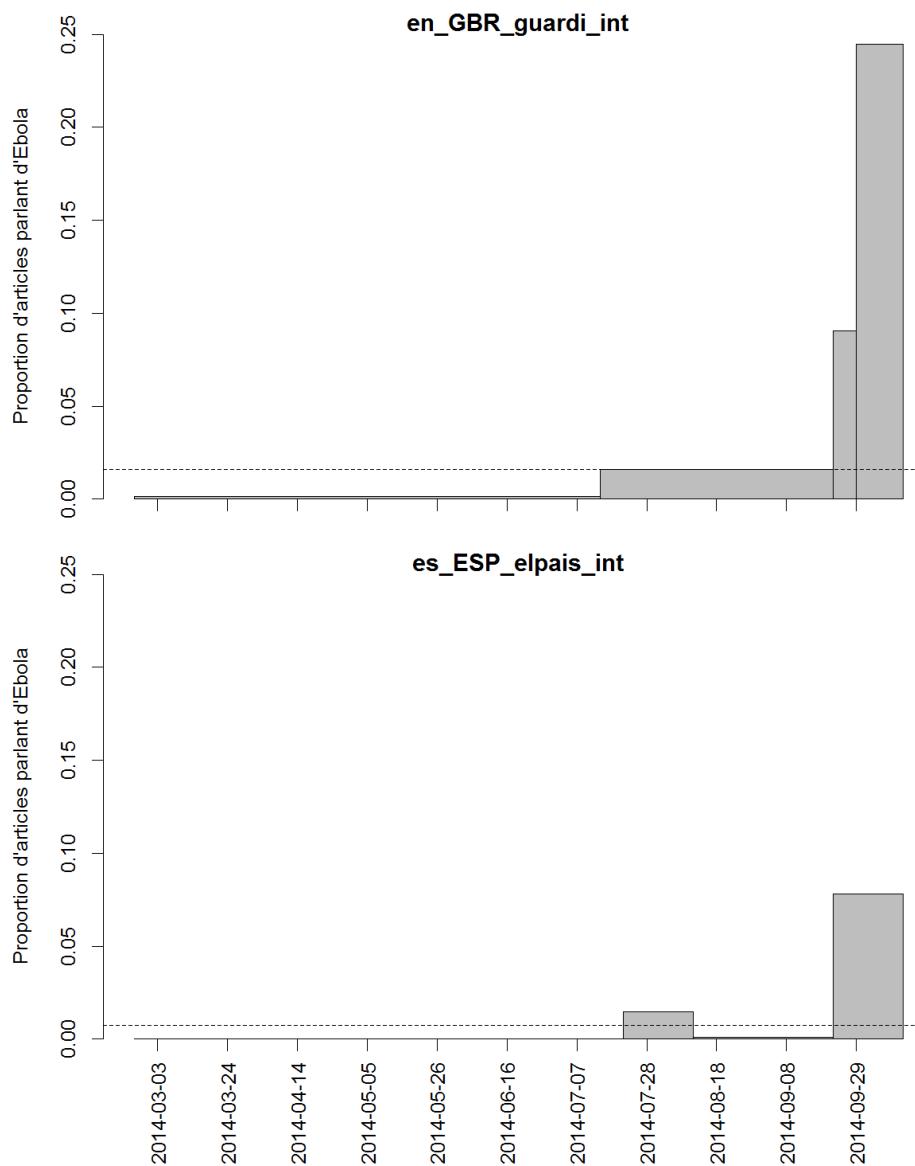
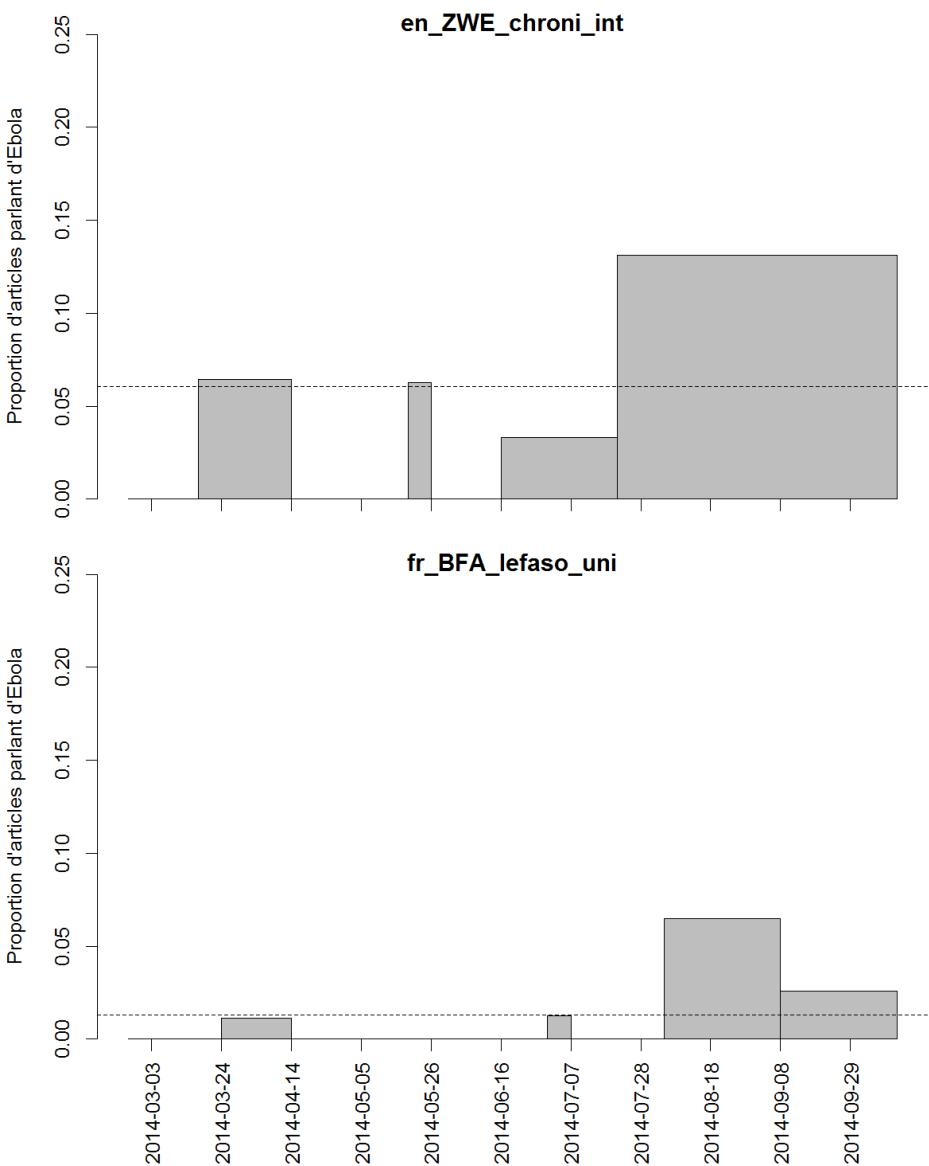
[Giraud, Grasland, Lamarche-Perrin *et al.*, ECTQG 2013]



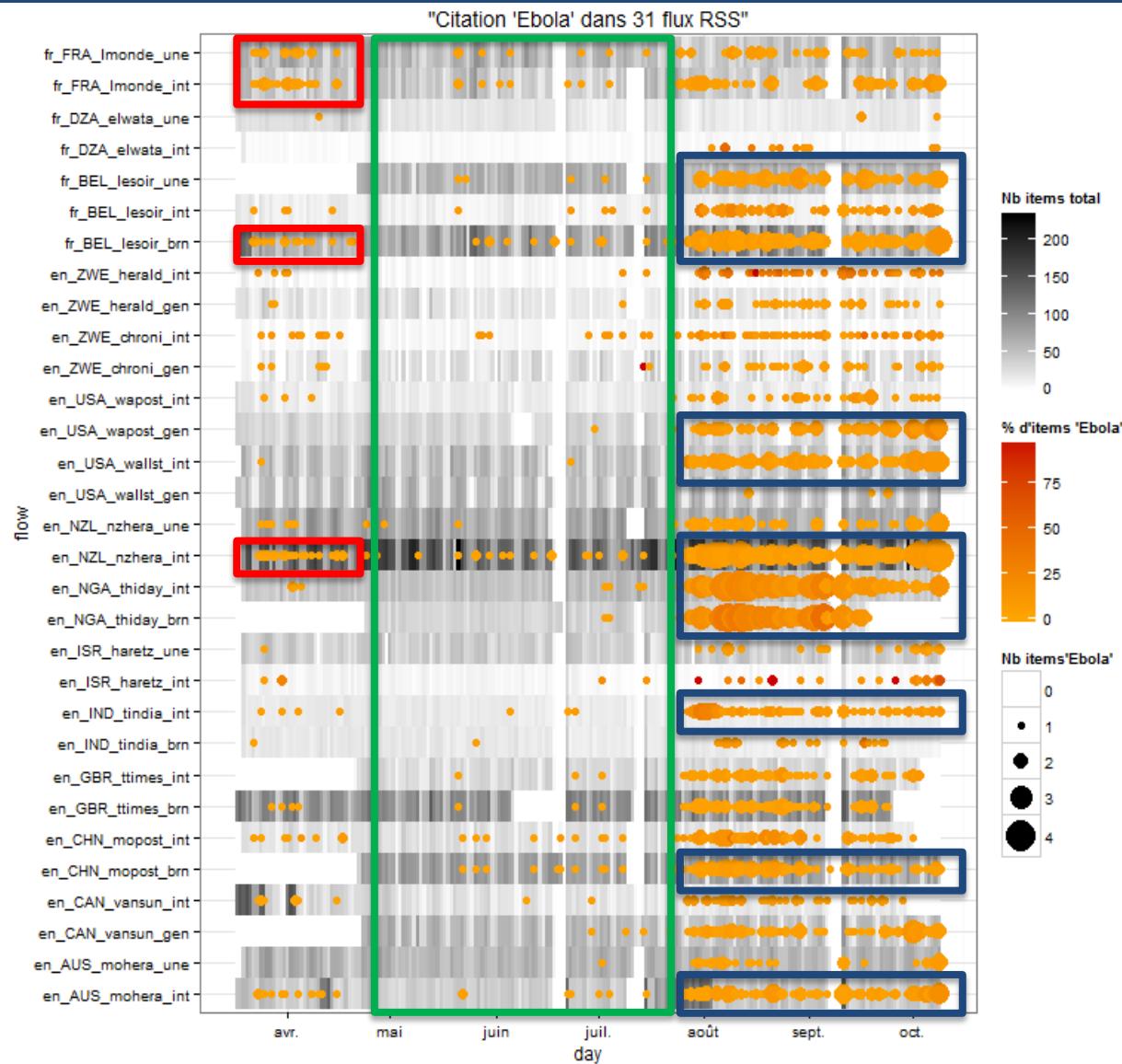
# Ebola vu par 8 journaux (90% d'information)



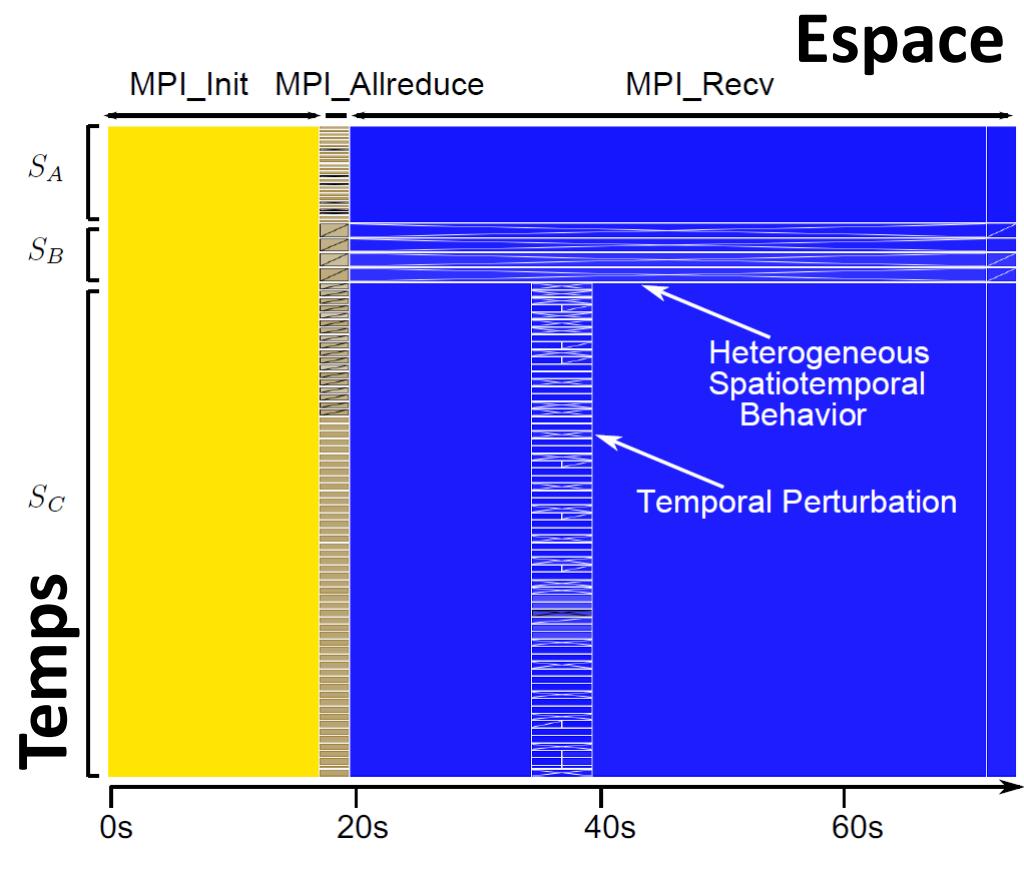
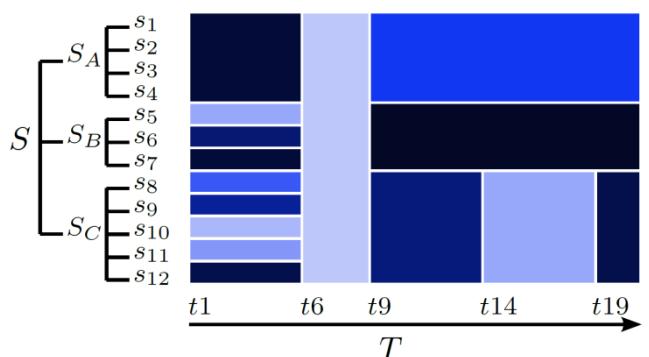
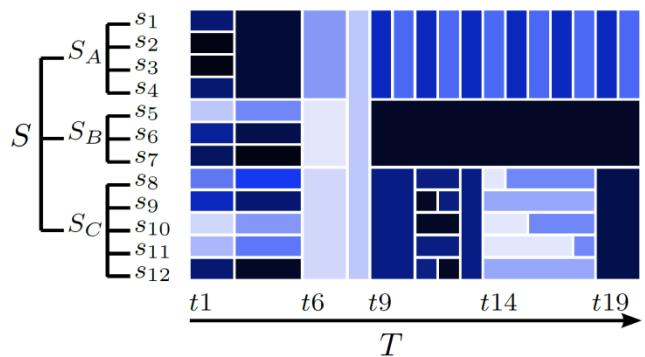
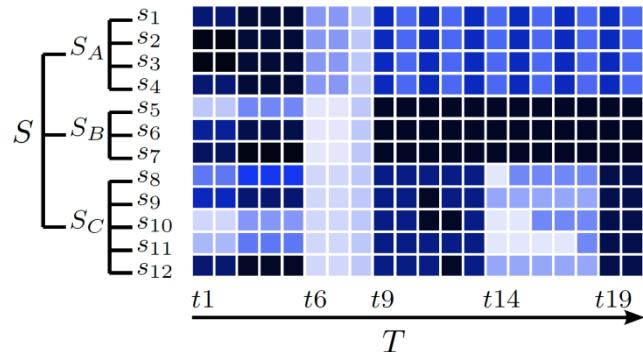
# Ebola vu par 8 journaux (90% d'information)



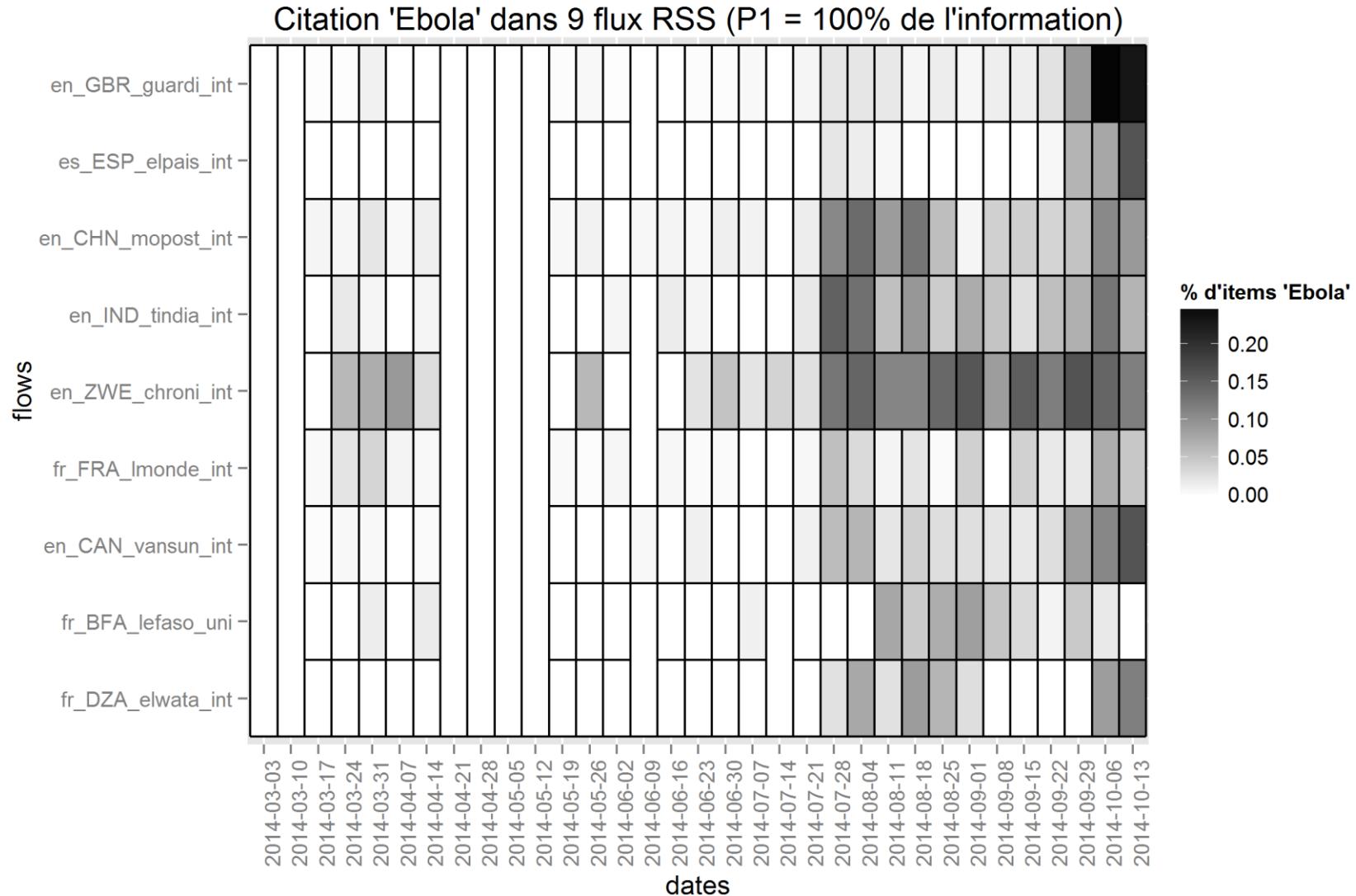
# Agrégation médiatique



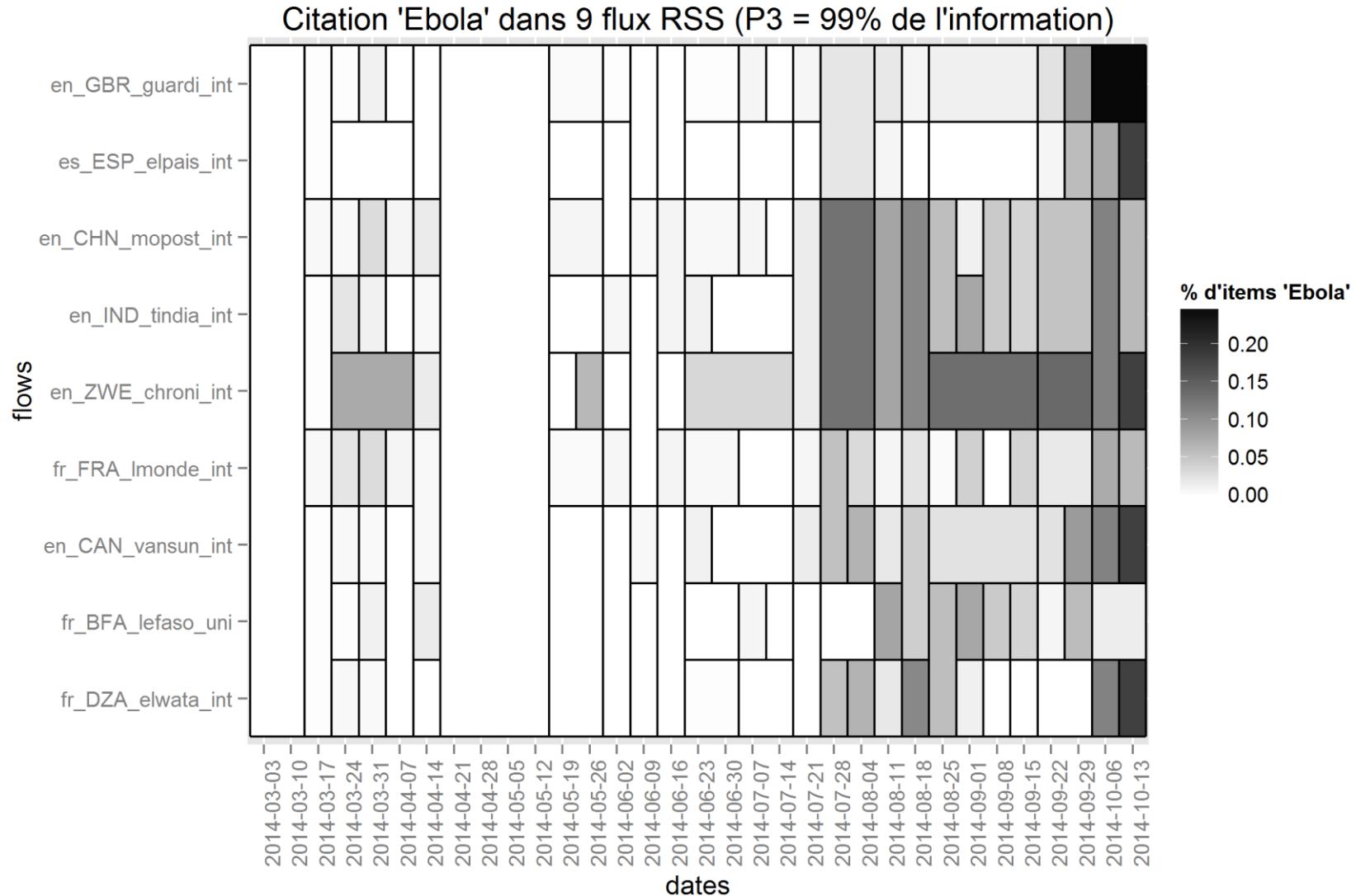
# Agrégation spatiotemporelle



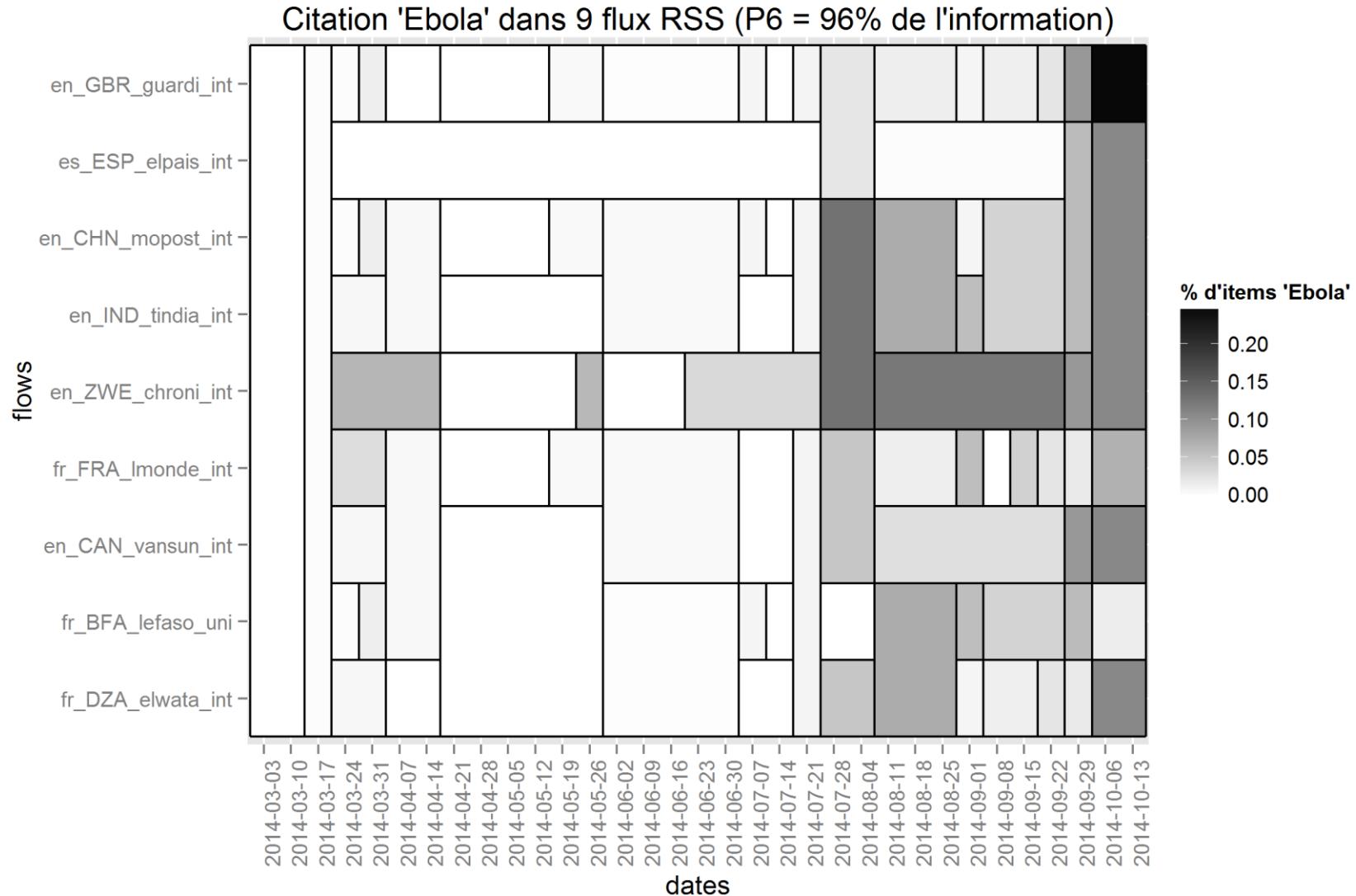
# Agrégation médiatique



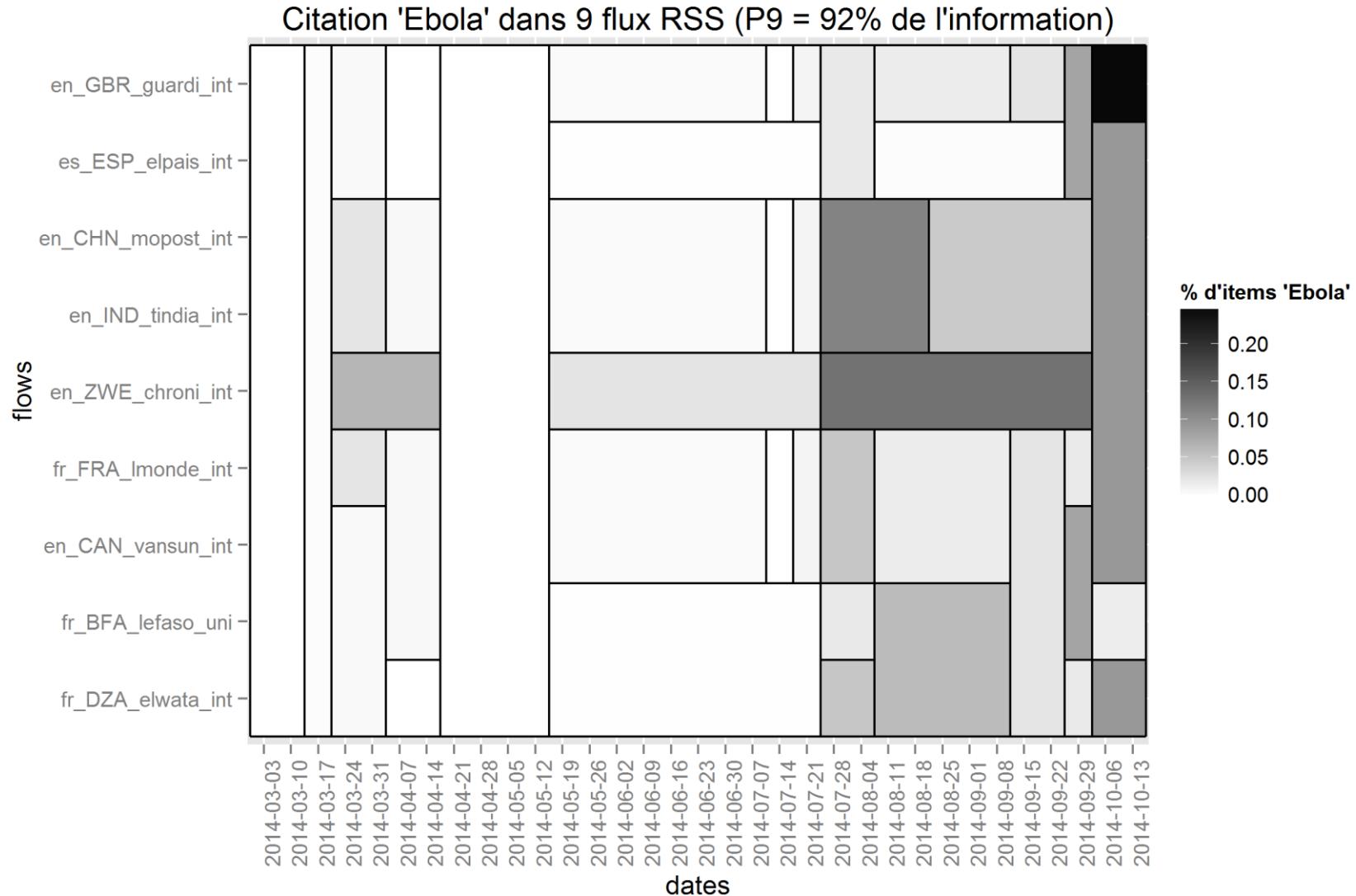
# Agrégation médiatique



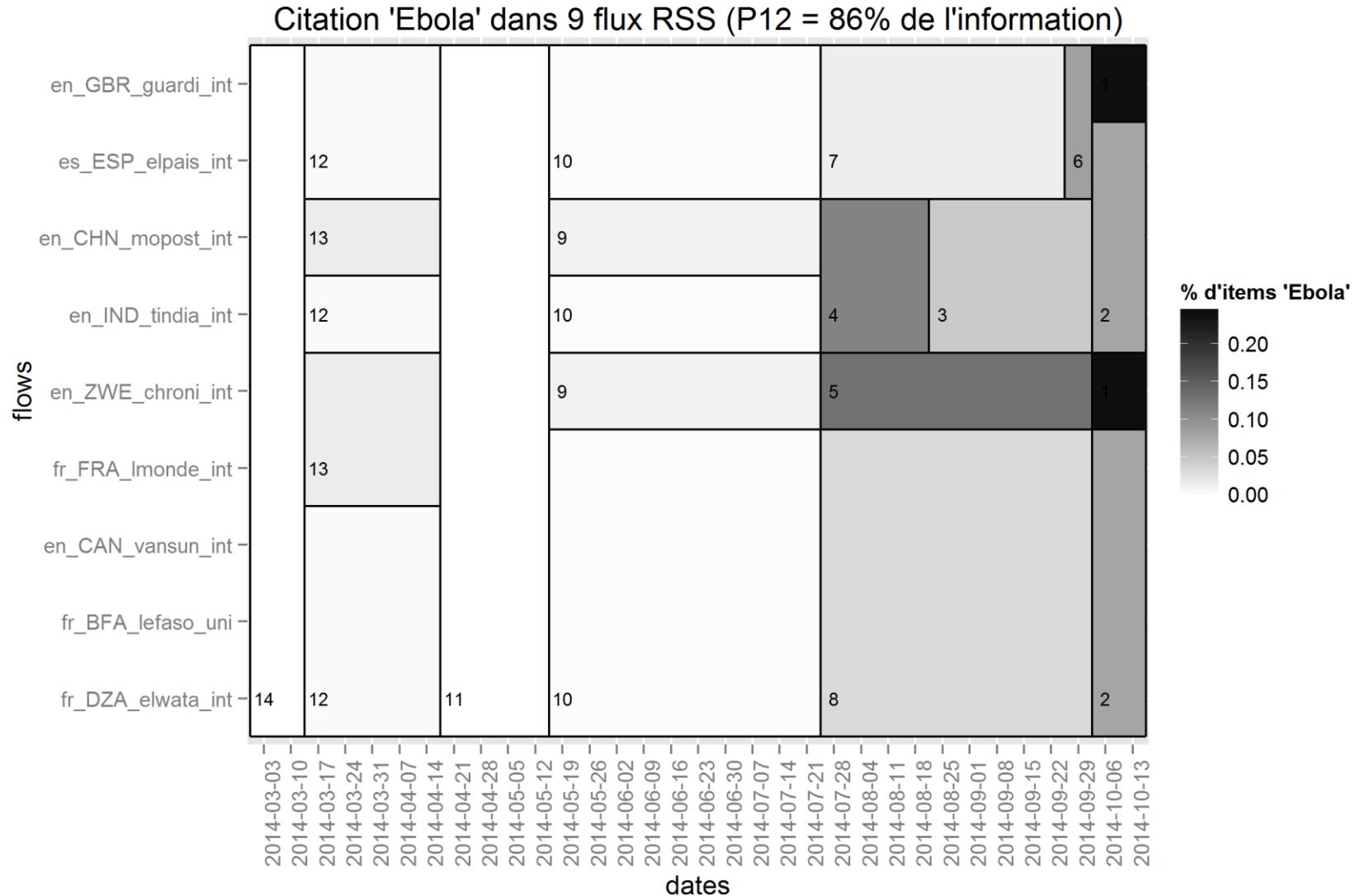
# Agrégation médiatique



# Agrégation médiatique



# Agrégation médiatique



# Agrégation et modèle de données

	Espace							
	$\pi$	USA	Libye	Syrie	France	Israël	...	Total
T	2 mai	25	12	11	10	4	...	142
	9 mai	14	6	12	12	5	...	108
	16 mai	20	11	12	6	9	...	142
	23 mai	15	9	6	13	5	...	120
	30 mai	10	16	17	9	...	...	137
$t_1$	6 juin	14	16	12	9	...	...	114
$t_2$	13 juin	15	14	11	9	52	...	119
$t_3$	20 juin	17	13	15	12	...	...	123
$t_4$	27 juin	7	6	14	20	2	...	103
	4 juill.	12	13	8	10	6	...	129
	11 juill.	21	10	10	14	3	...	107
	18 juill.	7	3	8	4	5	...	61
	25 juill.	16	7	6	13	4	...	128
	1 août	21	1	9	7	4	...	88
	...	...	...	...	...	...	...	...
	Total	423	308	260	248	153	...	3520

Citations réelles

$$v(\pi, t_1) = 12$$

$$v(\pi, t_2) = 11$$

$$v(\pi, t_3) = 15$$

$$v(\pi, t_4) = 14$$

Valeurs marginales

$$v(., t_1) = 12$$

$$v(., t_2) = 11$$

$$v(., t_3) = 15$$

$$v(., t_4) = 14$$

$$v(., T) = 459$$

Citations observées

$$v(\pi, T) = 52$$

Modèle 1

$$v^*(\pi, t) = \frac{v(\pi, T)}{|T|}$$

Modèle 2

$$v^*(\pi, t) = v(\pi, T) \frac{v(., t)}{v(., T)}$$

$$v^*(\pi, t_1) = 13$$

$$v^*(\pi, t_2) = 13$$

$$v^*(\pi, t_3) = 13$$

$$v^*(\pi, t_4) = 13$$

$$v^*(\pi, t_1) = 11.8$$

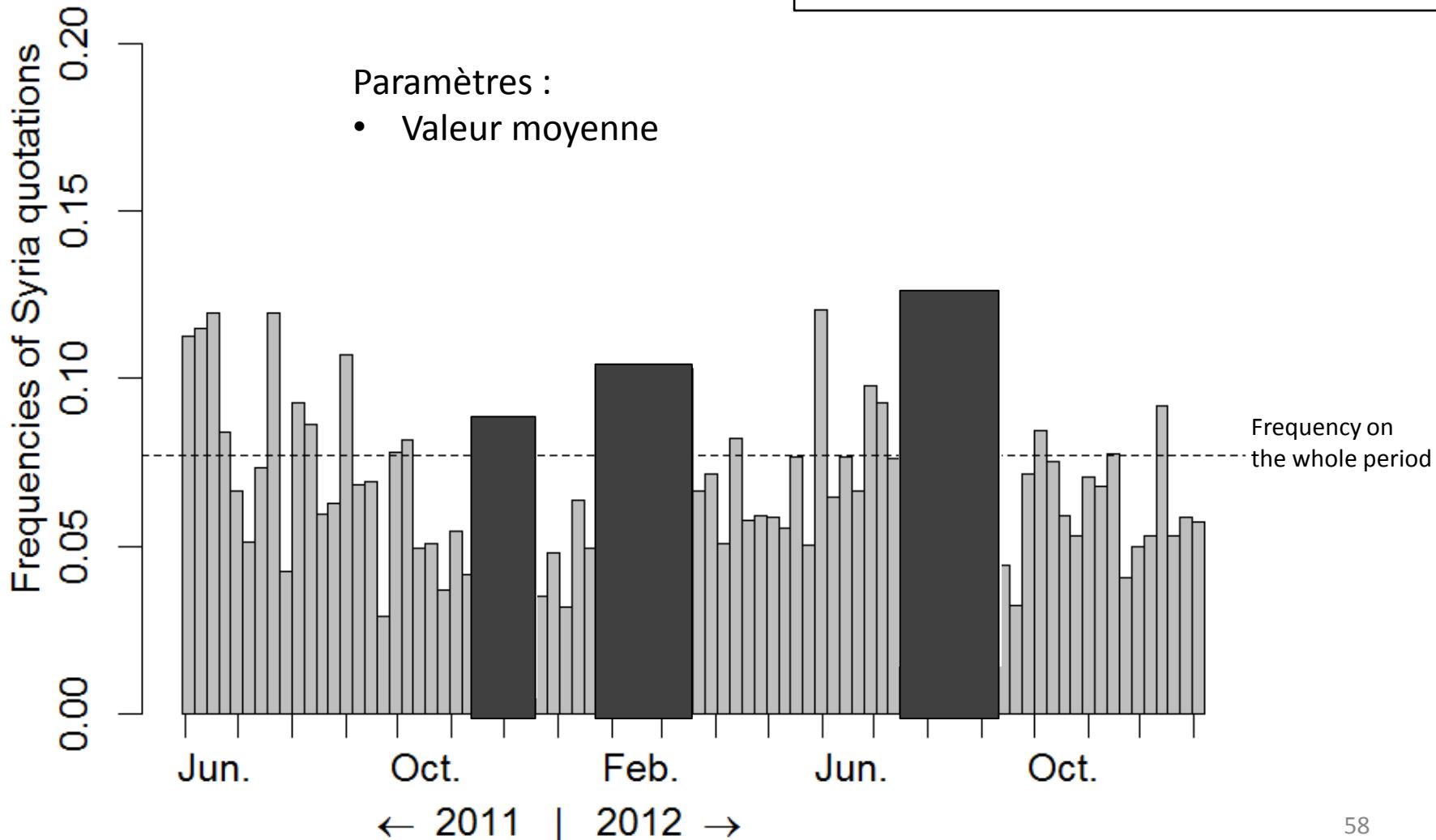
$$v^*(\pi, t_2) = 11.2$$

$$v^*(\pi, t_3) = 15.1$$

$$v^*(\pi, t_4) = 13.9^{57}$$

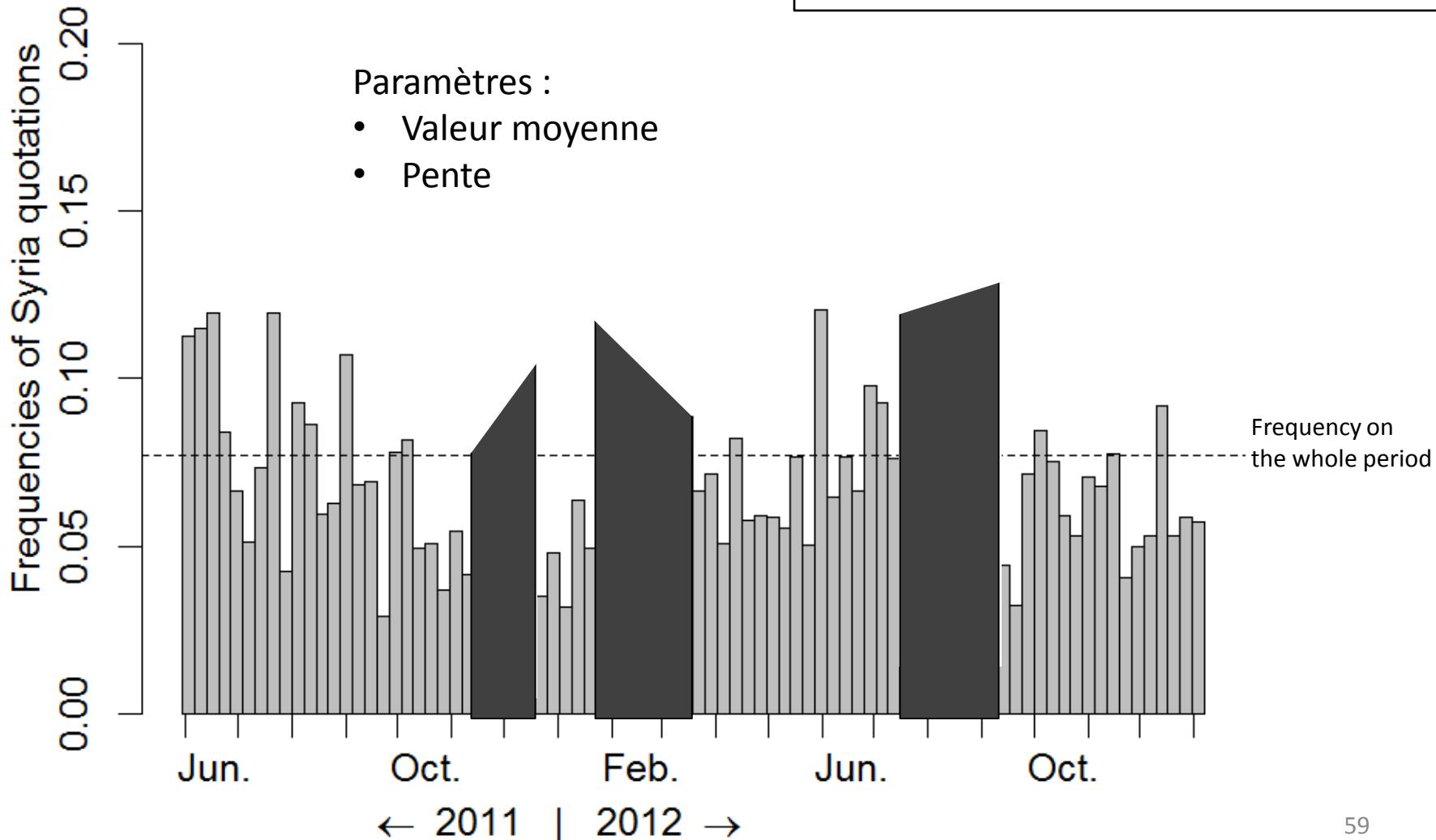
# Agrégation et modèle de données

Journal LE MONDE concernant la Syrie



# Agrégation et modèle de données

Journal LE MONDE concernant la Syrie



# ANR CORPUS GEOMEDIA

Paris, 25 novembre 2014

**MERCI POUR VOTRE ATTENTION**

**Courriel:** [Robin.Lamarche-Perrin@mis.mpg.de](mailto:Robin.Lamarche-Perrin@mis.mpg.de)

**Web:** [www.mis.mpg.de/jjost/members/robin-lamarche-perrin.html](http://www.mis.mpg.de/jjost/members/robin-lamarche-perrin.html)