

ATTRIBUTION D'AUTEUR : APPROCHE MULTILINGUE FONDÉE SUR LES RÉPÉTITIONS MAXIMALES

Romain Brixtel¹, Charlotte Lecluze², Gaël Lejeune³

(1) HEC - Dpt. de comportement organisationnel, Université de Lausanne

(2) GREYC, Université de Caen

(3) LINA, Université de Nantes

prenom.nom@unil.ch, unicaen.fr, univ-nantes.fr

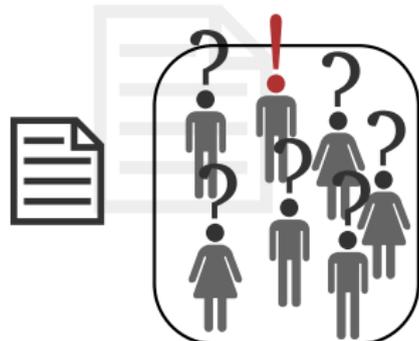
23 *Juin* 2015

ATTRIBUTION D'AUTEUR

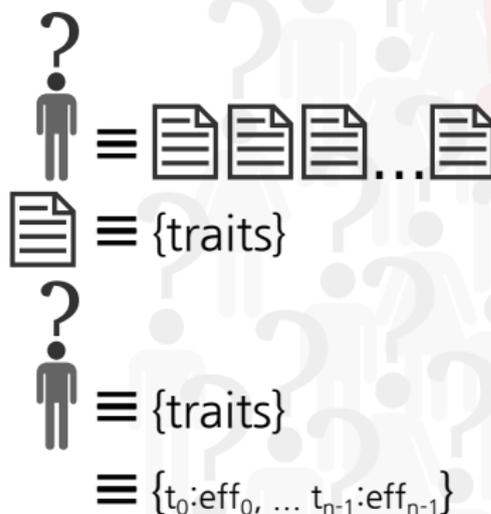
DÉFINITION

Caractériser un auteur en tant qu'individu à partir de ses écrits

Monde clos



Étude contrastive
Détection de Doppelgänger



ATTRIBUTION D'AUTEUR

DÉFINITION : ÉTAT DE L'ART

Utilisation des traits en contexte multilingue

- Apprentissage automatique : SVM
- Multilinguisme : n -grammes de caractères

Mhmm... n -grammes de caractères...

- ... n ?
- traits redondants, non ?

Proposition :

Répétitions maximales

RÉPÉTITIONS MAXIMALES (MOTIFS)

DÉFINITION

Les répétitions maximales (*motifs* dans (Ukkonen,2009)) sont des chaînes de caractères :

- **répétées** : apparaissent deux fois ou plus.
- **maximales** : pas incluses dans une autre chaîne de même effectif.



RÉPÉTITIONS MAXIMALES (MOTIFS)

EXEMPLES

motifs(A^n)

$\{A, AA, AAA, AAA, \dots, A^{n-1}\}$

motifs(HATTIVATTIAA)

$\{T, A, ATTI\}$

répétées : ni H ni V (hapaxes)

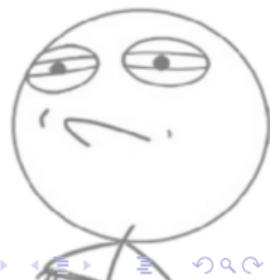
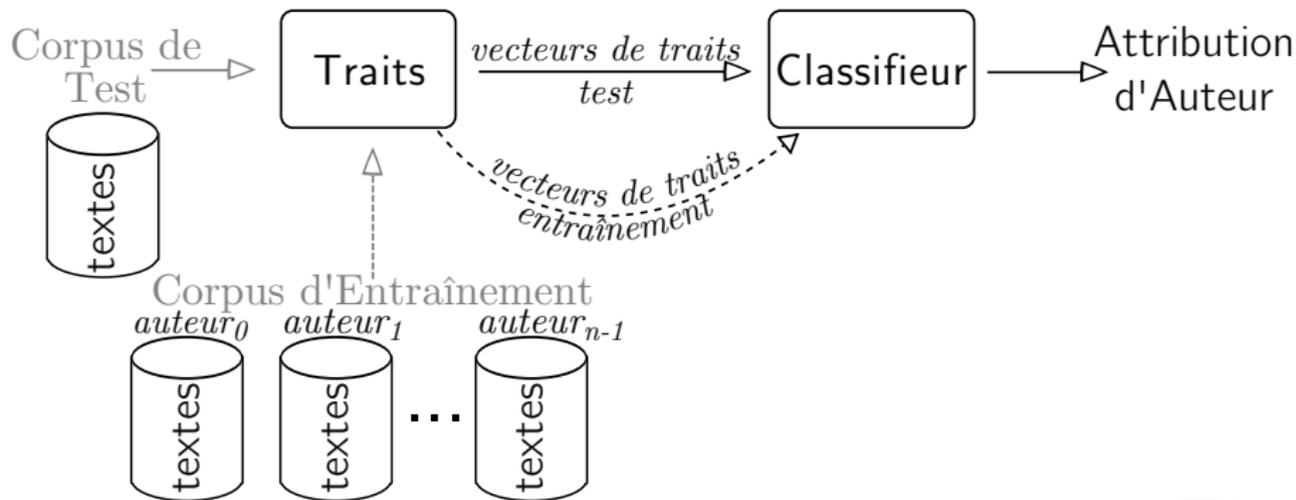
maximales : pas TT (inclu dans ATTI)

motifs(HATTIV, ATTIAA)

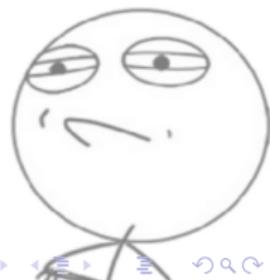
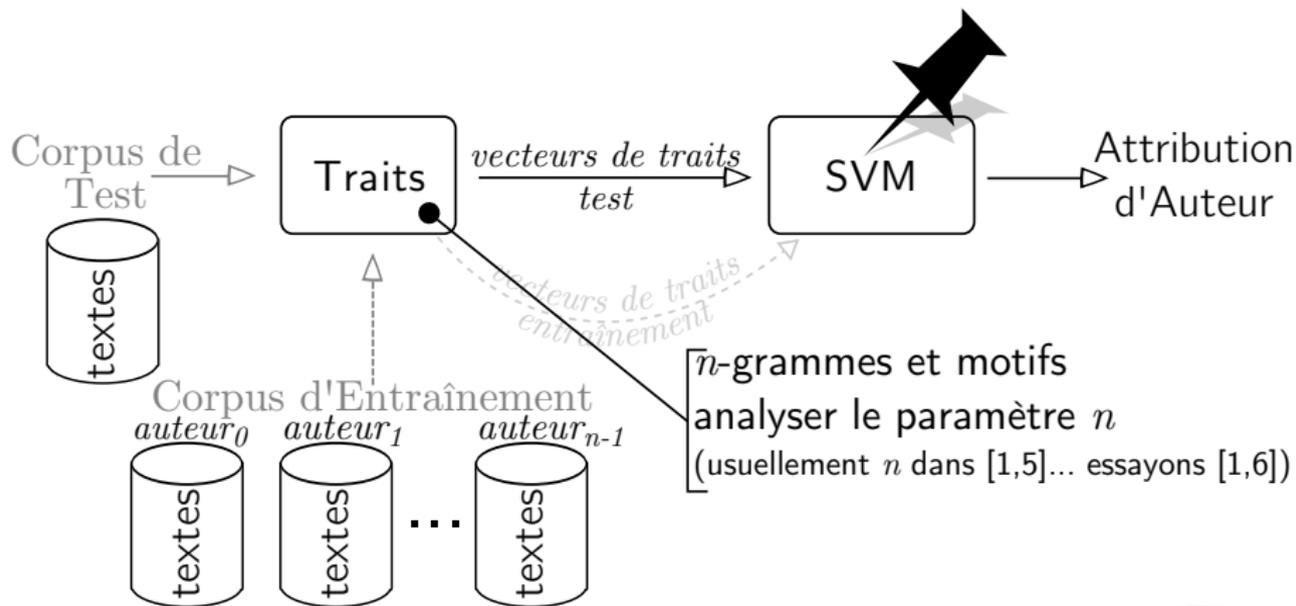
\equiv *motifs*(HATTIV\$₀ATTIAA\$₁)

\neq *motifs*(HATTIV) \cup *motifs*(ATTIAA)

CHAÎNE DE TRAITEMENT EXPÉRIMENTALE



CHAÎNE DE TRAITEMENT EXPÉRIMENTALE



CORPUS

EBG, LIB ET MIXT

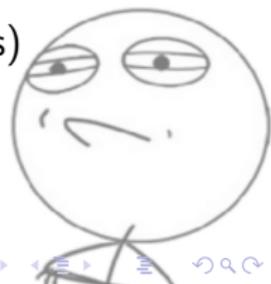
	#caractères	#textes	#auteurs
corpus	$1.9 \cdot 10^6$	631	40
auteurs (<i>moy. \pm ecart</i>)	$4.6 \cdot 10^4 \pm 8075$	15.8 ± 2.6	
textes (<i>moy. \pm ecart</i>)	2945.1 ± 178.5		

TABLE : Caractéristiques du corpus EBG (anglais)

	#caractères	#textes	#auteurs
corpus	$5.1 \cdot 10^6$	1247	40
auteurs (<i>moy. \pm ecart</i>)	$1.3 \cdot 10^5 \pm 2.6 \cdot 10^4$	31.2 ± 4.2	
textes (<i>moy. \pm ecart</i>)	4070.6 ± 1524.2		

TABLE : Caractéristiques du corpus LIB (français)

MIXT = EBG + LIB



RÉPÉTITIONS MAXIMALES ET n -GRAMMES

ÉVALUATION

Classifieur : SVM, Noyau linéaire, régularisation $C=1.0$

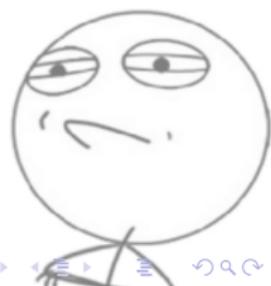
Traits : n -grammes et motifs, longueurs $\in [1, 6]$

Corpus :

- EBG (40 auteurs, anglais)
- LIB (40 auteurs, français)
- MIXT (80 auteurs, anglais et français)

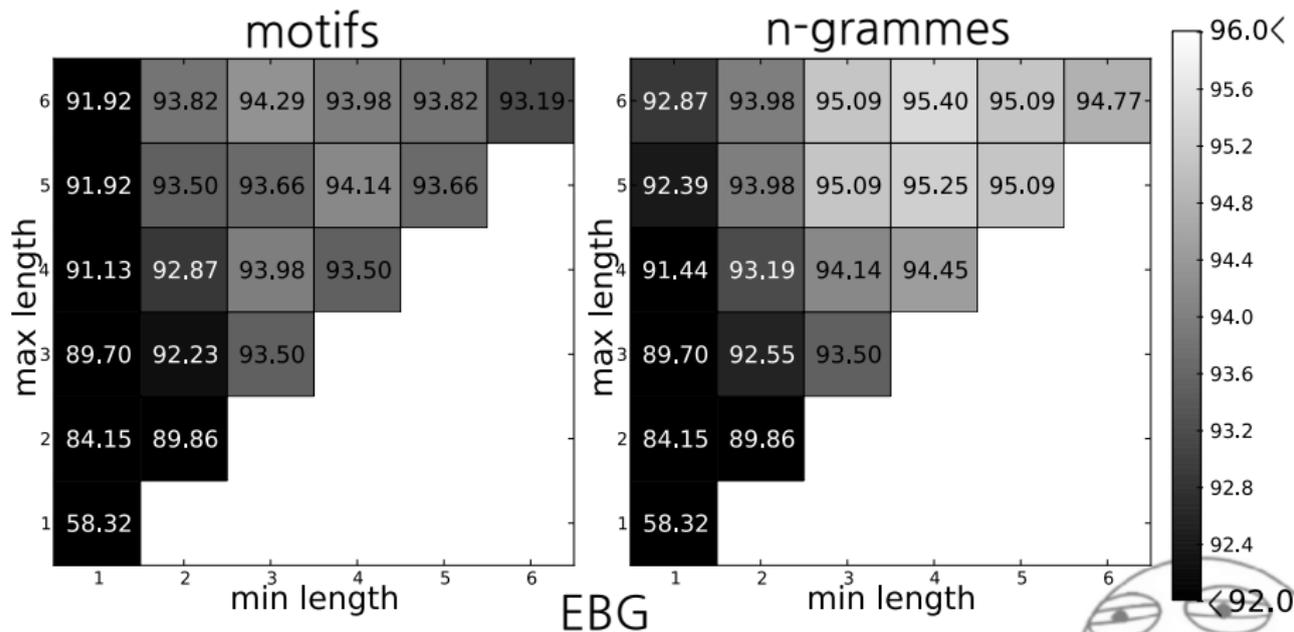
Validation croisée : stratified 10-fold

Score : $\frac{\text{\#couples (auteur, texte) correctement détectés}}{\text{\#couples (auteur, texte)}}$



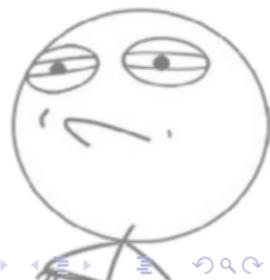
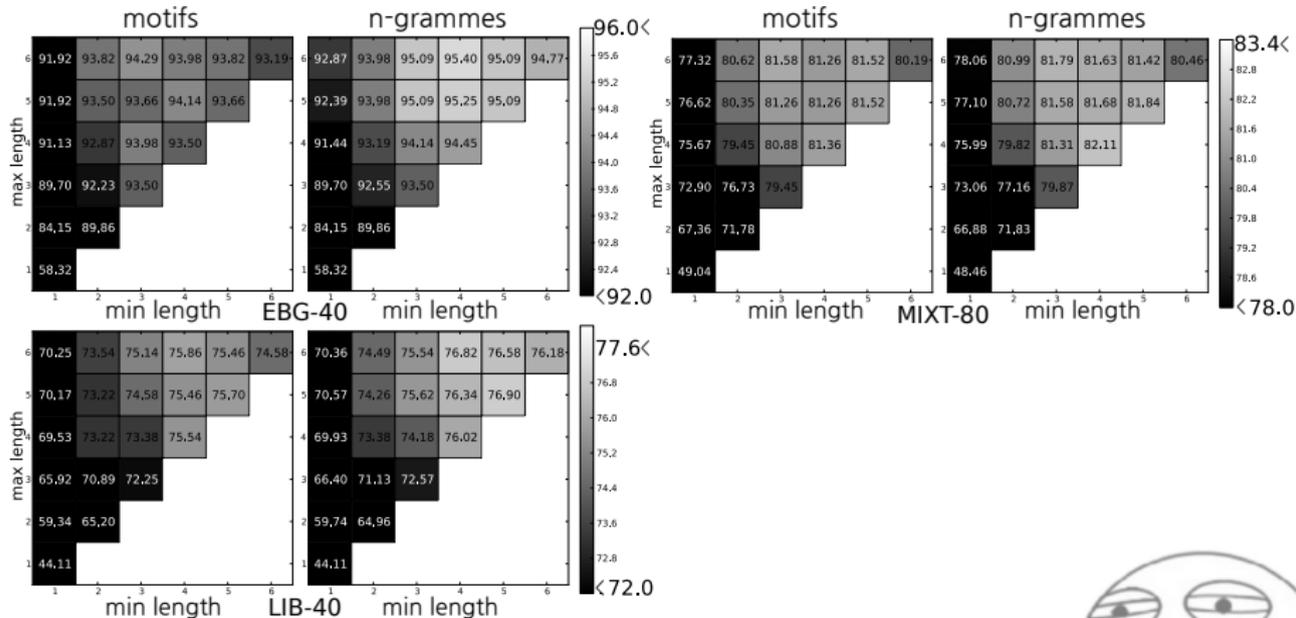
MOTIFS ET n -GRAMMES

PREMIERS RÉSULTATS...



MOTIFS ET n -GRAMMES

PREMIERS RÉSULTATS...



MOTIFS ET n -GRAMMES

PREMIERS RÉSULTATS... ET PREMIÈRES DÉCONVENUES



MOTIFS ET n -GRAMMES

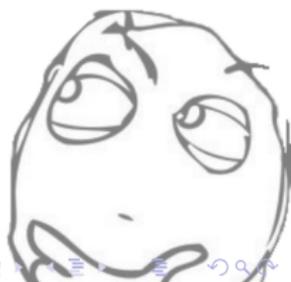
DIFFÉRENCES

Pour des longueurs comprises entre $[min, max]$:

- les motifs constituent un sous-ensemble des n -grammes
- les n -grammes « représentent » des chaînes de longueur $> max$

(min, max) -grammes :

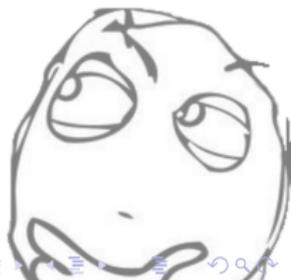
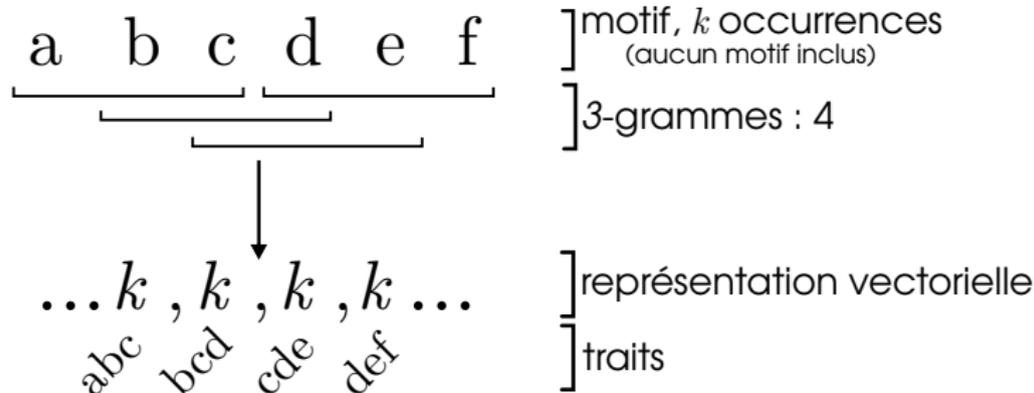
motifs] longueur $n \in (min, max)$
+ répétitions hors motifs	
+ hapaxes	



MOTIFS ET n -GRAMMES

DIFFÉRENCES

- les motifs constituent un sous-ensemble des n -grammes
- les n -grammes « représentent » des chaînes de longueur $> \max$

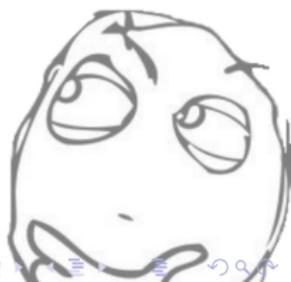


MOTIFS ET n -GRAMMES

EXPLOITATION DES DIFFÉRENCES

La taille compte : deux traits de même longueur ne sont pas pour autant d'égale importance

Les chaînes sont affectées par leurs sous-chaînes : les n -grammes « représentent » des chaînes de taille $> n$



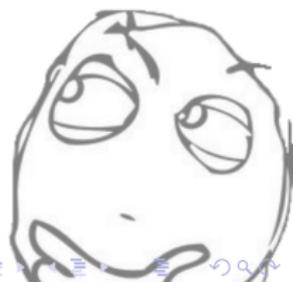
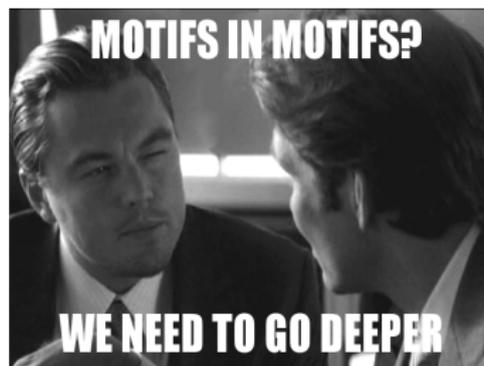
MOTIFS ET n -GRAMMES

EXPLOITATION DES DIFFÉRENCES

La taille compte : deux traits de même longueur ne sont pas pour autant d'égale importance

Les chaînes sont affectées par leurs sous-chaînes : les n -grammes « représentent » des chaînes de taille $> n$

Motifs de 2^{ème} ordre

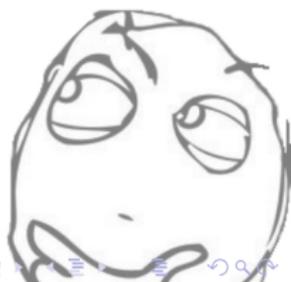
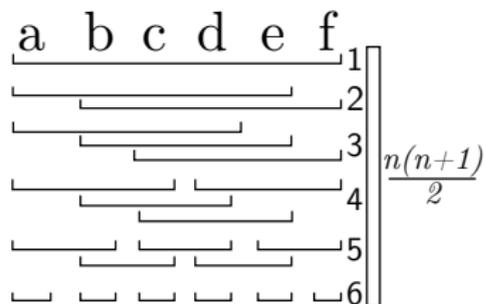


MOTIFS

PONDÉRÉS PAR DES MOTIFS DU 2^{ÈME} ORDRE

$$w_{2^{nd\ order}}(\text{trait}) = \# \text{sous-chaînes potentielles} - \# \text{motifs du } 2^{\text{ème}} \text{ ordre}$$

sous-chaînes potentielles :

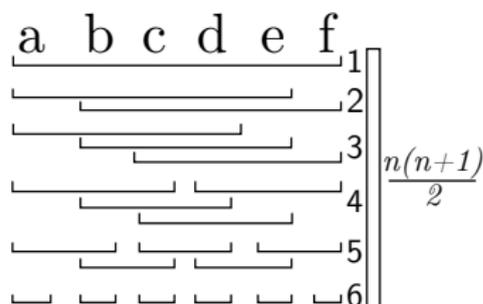


MOTIFS

PONDÉRÉS PAR DES MOTIFS DU 2^{ÈME} ORDRE

$$w_{2^{nd} order}(\text{trait}) = \# \text{sous-chaînes potentielles} - \# \text{motifs du } 2^{\text{ème}} \text{ ordre}$$

sous-chaînes potentielles :



sous-motifs (motifs du 2^{ème} ordre) :

$$\text{motifs } (\text{☰}, \text{☱}, \text{☲}, \dots, \text{☶}) = \{m_0, m_1, \dots, m_{n-1}\}$$

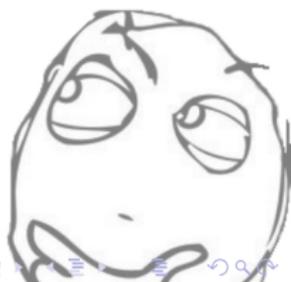
$$\text{motifs } (m_0, m_1, m_2, \dots, m_{n-1}) = \{r_0, r_1, \dots, r_{p-1}\}$$

("abcdef" est un des m_i)

$\{r_0, r_1, \dots, r_{p-1}\}$ un sous-ensemble de $\{m_0, m_1, \dots, m_{n-1}\}$

$p < n$

$$\equiv \{t_0 \cdot w(t_0) \cdot \text{eff}_0, \dots, t_{n-1} \cdot w(t_{n-1}) \cdot \text{eff}_{n-1}\}$$



MOTIFS (2^{ÈME} ORDRE) ET n -GRAMMES

ÉVALUATION

Classifieur : SVM, Noyau linéaire, régularisation $C=1.0$

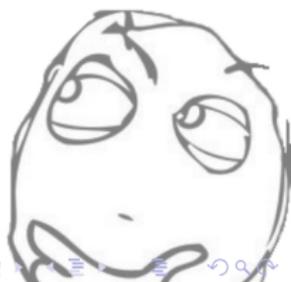
Traits : n -grammes, motifs (motifs du 2^{ème} ordre, longueur), longueur $\in [1, 6]$

Corpus :

- EBG (40 auteurs, anglais)
- LIB (40 auteurs, français)
- MIXT (80 auteurs, anglais et français)

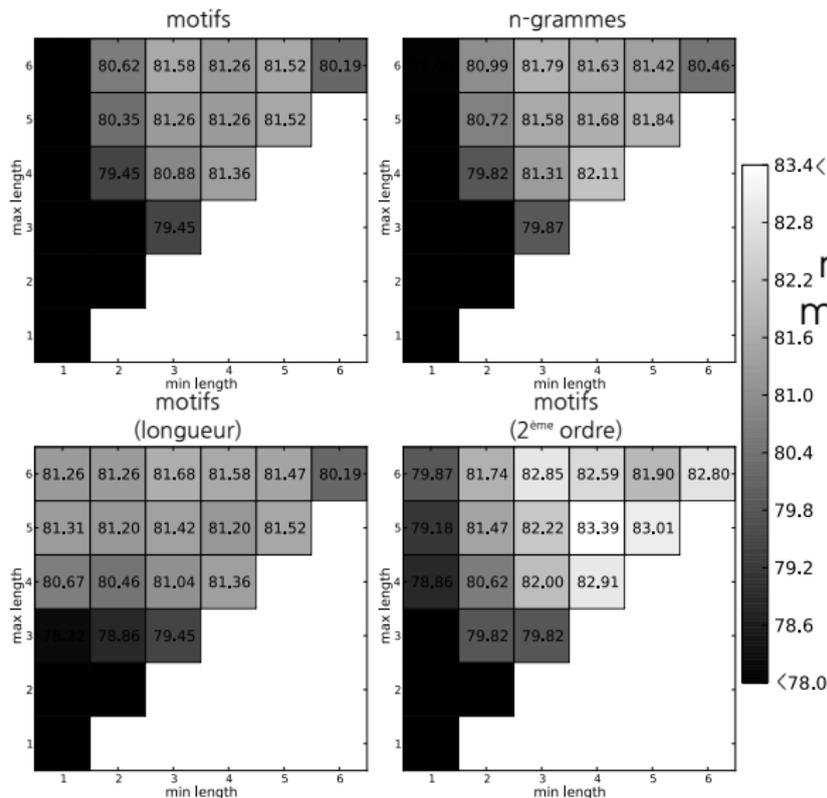
Validation croisée : stratified 10-fold

Score :
$$\frac{\# \text{couples (auteur, texte) correctement détectés}}{\# \text{couples (auteur, texte)}}$$



MOTIFS (2^{ÈME} ORDRE) ET *n*-GRAMMES

SECONDE MANCHE (MIXT)



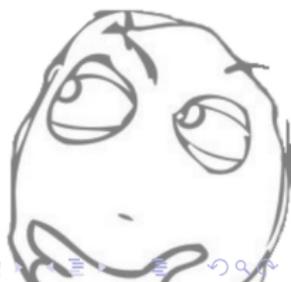
score moyen EBG / LIB / MIXT

n-grammes (4,6) 84.61%

motifs (4,6) 83.69%

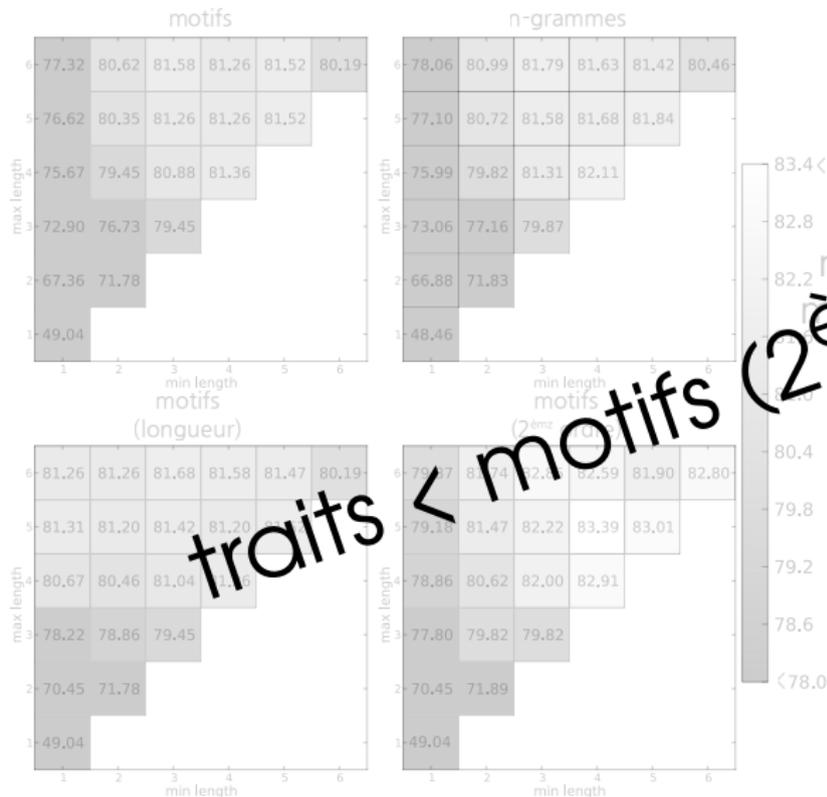
motifs (longueur) (4,6) 83.88%

motifs (2^{ÈME} ordre) **(4,5) 85.96%**



MOTIFS (2^{ÈME} ORDRE) ET *n*-GRAMMES

SECONDE MANCHE (MIXT)



traits < motifs (2ème ordre)

score moyen EBG / LIB / MIXT

n-grammes (4,6) 84.61%

motifs (4,6) 83.69%

motifs (longueur) (4,5) 83.88%

motifs (2ème ordre) (4,5) 85.96%



NOT BAD

n -GRAMMES

PONDÉRÉS PAR DES MOTIFS DU 2^{ÈME} ORDRE

Pas immédiat...

Les chevauchements entre n -grammes consécutifs créent naturellement des sous-motifs.

$3\text{-grammes}(abcdef) = \{abc, bcd, cde, def\}$

$\text{motifs}(abc, bcd, cde, def) = \{c, d, de, bc, cd\}$



MOTIFS (2^{ÈME} ORDRE) ET n -GRAMMES

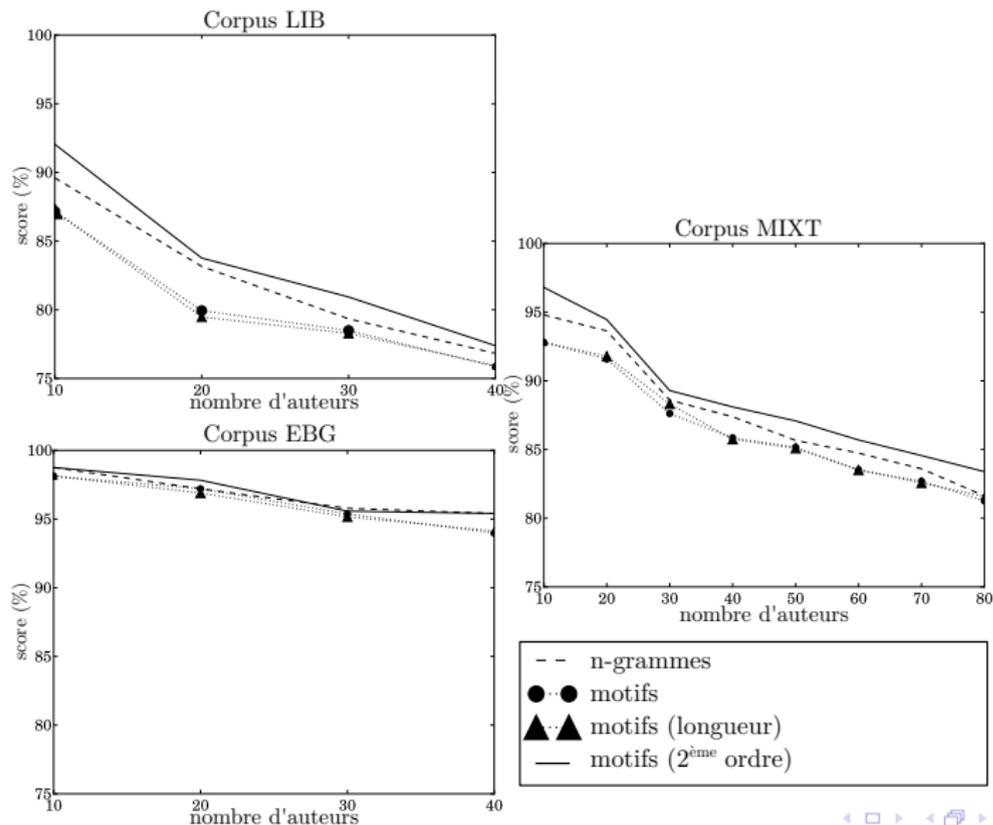
RÉSULTATS

Évaluation en fonction du nombre d'auteurs
- le meilleur score est le plus élevé -



MOTIFS (2^{ÈME} ORDRE) ET n -GRAMMES

RÉSULTATS



MOTIFS (2^{ÈME} ORDRE) ET n -GRAMMES

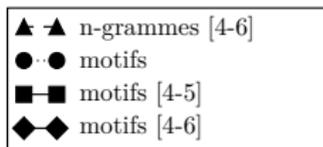
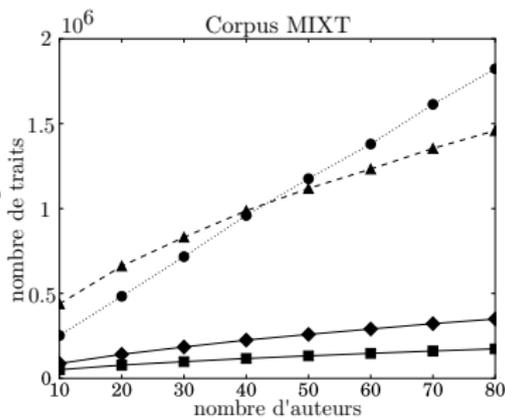
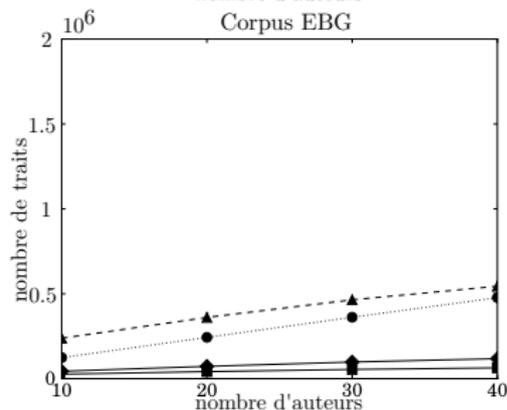
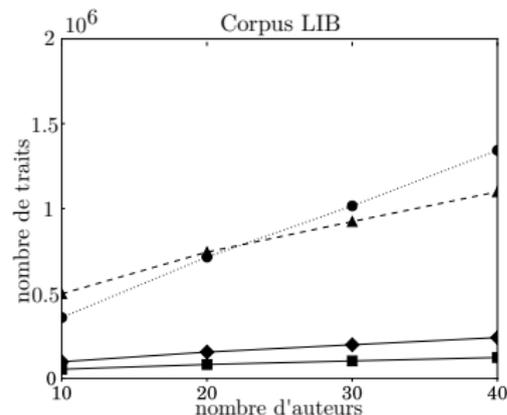
RÉSULTATS

Taille de l'espace de recherche décrit par les traits
- le meilleur score est le plus bas -



MOTIFS (2^{ÈME} ORDRE) ET n -GRAMMES

ÉVOLUTION DE LA TAILLE DE L'ESPACE DE RECHERCHE



MOTIFS (2^{ÈME} ORDRE) ET n -GRAMMES

PERSPECTIVES MULTILINGUES

Quel trait pour un corpus multilingue ?
(LIB & EBG) et (LIB & EBG à partir de MIXT)
- le meilleur score est le plus élevé -



MOTIFS (2^{ÈME} ORDRE) ET n -GRAMMES

PERSPECTIVES MULTILINGUES

n -grammes de longueur [4, 6]

nb. d'auteurs	EBG	EBG de MIXT	LIB	LIB de MIXT
10	98.75%	98.75%	89.60%	91.13%
20	97.20%	96.89%	83.15%	82.69%
30	95.79%	94.85%	79.34%	78.65%
40	95.40%	94.10%	76.82%	75.03%

Motifs de longueur [4, 5]
pondérés par des motifs du 2^{ème} ordre

nb. d'auteurs	EBG	EBG de MIXT	LIB	LIB de MIXT
10	98.75%	98.75%	92.01%	92.35%
20	97.83%	97.52%	83.77%	83.46%
30	95.59%	96.84%	80.93%	80.08%
40	95.40%	95.09%	77.38%	77.47%

TABLE : Scores sur LIB et EBG seuls et à partir de MIXT.

MOTIFS ET ATTRIBUTION D'AUTEUR

CONCLUSION

Motifs et n -grammes.

Mise en évidence de l'effet positif de la redondance des n -grammes

Pondération qui en découle adaptée aux motifs

Pour l'Attribution d'Auteur... et certainement ailleurs.

questions?

