



# **Observations diachroniques dans un corpus de presse avec le Lexicoscope**

Olivier Kraif et Sascha Diwersy  
Université Grenoble Alpes - Université Paul Valéry

*Colloque ConCorDiaL - 13-14 oct. 2022  
Lidilem & Litt&Arts, UGA*



# Plan

1. Contexte et objectif
2. Quelques outils grand public
3. Méthodes d'analyse diachronique
4. Présentation du Lexicoscope 2.0
5. Corpus *Imdiachro*
6. Fonctionnalité chronogramme et Étude de cas
7. Perspectives

# Contexte et objectifs

# Contexte et objectifs

- Intérêt des corpus arborés pour l'étude de la combinatoire
  - ~ Cooccurents syntaxiques (Evert, 2007)
  - ~ Wordsketches (Sketch Engine, Lexicoscope)
  - ~ Arbres lexicosyntaxiques récurrents (ALR)
- Nouveaux modèles performants sur différents états de langue
  - ~ Hops (Grobol & Crabbé, 2021), BERTrade (Crabbé et al., 2022)
- Projet PhraseoRoche : corpus diachronique du 13<sup>e</sup> au 17<sup>e</sup> siècle (romans de chevalerie) (Coavoux et al. 2022).
- Enrichissement du Lexicoscope pour l'étude de ces nouvelles données textuelles en diachronie.

# Outils pour l'exploration diachronique de corpus



# Quelques outils grand public

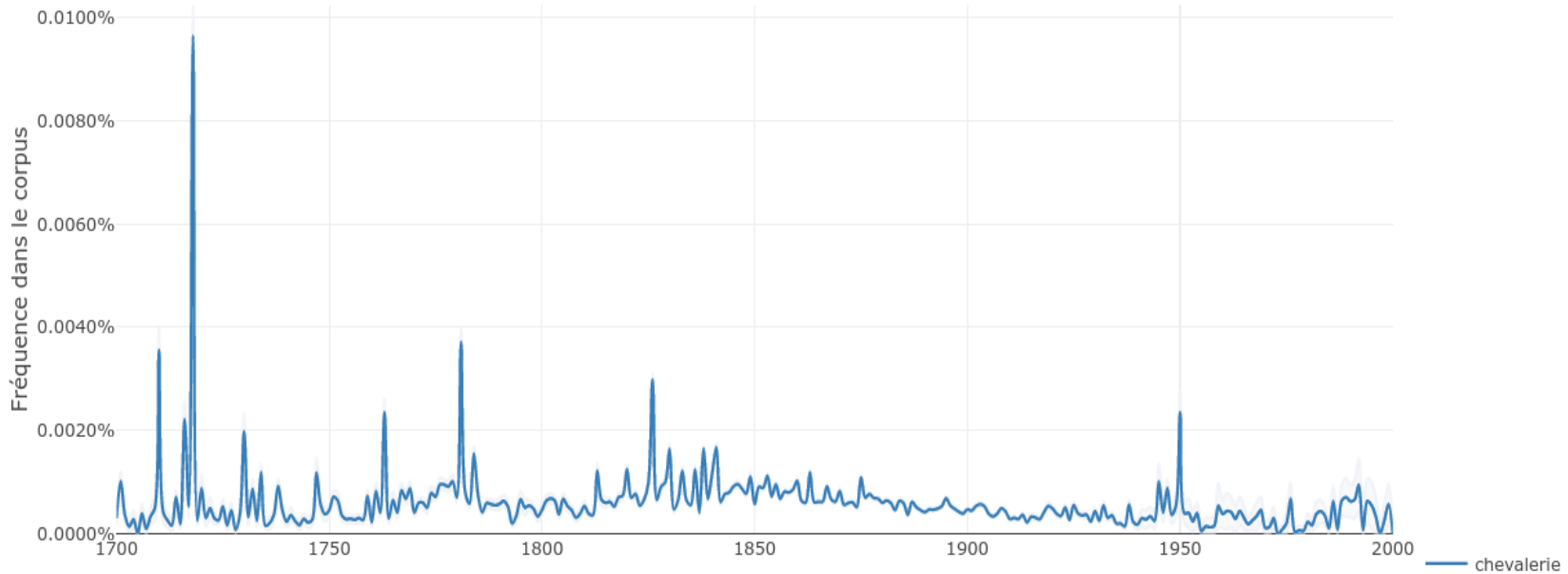
- Google Trend : popularité des requêtes (depuis 2004), mais pas de fréquence d'occurrence sur le Web.
- Sketch Engine : fonction Trends

Word	Trend ↓	Frequency	Sample
1	γ	1,033	
2	website	1,343	
3	annotators	13,723	
4	outperforms	7,918	
5	β	1,548	
6	baseline	35,179	
7	challenges	3,837	
8	annotator	6,036	
9	download	610	
10	outperformed	2,196	
11	outperform	4,004	
12	publicly	2,428	
13	cosine	4,541	
18	outliers	495	
19	leveraging	896	
20	rerank	290	
21	tokenised	184	
22	mining	4,460	
23	mined	725	
24	females	250	
25	toolkit	3,529	
26	custom	556	
27	median	1,218	
28	overfitting	1,486	
29	leverage	1,527	
30	metrics	14,021	
35	sampled	2,759	
36	outperforming	697	
37	priors	1,596	
38	impacted	149	
39	treebanks	2,705	
40	lemmatized	1,046	
41	clinical	3,769	
42	lemmatize	173	
43	θ	1,481	
44	α	6,208	
45	revisit	184	
46	φ	1,023	
47	λ	1,311	

# Quelques outils grand public

- Gallicagram

Evolution de la fréquence relative par année à travers les collections de Gallica (et d'autres bibliothèques françaises)



# Méthodes d'analyse diachronique





# Traitement lexico-statistique de données diachroniques

- Constat : Le traitement d'une variable chronologique doit tenir compte du fait que celle-ci relève d'une échelle d'intervalle.
- Différentes visées méthodologiques :
  - 1) Regroupement de mots par profils chronologiques
    - Mesures de tendance globale (Herman & Kovář, 2013 ; Hilpert & Gries, 2009, 388-390) ;
    - Spécificités diachroniques (Salem 1988, 126-131 ; Lebart et al. 1998, 155-161)
    - Calcul du barycentre (Salem 1988, 134-139)
    - Curve clustering (Trevisani & Tuzzi, 2016)
  - 2) Diagnostiquer l'influence du facteur « temps » au sein d'une série préconstruite selon des critères chronologiques
    - AFC (Salem 1988, 110-118) ;

# Traitement lexico-statistique de données diachroniques

- Différentes visées méthodologiques (suite) :
  - 3) Périodisation automatique
    - Variability based neighbor clustering (VNC, Gries & Hilpert 2008)
  - 4) Identification de tournants
    - Modèles additifs généralisés (Gabrielatos et al. 2012)

# Présentation du Lexicoscope 2.0



# Le Lexicoscope 2.0

- Une conception proche du Sketch Engine

- ~ Concordances
- ~ Spécificités (*Keywords*)
- ~ Index hiérarchiques (*Word lists*)
- ~ Expressions polylexicales et ALR (*n-grams*)
- ~ Extraction de cooccurrents syntaxiques et de Wordsketches

Adjectifs modifieurs	Coordination : pivot et/ou...	Déterminants	Noms modifiant le pivot	Noms modifiés par le pivot
<a href="#">nu_ADJ</a> 14 299.26	<a href="#">poignard_NOUN</a> 10 207.04	<a href="#">de_DET</a> 11 235.12	<a href="#">côté_NOUN</a> 75 562.18	<a href="#">fil_NOUN</a> 68 761.89
<a href="#">court_ADJ</a> 25 281.92	<a href="#">lance_NOUN</a> 10 187.46		<a href="#">Albaicin_NOUN</a> 28 515.73	<a href="#">pointe_NOUN</a> 37 326.6
<a href="#">haut_ADJ</a> 27 253.72	<a href="#">hache_NOUN</a> 7 137.39		<a href="#">Nudd_NOUN</a> 27 473.33	<a href="#">maniement_NOUN</a> 13 277.88
<a href="#">grand_ADJ</a> 10 213.74	<a href="#">dague_NOUN</a> 16 128.97		<a href="#">Lachlann_NOUN</a> 22 396.29	<a href="#">pommeau_NOUN</a> 22 265.37
<a href="#">long_ADJ</a> 24 189.86	<a href="#">mousquet_NOUN</a> 4 85.49		<a href="#">poing_NOUN</a> 41 374.66	<a href="#">plat_NOUN</a> 10 213.74
<a href="#">seul_ADJ</a> 8 170.99	<a href="#">couteau_NOUN</a> 3 59.61		<a href="#">fourreau_NOUN</a> 32 354.69	<a href="#">poignée_NOUN</a> 21 201
<a href="#">brisé_ADJ</a> 13 160.72	<a href="#">regard_NOUN</a> 3 59.61		<a href="#">père_NOUN</a> 13 277.88	<a href="#">tranchant_NOUN</a> 9 192.36
<a href="#">double_ADJ</a> 11 151.95	<a href="#">pourpoint_NOUN</a> 3 55.8		<a href="#">Kradath_NOUN</a> 18 272.72	<a href="#">garde_NOUN</a> 26 183.1
<a href="#">tendu_ADJ</a> 7 149.61	<a href="#">pistolet_NOUN</a> 6 46.65		<a href="#">roi_NOUN</a> 12 256.5	<a href="#">nom_NOUN</a> 7 149.61
<a href="#">magique_ADJ</a> 7 149.61	<a href="#">couronne_NOUN</a> 2 42.74		<a href="#">guerrier_NOUN</a> 9 192.36	<a href="#">fragment_NOUN</a> 5 106.86
<a href="#">noir_ADJ</a> 7 149.61	<a href="#">Elmwood_NOUN</a> 2 42.74		<a href="#">Georges_NOUN</a> 8 170.99	<a href="#">combat_NOUN</a> 14 92.9

# Le Lexicoscope 2.0

## • Exemple d'extraction d'ALR

pointe_NOUN	6	<pre>&lt;c=DET,l=le,#1&gt;&amp;&amp; &lt;c=DET,l=le,#2&gt;&amp;&amp; &lt;c=NOUN,l=garde  pointe,#3&gt;&amp;&amp; &lt;c=NOUN,l=épée,#4&gt;&amp;&amp; &lt;c=PREP,#5&gt;&amp;&amp; &lt;c=PREP,l=de,#6&gt;:: (DETERM_DEF,3,1) (NMOD_POSIT1,3,4) (PREOBJ,3,5) (DETERM_DEF,4,2) (PREOBJ,4,6)</pre>	de_PREP le_DET garde pointe_NOUN de_PREP le_DET épée_NOUN	Il ouvrait la bouche et serrait les dents sur la pointe de l'épée qui les disjoignait.	25
fil_NOUN	6	<pre>&lt;c=DET,l=le,#1&gt;&amp;&amp; &lt;c=NOUN,l=fil,#2&gt;&amp;&amp; &lt;c=NOUN,l=épée,#3&gt;&amp;&amp; &lt;c=PREP,l=de,#4&gt;&amp;&amp; &lt;c=PREP,l=à,#5&gt;&amp;&amp; &lt;c=VERB,l=passer,#6&gt;:: (NMOD_POSIT1,2,3) (PREOBJ,2,5) (DETERM_DEF,3,1) (PREOBJ,3,4) (VMOD_POSIT1,6,2)</pre>	passer_VERB à_PREP fil_NOUN de_PREP le_DET épée_NOUN	Emportant l'odeur de la mort, les cris des femmes que les assassins passent au fil de l'épée, les silhouettes grotesques des enfants pendus au gibet de Bénarès, leur cocarde autour du cou.	25
dégainer_VERB	3	<pre>&lt;c=DET,l=le,#1&gt;&amp;&amp; &lt;c=NOUN,l=épée,#2&gt;&amp;&amp; &lt;c=VERB,l=brandir  dégainer,#3&gt;:: (DETERM_DEF,2,1) (NMOD_POSIT1,2,3)</pre>	le_DET épée_NOUN brandir dégainer_VERB	L'escarcelle derechef garnie, je quittai l'hôtel de Joyeuse tout ému de gratitude pour celle qui en était l'âme, regagnant mon logis sur mon Accla galopante, Giacomo et Miroul eux aussi montés et me flanquant, l'épée dégainée, pour ce que Cossolat m'avait avisé de ne plus cheminer à pied nuitamment par les rues de Montpellier, craignant pour moi quelque embûche et revanchement de la truanderie.\n	25
lance_NOUN	4	<pre>&lt;c=DET,l=le,#1&gt;&amp;&amp; &lt;c=DET,l=le,#2&gt;&amp;&amp; &lt;c=NOUN,l=dague  lance,#3&gt;&amp;&amp; &lt;c=NOUN,l=épée,#4&gt;:: (DETERM_DEF,3,1) (DETERM_DEF,4,2) (COORDITEMS 4 3)</pre>	le_DET épée_NOUN le_DET dague lance_NOUN	Il montra à Felix l'épée et la lance avec lesquels il avait tué plus de mille taureaux.	24

# Le Lexicoscope 2.0

- Mais avec une orientation « recherche »
  - ~ Développement d'expressions complexes et de patterns syntaxiques
  - ~ Étude contrastive de sous-corpus (partitions comme dans TXM)
  - ~ Construction automatique des requêtes complexes (requêtage basé sur l'exemple, Augustinus et al., 2016)
  - ~ Intermédiaire entre Sketch Engine (Kilgarriff et al., 2004) et Grew Match (Guillaume, 2021)

**Suggestion de requêtes avancées**

Cliquez sur une des suggestions ci-dessous pour générer automatiquement une requête avancée

Requête : <l=le,c=DET,#1>&&<l=tranchant,c=NOUN,#2>&&<l=de,c=PREP,#3>&&<l=son,c=DET,#4>&&<l=épée,c=NOUN,#5>:(DERM\_DEF,2,1) (DERM\_POSS,5,4) (NMOD\_POSIT1,2,5) (PREPOBJ,5,3)

```
graph TD; A(tranchant_NOUN) -- DERM_DEF --> B(le_DET); A -- NMOD_POSIT1 --> C(épée_NOUN); C -- PREPOBJ --> D(de_PREP); C -- DERM_POSS --> E(son_DET);
```

Fréquence estimée : 5

Exemple : ...

Fayolle pique son cheval dans la même direction, se baisse sur l'encolure, coupe la lanière de l'étrier avec [[le#1]] [[tranchant#2]] [[de#3]] [[son#4]] [[épée#5]].

# Corpus Imdiachro



# Corpus *Imdiachro*

- Compilé par Sascha Diwersy, consultable sur le Lexicoscope 2.0
- Echantillon d'articles du Monde de 1944 à 2015 (2 à 3 éditions par semaine)
- Traitement avec le pipeline Bonsai (Candito et al., 2010)
- Découpé par décades
- Imdiachro V2 : complété jusqu'à 2019. Traité avec Stanza.

Décade	Editions	Articles	Tokens
1944-49	410	23 089	8 767 430
1950-59	827	71 069	29 048 869
1960-69	785	85 160	37 517 057
1970-79	789	93 415	50 826 934
1980-89	810	91 392	55 097 518
1990-99	830	125 020	66 990 080
2000-09	827	116 726	64 344 690
2010-15	499	66 072	35 905 223
Total	5 777	671 943	348 497 801



# Fonctionnalité

# Chronogramme

## étude de cas



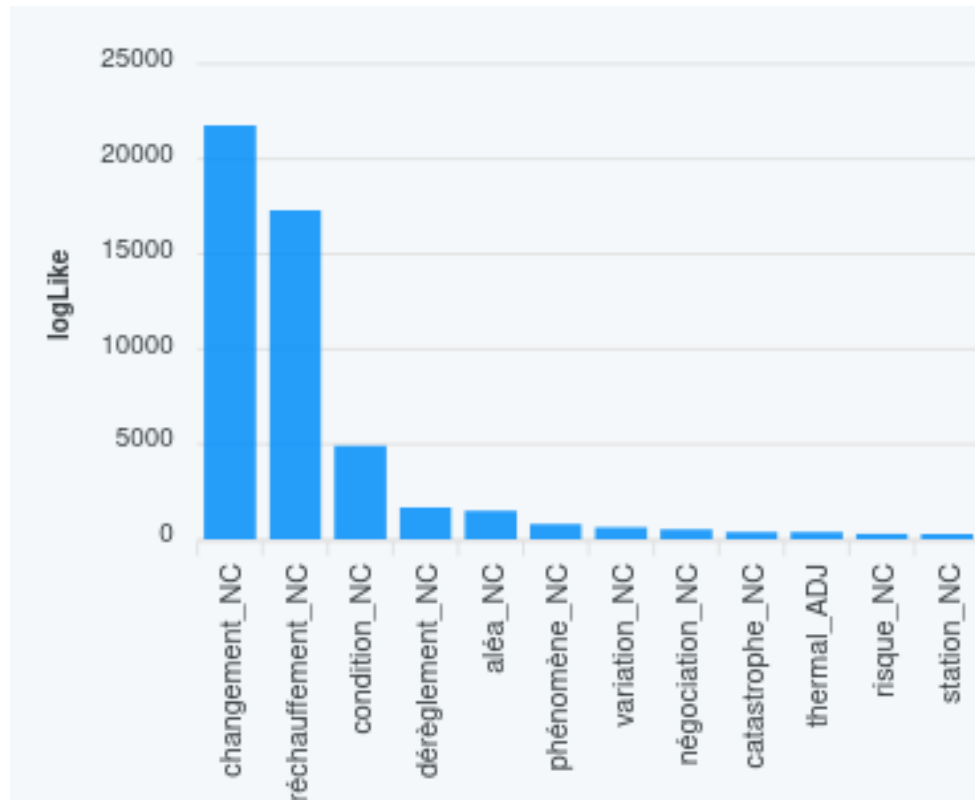
# Fonctionnalité chronogramme et étude de cas

- Étude de l'adjectif *climatique* (5117 occ.)
- Le chronogramme donne l'évolution des fréquences (absolues, relatives et spécificité) pour la taille d'empan choisie (mois, année, décade, etc.)



# Fonctionnalité chronogramme et étude de cas

- Cooccurents de *climatique*

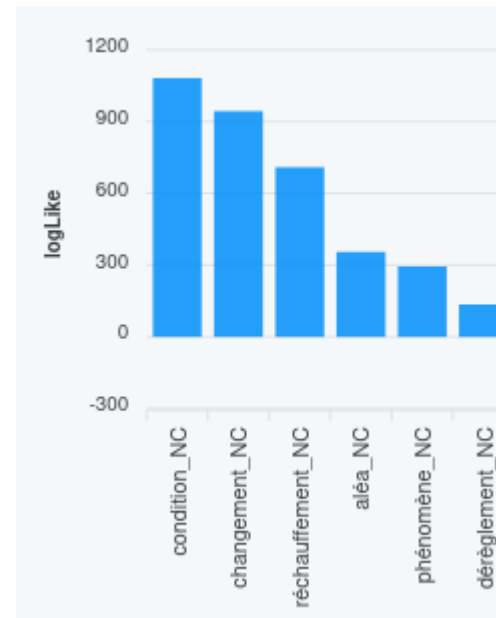
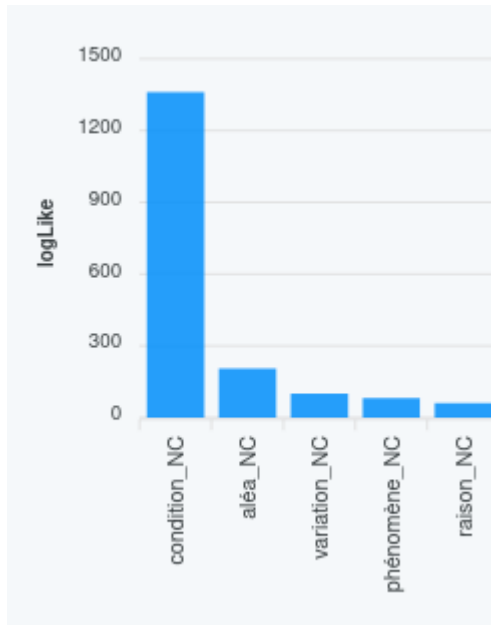


# Fonctionnalité chronogramme et étude de cas

- Deux séries de cooccurrents :

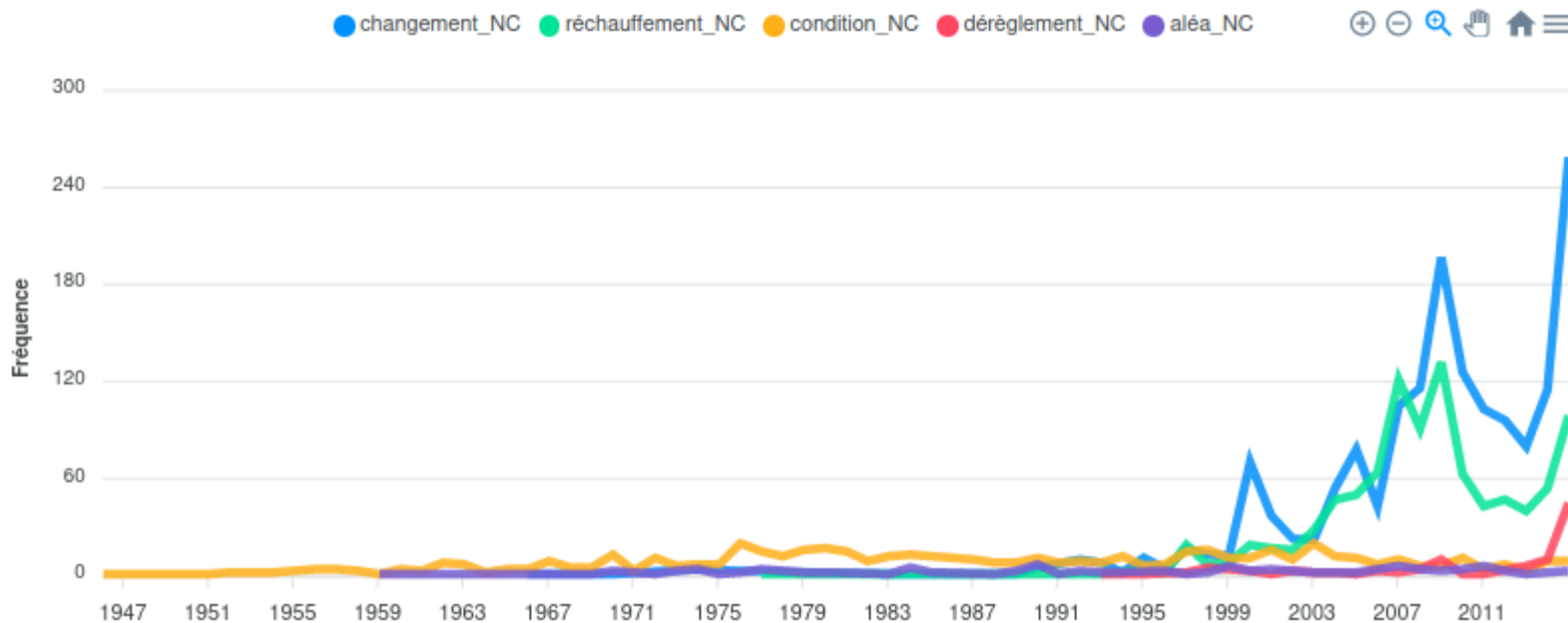
*changement, réchauffement, condition, dérèglement, aléa, phénomène, variation, négociation, catastrophe, thermal, risque, station*

une étude par décennie montre une bascule entre les années 1980 et 1990



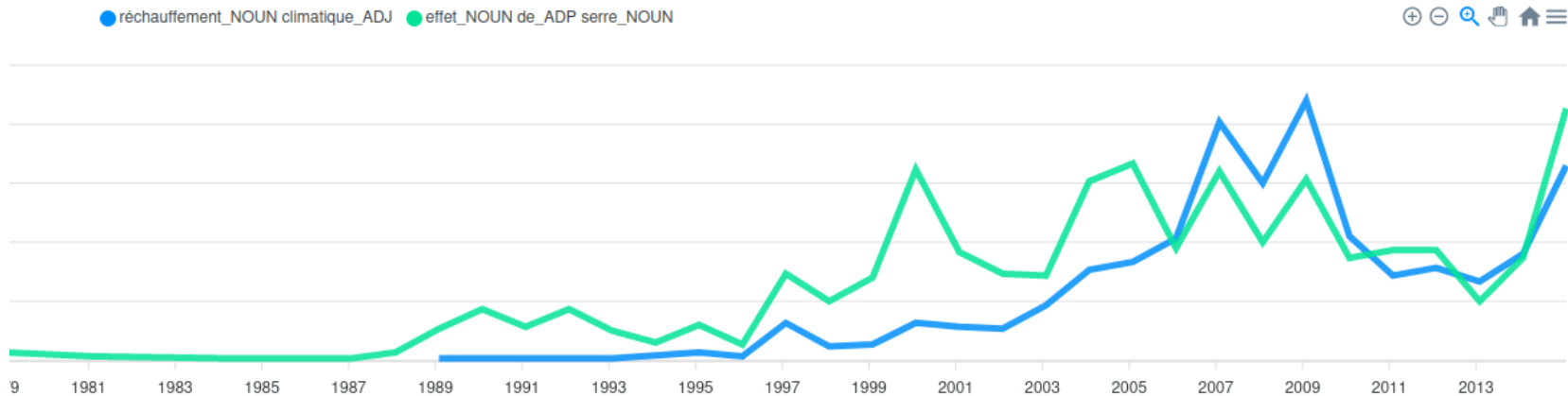
# Fonctionnalité chronogramme et étude de cas

- Possibilité de croiser chronogramme et cooccurrents



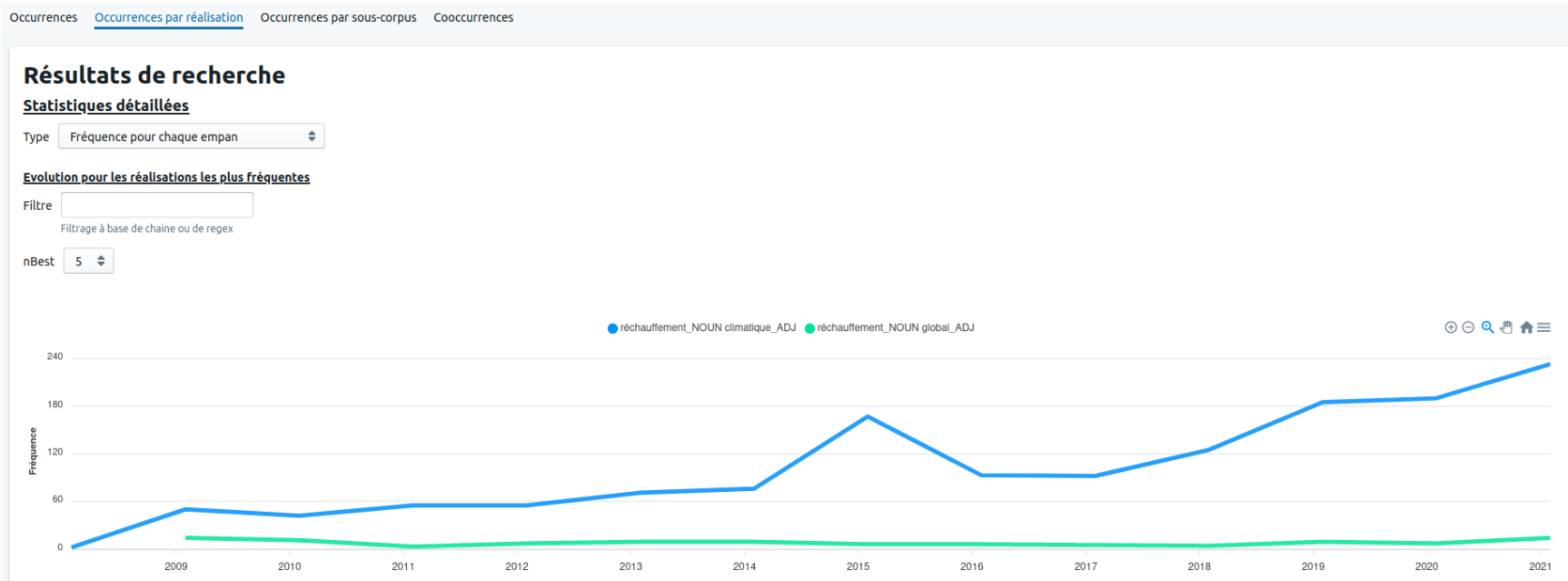
# Fonctionnalité chronogramme et étude de cas

- Possibilité de croiser chronogramme et réalisations



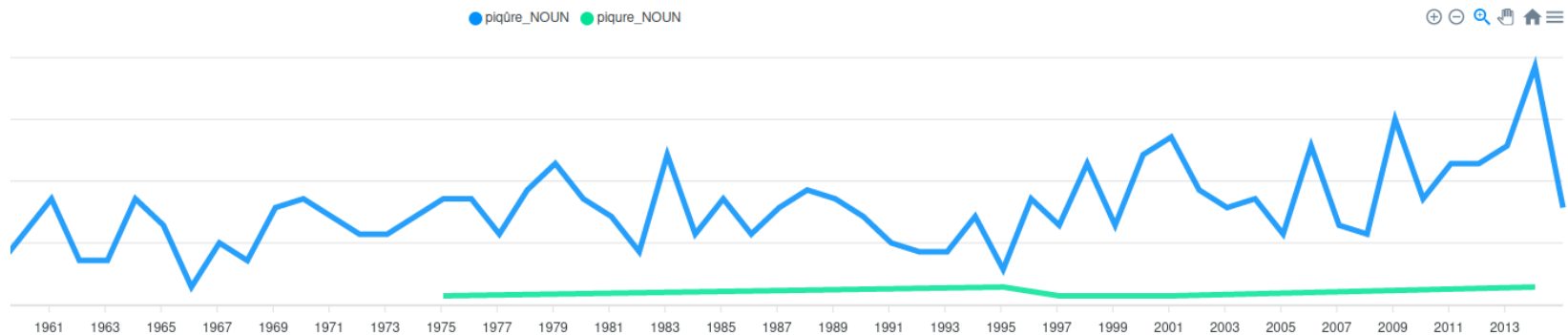
# Fonctionnalité chronogramme et étude de cas

- Possibilité de croiser chronogramme et réalisations



# Fonctionnalité chronogramme et étude de cas

- Possibilité de croiser chronogramme et réalisations : ex. d'évolution orthographique





# Fonctionnalité chronogramme et étude de cas

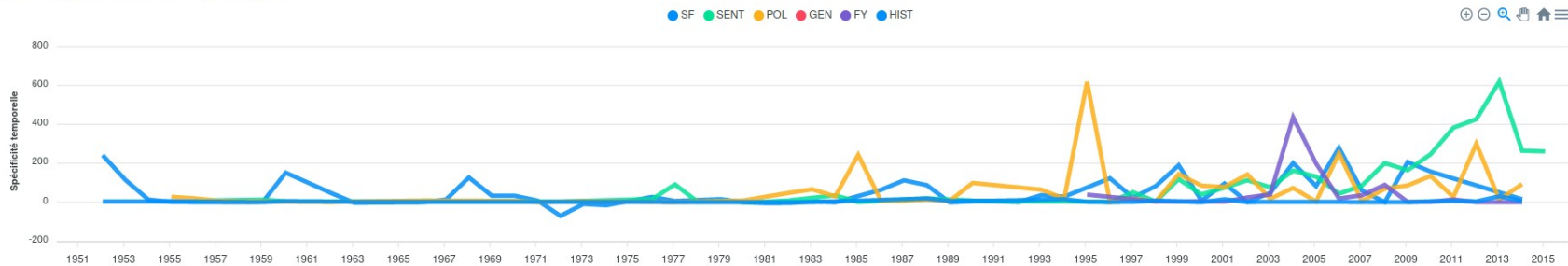
- Possibilité de croiser chronogramme et sous-corpus

## Résultats de recherche

### Statistiques détaillées

Type

#### Evolution pour les réalisations les plus fréquentes



évolution des occurrences de *écran*  
(spécificité temporelle)

# Fonctionnalité chronogramme et étude de cas

- Affichage des tendances

## Cooccurrents en augmentation ↗

- **changement\_NOUN** : Première attestation= 1966 (freq= 1613, coeff= 20.69)
- **réchauffement\_NOUN** : Première attestation= 1989 (freq= 970, coeff= 11.53)
- **risque\_NOUN** : Première attestation= 1955 (freq= 43, coeff= 3.86)
- **modèle\_NOUN** : Première attestation= 1993 (freq= 37, coeff= 2.5)
- **archive\_NOUN** : Première attestation= 1979 (freq= 10, coeff= 2.33)
- **aléa\_NOUN** : Première attestation= 1959 (freq= 112, coeff= 1.88)
- **problème\_NOUN** : Première attestation= 1970 (freq= 22, coeff= 1.83)
- **donnée\_NOUN** : Première attestation= 1974 (freq= 15, coeff= 1.8)
- **catastrophe\_NOUN** : Première attestation= 1973 (freq= 45, coeff= 1.62)
- **raison\_NOUN** : Première attestation= 1961 (freq= 31, coeff= 1.55)

## Cooccurrents en régression puis en augmentation ↘↗

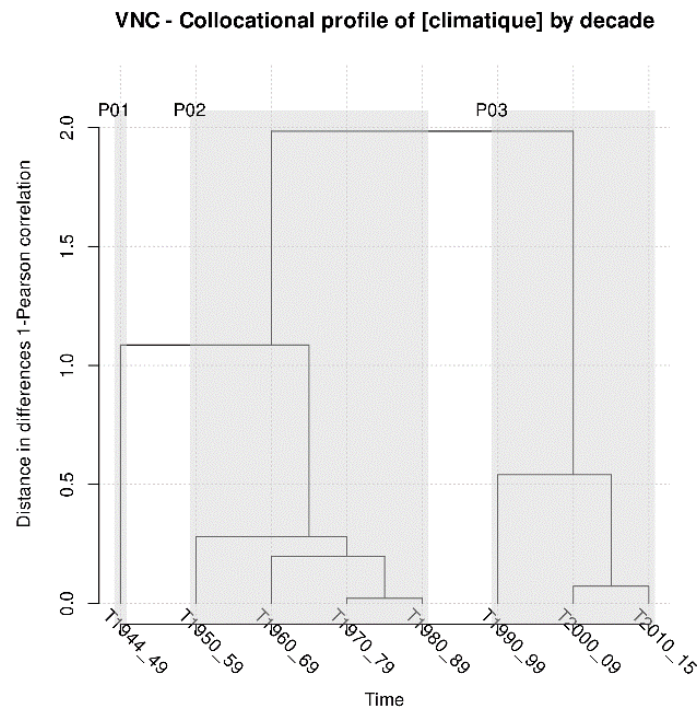
- **conférence\_NOUN** (freq= 11, coeff= 10)
- **désordre\_NOUN** (freq= 12, coeff= 2)
- **politique\_NOUN** (freq= 43, coeff= 1.82)

# Perspectives



# Perspectives

- Ajout de mesures de tendance
  - ~ Estimateur Theil Sen (Herman & Kovář, 2013)
  - ~ Correlation des rangs de Kendall (Hilpert & Gries, 2009, 388-390)
    - exemple *climatique* :  $\tau = 0.8571$ , p-valeur = 0.0044 → tendance à la hausse à travers la série chronologique
- Périodisation automatique
  - ~ Classification ascendante hiérarchique par contiguïté - CAHC (Diwersy et al., 2017) ou VNC (Gries & Hilpert, 2008)



Merci de votre attention !

# Références

- Augustinus, L., Vandeghinste, V., Vanallemeersch, T. (2016). Poly-GrETEL: Cross-Lingual Example-based Querying of Syntactic Constructions. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, pages 3549–3554, Portorož, Slovenia. ELRA.
- Candito M., Nivre J., Denis P. and Henestroza Anguiano E. (2010). Benchmarking of Statistical Dependency Parsers for French. In *Proceedings of The 23rd International Conference on Computational Linguistics (COLING 2010)*, Beijing, China.
- Coavoux M., Denoyelle, C. Kraif, O., Sorba J. (2022) Phraséologie du roman médiéval en prose. In *Actes du Colloque international DIACHRO X, Le français en diachronie*, Paris, Sorbonne Université, 17-19 janvier 2022.
- Diwersy, S., Falaise A., Lay, M.-H., Souvay, G. (2017) Ressources et méthodes pour l'analyse diachronique. *Langages 2017/2* (N° 206), 21-44.
- Evert, S. (2007). Corpora and collocations. In A. Lüdeling & M. Kytö (Eds.), *Corpus Linguistics. An International Handbook*. Berlin: Mouton de Gruyter, pp. 1212–1248.
- Gabrielatos, C., McEnery, T., Diggle, P. J. & Baker, P. (2012). The peaks and troughs of corpus-based contextual analysis. *International Journal of Corpus Linguistics* 17(2). 151–175. <https://doi.org/10.1075/ijcl.17.2.01gab>.

# Références

- Grobol, L. and Crabbé, B. (2021). Analyse en dépendances du français avec des plongements contextualisés. In *Actes de la 28e Conférence sur le Traitement Automatique des Langues Naturelles*, pages 106–114. Lille, France, June.
- Grobol, L., Regnault, M., Oretiz Suarez P., Sagot, B., Romary, L., Crabbé, B. (2022). BERTrade: Using Contextual Embeddings to Parse Old French. LREC 2022, Marseille : 1104-1113.
- Gries S. Th. and Hilpert M. (2008). The identification of stages in diachronic data: variability-based neighbour clustering. *Corpora*, 3, 59–81.
- Guillaume B. (2021). Graph Matching and Graph Rewriting: GREW tools for corpus exploration, maintenance and conversion. *EACL (System Demonstrations) 2021*: 168-175.
- Herman O. and Kovář V. (2013). Methods for Detection of Word Usage over Time. In *Seventh Workshop on Recent Advances in Slavonic Natural Language Processing, RASLAN 2013*, 79–85.
- Kilgariff A., Rychly P., Smrz P., Tugwell D. (2004). The Sketch Engine". *Proceedings of the Eleventh EURALEX International Congress*, pages 105–116. Lorient, France. Lebart L., Salem A. and Berry L. (1998). *Exploring Textual Data*. Kluwer Academic Publisher.
- Salem A. (1988). Approches du temps lexical. *Mots*, 17, 105–143. <https://doi.org/10.3406/mots.1988.1401>
- Trevisani, M. & Tuzzi, A. (2016). Analisi di dati testuali cronologici in corpora diacronici: effetti della normalizzazione sul curve clustering. In *JADT 2016: 13ème Journées Internationales d'Analyse Statistique des Données Textuelles*. Nice, France. <http://lexicometrica.univ-paris3.fr/jadt/jadt2016/01-ACTES/82630/82630.pdf>.