

# Réduire les biais dans la collecte de tweets

**B. Mazoyer, N. Hervé, C. Hudelot, J. Cagé**

# Obtenir des données de Twitter

Dans quel but, contexte de l'étude :

- Étude du graphe social ou des interactions ?
- Domaine d'étude bien circonscrit et limité ? (comptes, mots-clés, période temporelle, ...)
- Captation *a posteriori* ou en temps réel ?
- Sampling ou exhaustivité ?

Dans notre cas :

- Analyse de la propagation d'information en français
- Comparaison avec les médias traditionnels

# Approches possibles

## Accès aux tweets :

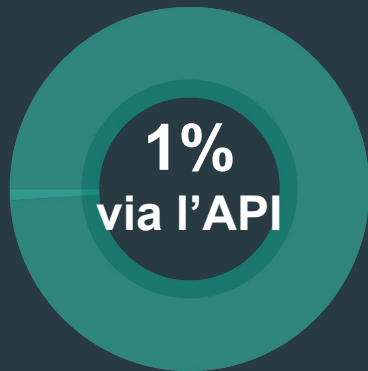
- Payant : firehose, broker et API premium / entreprise
- Gratuit : API restreinte et scraping web

## Les API :

- search : 7 jours d'historique
- streaming sample : 1% de tous les tweets
- streaming filter : critères de filtre sur requête

# Volumétrie avec l'API sample

500 millions de tweets par jour



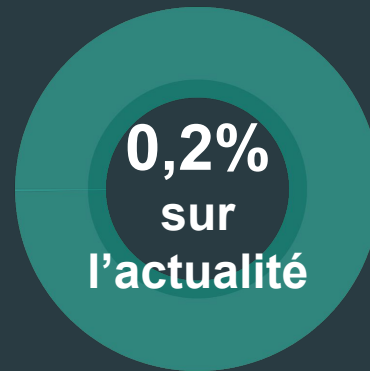
5 millions de tweets collectés

5 millions de tweets collectés



90 000 tweets en français

90 000 tweets en français



180 tweets par jour

# Approche proposée : optimisation de l'API filter

Objectif : couverture maximale des tweets en français sans biais de captation

- Utilisation de plusieurs clés d'accès
  - Dans la limite de ce qui est autorisé par Twitter (1 clé par chercheur sur le projet)
  - Dans la limite de ce qui est gérable pour les ressources machine
- Requête de mots clés neutres (stop words)
  - Détermination de ces mots clés par analyse statistique du flux sample
- Répartition des mots clés de filtre sur les clés d'API disponibles
  - Aléatoirement
  - En tenant compte de leur cooccurrence
- Nombre maximal de tweets atteignable pour une clé d'API si requête idéale
  - $1\%/1.8\% = 55\%$  des tweets français (~ 5 millions sur les 9 millions émis par jour)

# Répartition aléatoire

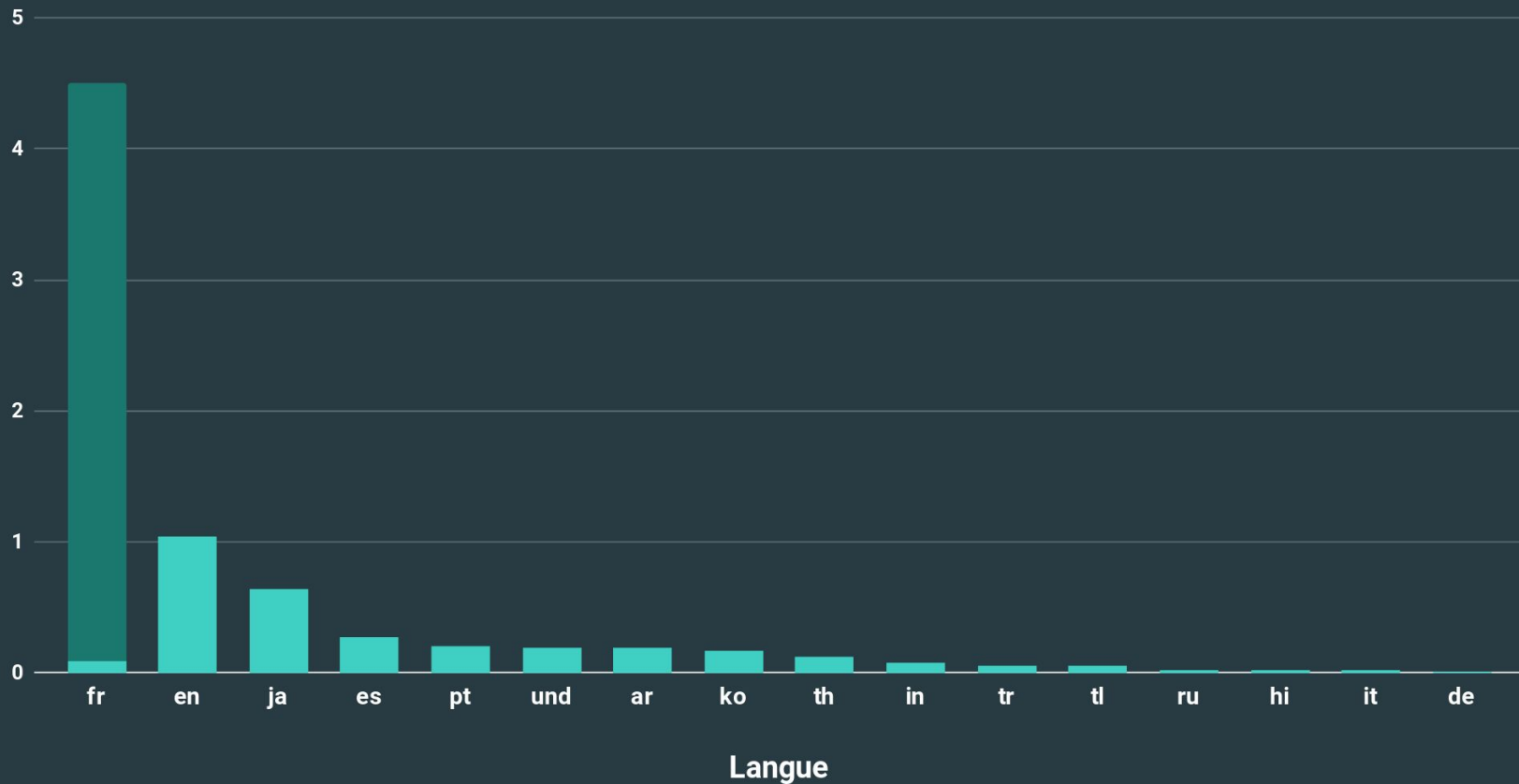


# Répartition selon la cooccurrence



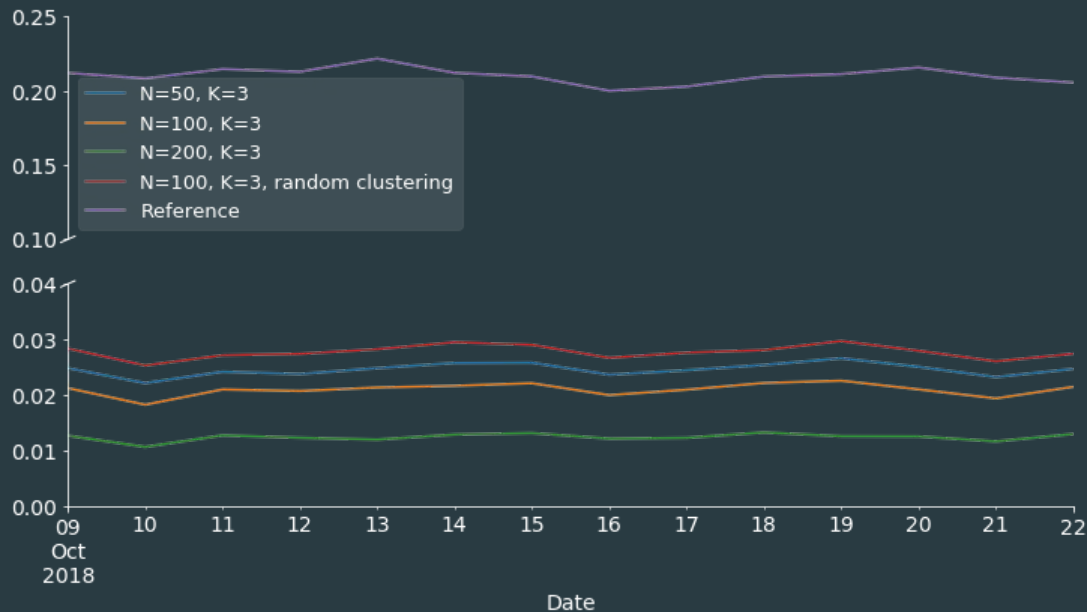
# Nombre de tweets moyen par jour (en millions)

■ Notre méthode de collecte ■ Sample Twitter





# Approche proposée : évaluation



Évolution jour par jour de la divergence de Kullback-Leibler entre la distribution des mots dans l'échantillon  $C_1$  et la distribution des mots dans chaque échantillon collecté

$$D_{K,L}(P||Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

# Approche proposée : évaluation

	$C_{1French}$	$C_{3,200} \cup C_{1French}$	$C_{3,100} \cup C_{1French}$	$C_{3,50} \cup C_{1French}$	$R_{3,100} \cup C_{1French}$
Number of characters	<b>116</b>	7*** (0)	8*** (0)	10*** (0)	9*** (0)
Share of retweets	<b>0.61</b>	0.02*** (0.00)	0.02*** (0.00)	0.02*** (0.00)	0.02*** (0.00)
Share of quotes	<b>0.18</b>	0.01*** (0.00)	0.01*** (0.00)	0.02*** (0.00)	0.01*** (0.00)
Share of replies	<b>0.19</b>	-0.02*** (0.00)	-0.01*** (0.00)	-0.01*** (0.00)	-0.02*** (0.00)
Number of URLs	<b>0.24</b>	0.01*** (0.00)	0.01*** (0.00)	0.01*** (0.00)	0.01*** (0.00)
Number of hashtags	<b>0.24</b>	-0.01*** (0.00)	-0.05*** (0.00)	-0.06*** (0.00)	-0.05*** (0.00)
Share of verified users	<b>0.01</b>	-0.00* (0.00)	-0.00 (0.00)	-0.00** (0.00)	0.00*** (0.00)
Number of followers	<b>3,073</b>	-68 (74)	-72 (74)	-165* (72)	59 (76)
Number of friends	<b>692</b>	-9* (4)	-34*** (4)	-44*** (3)	-30*** (4)
Number of lists	<b>38</b>	0 (0)	0 (0)	-0 (0)	2*** (0)
Observations	<b>875,630</b>	63,296,610	59,118,730	64,272,167	55,405,602

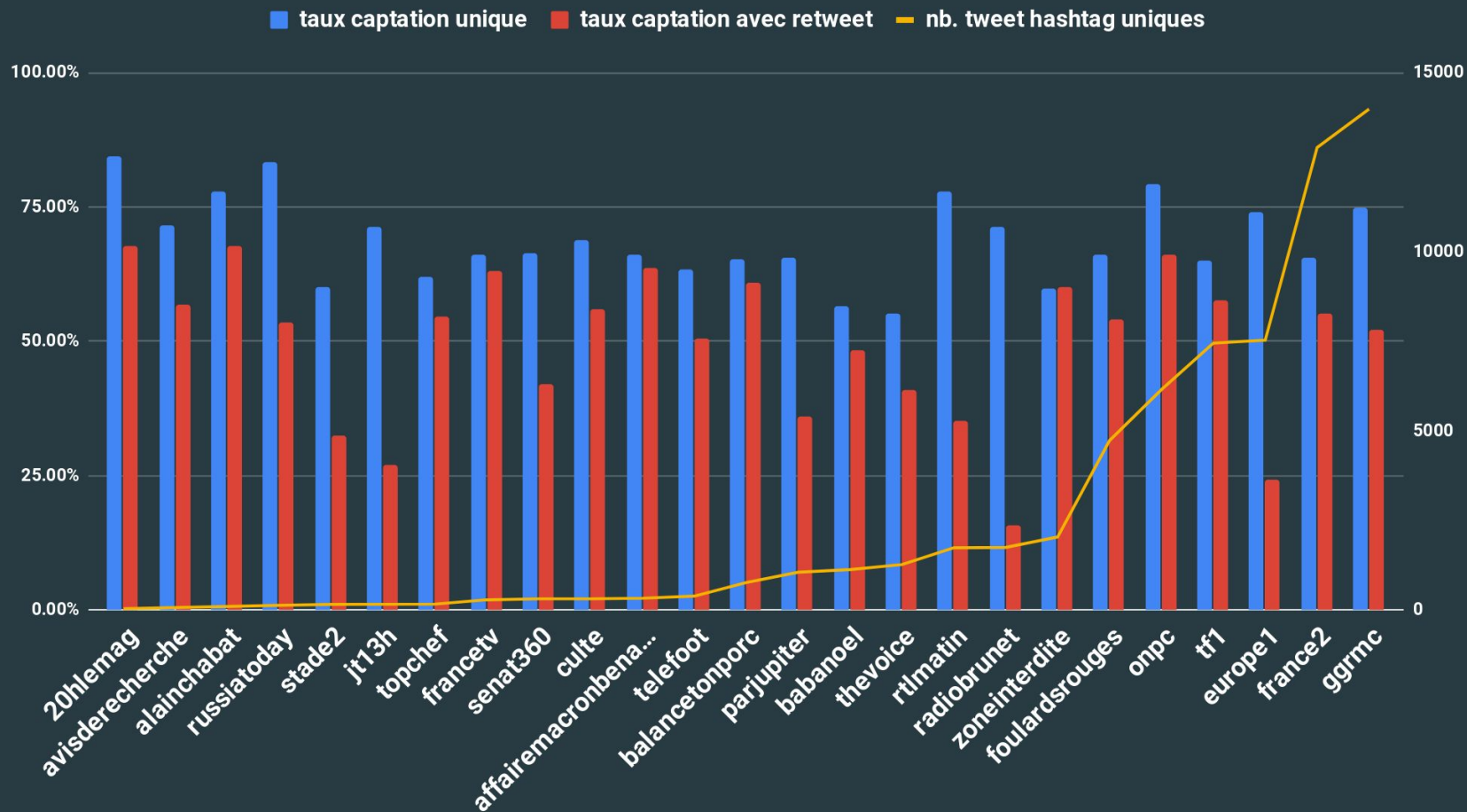
\*  $p < 0.1$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.001$ . Standard errors in parentheses.

# Approche proposée : évaluation

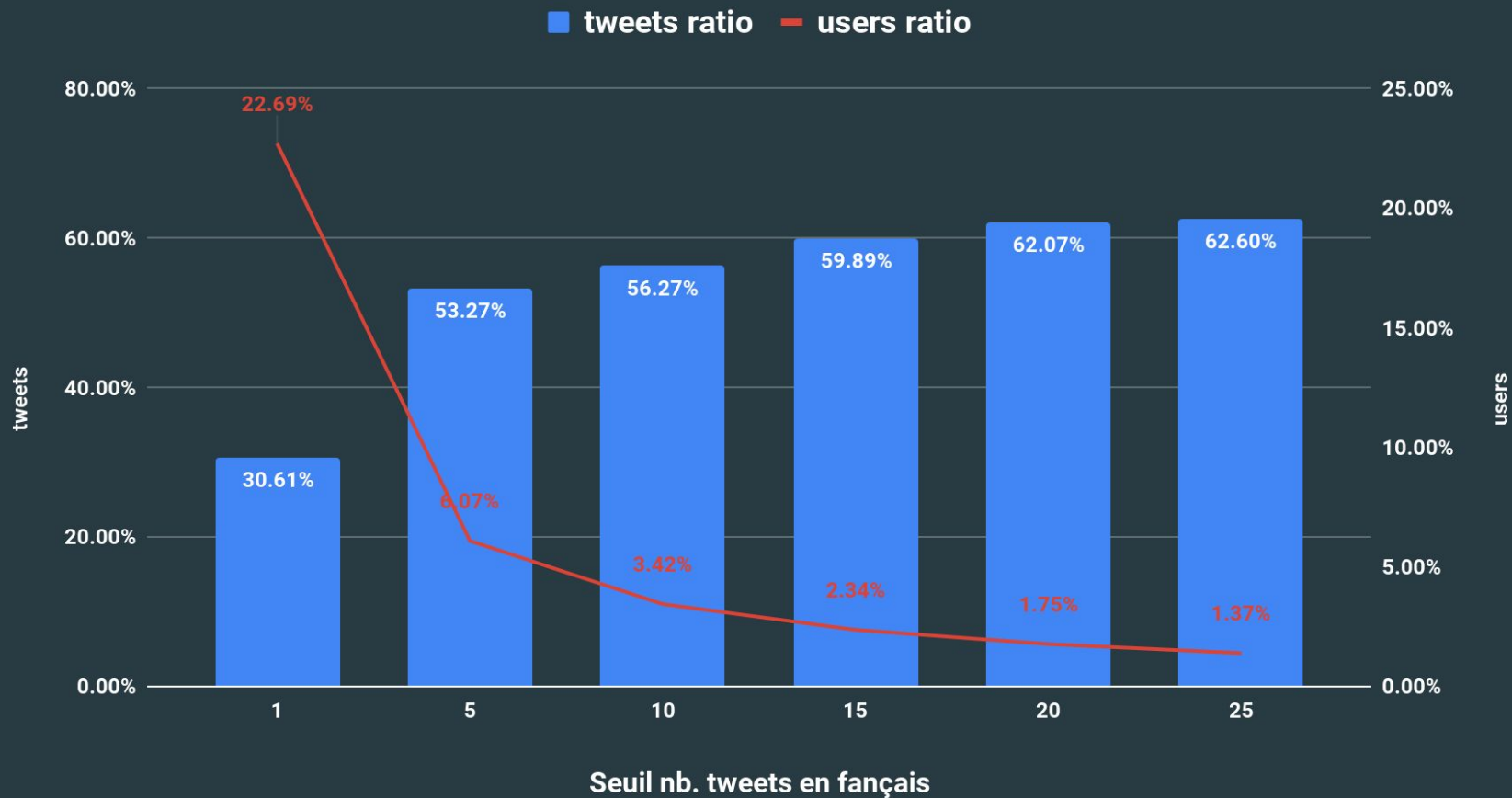
Comment évaluer la couverture de la captation ?

- On ne connaît pas le volume exact de tweets en français
- Tenir compte ou pas des retweets ?
- Étude comparée avec d'autres approches de captation
- Estimation par utilisateur

## Taux de captation comparé à des hashtags (total : 45,59%, unique : 69,82%)



# Taux de captation par utilisateur (~ 25 millions de tweets)



# Approche proposée : limitations

- Tweets sans texte (y compris via sample)
- Identification de la langue parfois incertaine
- Capacité de gestion de gros volumes de tweets :-)

# Architecture de la solution de captation

