

Complexification des données et des techniques en linguistique : contributions du TAL aux solutions et aux problèmes

Soutenance en vue de l'obtention de
l'Habilitation à Diriger des Recherches
Spécialité : Linguistique



Ludovic Tanguy



Anne Condamines
Benoît Habert
Marie-Paule Péry-Woodley
Pascale Sébillot
Mathieu Valette
François Yvon

CLLE - Université de Toulouse promotrice
ICAR - ENS de Lyon président
CLLE - Université de Toulouse examinatrice
IRISA - INSA de Rennes examinatrice
ERTIM - INaLCO rapporteur
LIMSI - Université de Paris Sud rapporteur

Plan de la présentation

- 1 Vue d'ensemble
- 2 Outillage de l'exploitation des corpus
- 3 Utilisation du Web comme corpus
- 4 Visualisation et interprétation des données linguistiques
- 5 Apprentissage artificiel et articulation avec la linguistique
- 6 Discussion sur le TAL en linguistique
 - Bienfaits et méfaits
 - Enseignement du TAL en sciences du langage

- 1 Vue d'ensemble
- 2 Outillage de l'exploitation des corpus
- 3 Utilisation du Web comme corpus
- 4 Visualisation et interprétation des données linguistiques
- 5 Apprentissage artificiel et articulation avec la linguistique
- 6 Discussion sur le TAL en linguistique
 - Bienfaits et méfaits
 - Enseignement du TAL en sciences du langage

Réflexions sur l'évolution du TAL et de la linguistique outillée

Augmentation de la masse et de la diversité des données utilisées

Développements de grands corpus, exploitation du Web

Montée en puissance des outils informatiques

Incontournables pour acquérir, annoter et interroger ces données massives

Complexification de l'outillage

Modes d'accès et d'interrogation, méthodes statistiques, techniques d'apprentissage artificiel et de fouille de données

Risques liés à cette évolution

Éloignement des données, creusement du fossé entre la linguistique et l'informatique au sein du TAL

Ma position dans tout ça

Jouer un rôle de passeur

Médiation entre les préoccupations empiriques et théoriques sur les données langagières et les techniques informatiques permettant d'y répondre

Promouvoir les techniques informatisées

Proposer de nouvelles méthodes pour exploiter les données, valoriser les annotations, mettre à portée des utilisateurs les techniques d'analyse quantitative

Remettre de la linguistique dans le TAL

Utiliser les connaissances et les ressources linguistiques dans des applications qui ont tendance à les ignorer

Formation

Diplôme d'ingénieur de l'École Nationale Supérieure des Télécommunications de Bretagne (1993)

Spécialité : Intelligence Artificielle et Sciences Cognitives

DEA d'informatique de l'Université de Rennes 1 (1993)

Mémoire : *Apprentissage de structures hiérarchiques de représentation des connaissances par algorithmes génétiques*

Doctorat en informatique de l'Université de Rennes 1 (1997)

Traitement automatique de la langue naturelle et interprétation : contribution à l'élaboration d'un modèle informatique de la sémantique interprétative

Laboratoire d'Intelligence Artificielle et Sciences Cognitives (ENSTBr)

Directeurs : J.-P. Barthélemy et I. Kanellos

Postes occupés

1996–98

ATER en Informatique à l'Université de Bretagne Occidentale (faculté des lettres et sciences humaines)

1998–99

Chercheur en linguistique informatique à l'ISSCO (Université de Genève)

1999–

Maître de conférences en sciences du langage à l'Université de Toulouse 2
Département des sciences du langage
Laboratoire CLLE-ERSS

Grandes lignes de mes travaux

Linguistique de corpus

Acquisition, annotation, fouille et analyse de données textuelles

Linguistique appliquée

Proposition de méthodes spécifiques d'investigation des données

Aspects linguistiques du TAL

Recherche et extraction d'information, classification de textes

Collaborations

Autres chercheurs en linguistique

- Morphologie
- Lexique
- Discours

Autres disciplines académiques

- Informatique
- Psychologie
- Sociologie
- Médecine

Entreprises

- Développement d'applications en TAL
- Applications visant les documents professionnels

Enseignements

Traitement automatique des langues

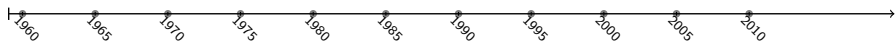
- Programmation pour le TAL
- Développement d'applications

Outillage de la linguistique

- Méthodologie de l'utilisation des données
- Manipulation de textes numériques
- Statistique pour la linguistique

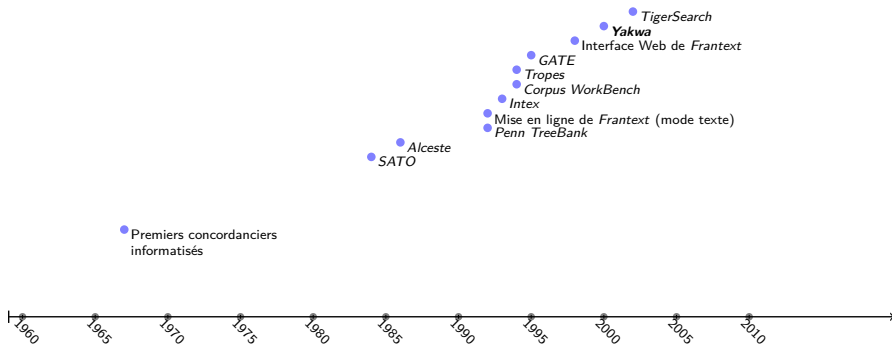
- 1 Vue d'ensemble
- 2 Outillage de l'exploitation des corpus**
- 3 Utilisation du Web comme corpus
- 4 Visualisation et interprétation des données linguistiques
- 5 Apprentissage artificiel et articulation avec la linguistique
- 6 Discussion sur le TAL en linguistique
 - Bienfaits et méfaits
 - Enseignement du TAL en sciences du langage

Frise chronologique de l'outillage des corpus



Frise chronologique de l'outillage des corpus

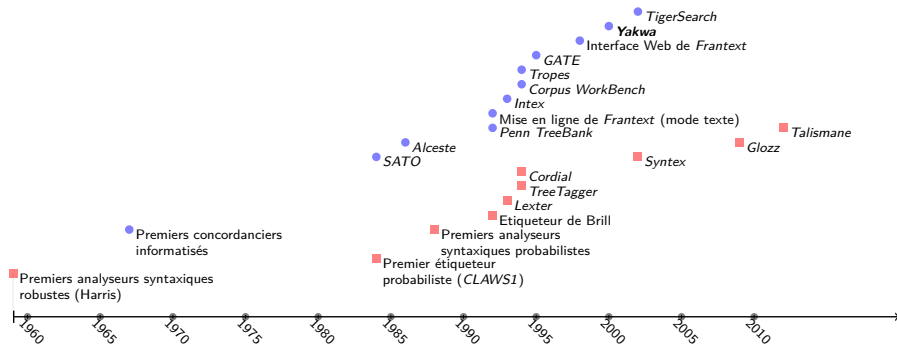
● Interrogation de corpus



Frise chronologique de l'outillage des corpus

● Interrogation de corpus

■ Étiquetage de corpus

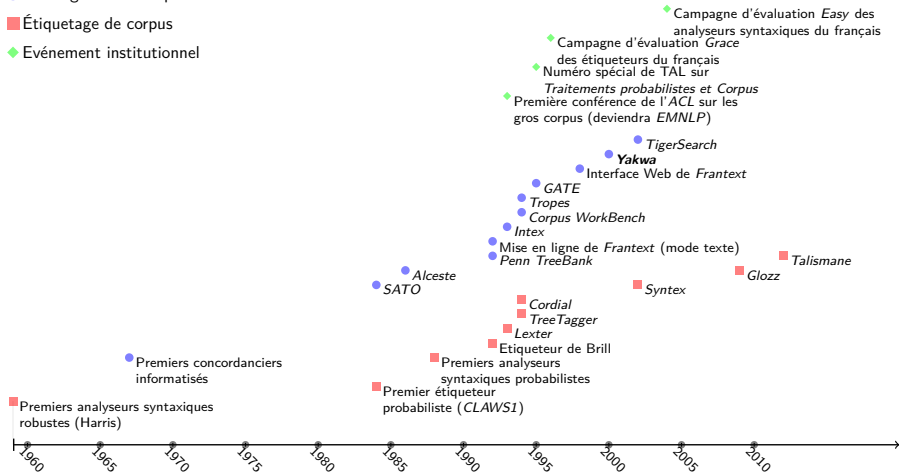


Frise chronologique de l'outillage des corpus

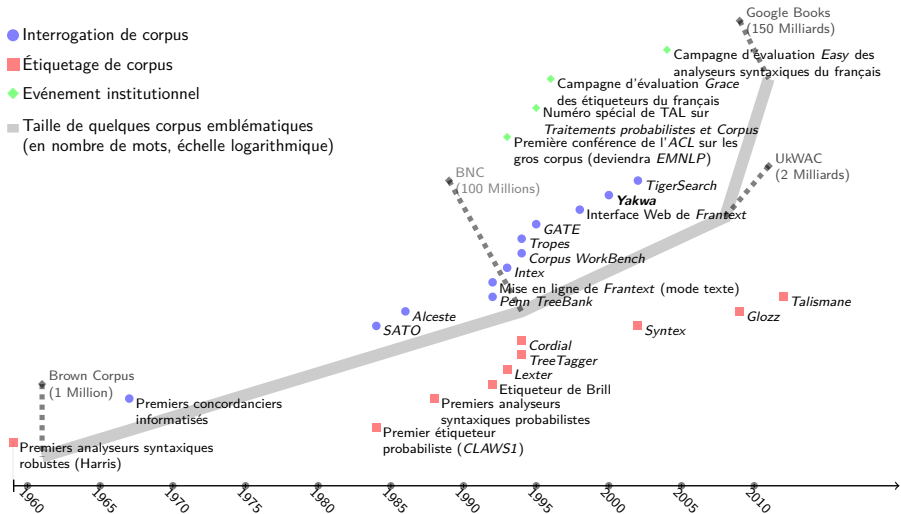
● Interrogation de corpus

■ Étiquetage de corpus

◆ Événement institutionnel



Frise chronologique de l'outillage des corpus



Mes travaux

L'ERSS : un terrain propice

Comment utiliser les corpus (et pas pourquoi) ?

Des besoins et des envies

Passage du texte nu au texte annoté (morphosyntaxiquement, puis syntaxiquement)

Mon rôle

- Comprendre les besoins et les enjeux
- Concilier les problématiques linguistiques et les considérations techniques
- Proposer de nouveaux modes d'accès aux corpus

Yakwa

Un concordancier sur corpus étiquetés morphosyntaxiquement

Conception du moteur et de l'interface.

Définition et utilisation de patrons morphosyntaxiques, notamment pour :

- Les énoncés définitoires (J. Rebeyrolle)
- Les expressions du dysfonctionnement (P. Vergely)
- Divers marqueurs de relations sémantiques (A. Condamines, N. Aussenac-Gilles)

Bilan

- Un outil rapidement pris en main par les linguistes (interface)
- Qui a permis de démontrer l'avantage (et le faible coût) des patrons morphosyntaxiques
- Un mode d'indexation ne permettant pas des annotations par l'utilisateur
- Une machinerie trop lourde à maintenir

Travaux autour de Syntex

Syntex

- Disponibilité d'un analyseur syntaxique opérationnel en dépendances (Bourigault 2007)
- Développé pour l'extraction terminologique et la construction de ressources sémantiques distributionnelles

Mon rôle : développement d'outils d'interrogation de corpus

- Changement de formalisme : arbres syntaxiques à partir des dépendances
- Utilisation des outils de treebanking (*TigerSearch*)
- Promotion auprès de collègues syntacticiens (ERSS, BCL)

Bilan pour les corpus annotés syntaxiquement

Un échec à causes multiples

- Raisons techniques : complexité des modes d'interrogation
- Raisons culturelles : représentation trop proche des grammaires génératives
- Problèmes de cible : les phénomènes syntaxiques complexes les plus recherchés sont les plus difficiles à traiter automatiquement

Un surcoût trop élevé

- La puissance de calcul nécessaire ne suit pas la taille des corpus
- Les relations simples sont aussi bien traitées par des patrons de surface

Réflexions sur l'outillage des corpus

Adapter l'outillage au besoin

- La simple recherche de séquences peut se contenter des concordanciers simples
- Garder les moteurs complexes pour les développements en TAL
- Le gain d'une utilisation directe des analyseurs syntaxiques pour la linguistique descriptive reste à démontrer

Ne pas négliger les compétences nécessaires

- Formalisme de requête (des expressions régulières aux langages complexes)
- Connaissances sur les traitements appliqués aux corpus
- Acceptation et gestion des erreurs d'étiquetage
- Stratégies spécifiques apprises sur le tas

De nouveaux palliers de complexification

Les annotations multiples

- Cumul possible par les formalismes XML déportés
- Une exploitation très complexe, même pour le TAL actuel

Les annotations structurelles

- Organisation logique des documents : comment les rendre facilement exploitables ?
- Structures discursives : comment les représenter et les rendre accessibles ?

Ne plus décorrélérer l'annotation de l'interrogation

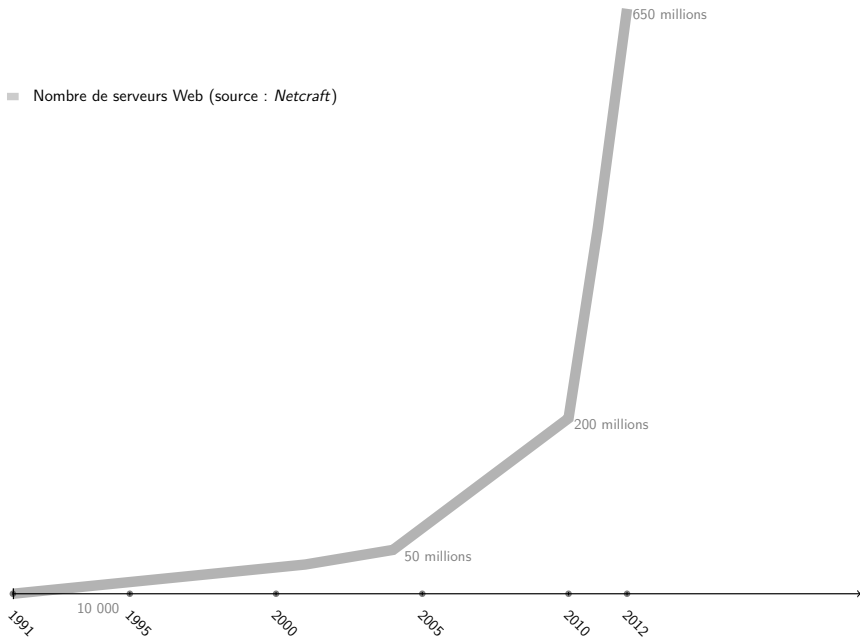
- L'exemple d'Annodis : projection d'indices pour assister l'annotation
- L'investigation des traits linguistiques : définition, projection et analyse de marqueurs

- 1 Vue d'ensemble
- 2 Outillage de l'exploitation des corpus
- 3 Utilisation du Web comme corpus**
- 4 Visualisation et interprétation des données linguistiques
- 5 Apprentissage artificiel et articulation avec la linguistique
- 6 Discussion sur le TAL en linguistique
 - Bienfaits et méfaits
 - Enseignement du TAL en sciences du langage

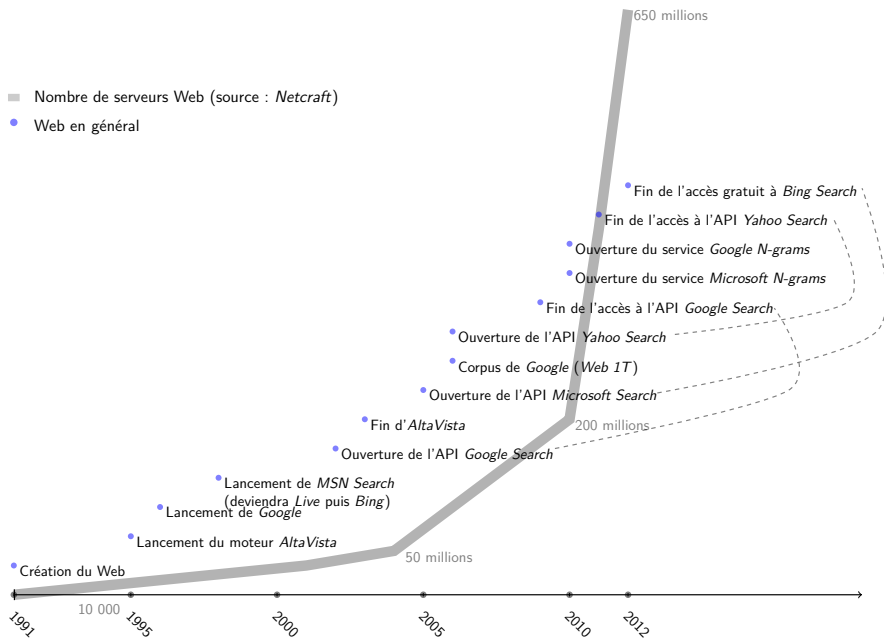
Chronologie des usages du Web en linguistique et en TAL



Chronologie des usages du Web en linguistique et en TAL



Chronologie des usages du Web en linguistique et en TAL

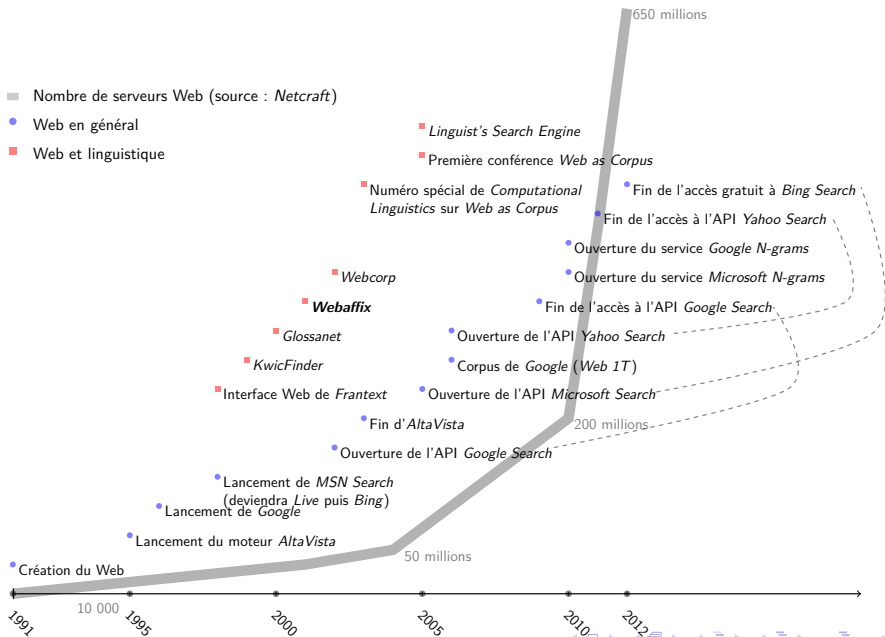


Chronologie des usages du Web en linguistique et en TAL

■ Nombre de serveurs Web (source : Netcraft)

● Web en général

■ Web et linguistique



Vue d'ensemble des usages

Web as corpus / Web for corpus

- *As corpus* : recherche d'attestations *via* les moteurs de recherche
- *For corpus* : construction de corpus par moissonnage de pages Web

Le Web : un corpus ou pas ?

- **Aucun contrôle sur la nature, le nombre, le type des documents, le statut des énoncés**
- **Une masse de documents d'une taille inégalée, prête à l'emploi, facile d'accès, variée, comprenant des genres propres et des productions spontanées**
- Au final, un phénomène impossible à ignorer, mais à aborder de face

Principales utilisations en linguistique et en TAL

- Recherche d'énoncés (généralement atypiques)
- Comparaison des fréquences (approximative)
- Acquisition de ressources
- Accès à la masse indifférenciée pour des approches par apprentissage

Mes travaux

Principales réalisations

Repérage et analyse des créations lexicales par suffixation (*Webaffix*, avec N. Hathout)

- Étude de suffixes ciblés (*-esque*, *-este*, *-able*, etc.)
- Acquisition de lexique
- Étude de la concurrence suffixale

Utilisation ponctuelle de gros corpus ou de données issus du Web

Évolution des techniques au fil des décisions des moteurs de recherche

- *Hacking* direct des premiers moteurs (dont *AltaVista*)
- Utilisation des *API* proposées par les principaux moteurs de recherche
- Utilisation de corpus statiques (*Exalead*, *WAC*) ou de ressources préconstruites (N-grammes de *Google*)
- Tentatives malheureuses de moissonnage artisanal (avec F. Sajous)

Quelques grands succès

Suffixation en *-este* (avec M. Plénat)

Source	Date	Dérivés	Nb
Pichon	1940	<i>silvio-pelliqueste</i>	1
Plénat	1997	<i>astraqueste, grandiloqueste, dingueste</i>	4
Webaffix (1 ^{re} campagne)	2002	<i>dingueste, hageste, iagueste, langueste, mar- queste, pragmatiqueste, punkeste, titaniqueste</i>	12
Webaffix (2 ^e campagne) complété par Plénat	2004	<i>algueste, bangueste, big-bangueste, blagueste, blogueste, borgueste, bouledogueste, bouy- guesta, cirqueste, darkeste, dukeste, fiasqueste, fliqueste, gagueste, gangueste, etc.</i>	44

Extension de *Verbaction*

Initialement (Hathout et al. 2002) 6 471 couples noms-verbos

Avec *Webaffix* et validation manuelle, 9 393 couples

Dérivés en *-able*

Dans les dictionnaires, 1400 adjectifs ; avec *Webaffix*, plus de 5000

Quelques réflexions

Les raisons et la nature du succès

- Un accès très rapide aux données tant que l'on cible des unités lexicales
- Les phénomènes rares demandent une masse de données inaccessible ailleurs (moins d'une page Web sur 200 pour trouver un nouveau déverbal d'action, soit moins d'un mot sur 300 000)
- L'accès à des données extensives permet de renouveler les points de vue en identifiant de nouveaux types d'emplois

L'envers du décor

- Même automatisé, le dépouillement des données recueillies est fastidieux, et nécessite de développer des compétences spécifiques
- Les évolutions constantes des moteurs et des technologies sont décourageantes
- Les moteurs de recherche ne sont actuellement plus du tout coopératifs (fermeture ou restriction drastique des API)

Quelques remarques sur les attestations trouvées sur le Web

Sources de bruit

- Détectables : orthographe et (dys-)typographie, noms propres, langues proches, codes informatiques, etc.
- Indétectables : scripteur non natif, traduction automatique, textes artificiels, contextes métalinguistiques

Autres phénomènes locaux

- Jeux de langage

*Hier, notre **candidateuse socialistique** à la **présidenciation** de la **Républiquitude** était en **baladage** chez les Chinetoques.*

- Rafales suffixales et amorçage

*Niveau **élégance prestance classance** et **distinctance**, je reste sur mes positions*

Guider l'interprétation des occurrences

Viser des contextes productifs

Les campagnes passées ont permis d'identifier des lieux de haute densité en création suffixale :

doctissimo, aufeminin, etc.

Caractériser les contextes

- Dans la lignée des travaux sur les genres du Web
- Développement de caractérisations calculables permettant une aide au jugement sur le statut des créations lexicales
- Utilisable à la volée, ou bien pour identifier des zones prometteuses (investir dans la prospection puisque le forage est coûteux)

- 1 Vue d'ensemble
- 2 Outillage de l'exploitation des corpus
- 3 Utilisation du Web comme corpus
- 4 Visualisation et interprétation des données linguistiques**
- 5 Apprentissage artificiel et articulation avec la linguistique
- 6 Discussion sur le TAL en linguistique
 - Bienfaits et méfaits
 - Enseignement du TAL en sciences du langage

État des besoins en exploration des données langagières

Des données exigeantes

- En grandes quantités
- Des annotations multiples et variées
- Des structures linguistiques complexes

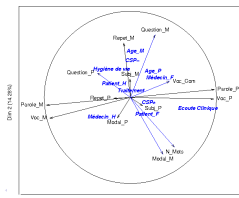
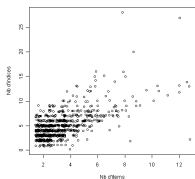
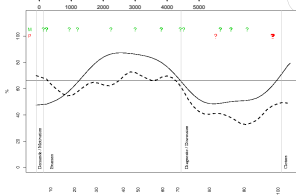
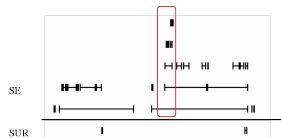
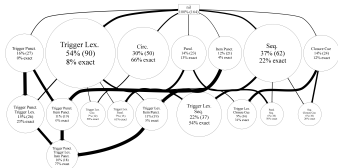
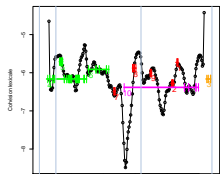
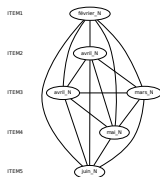
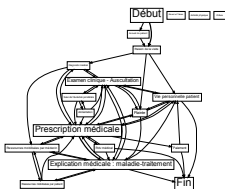
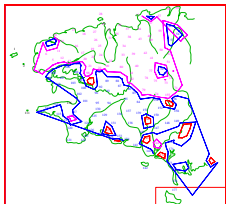
Des questions de divers types

- Avoir une vue d'ensemble (pas de question précise)
- Identifier des régularités
- Isoler des cas particuliers
- Croiser des phénomènes de types différents

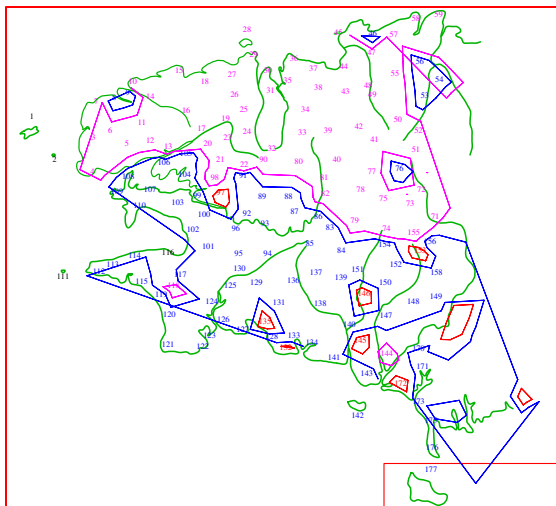
Un travail collectif

- Des questions souvent interdisciplinaires, avec des regards croisés
- Instauration d'un dialogue autour des données avec leurs concepteurs et leurs utilisateurs

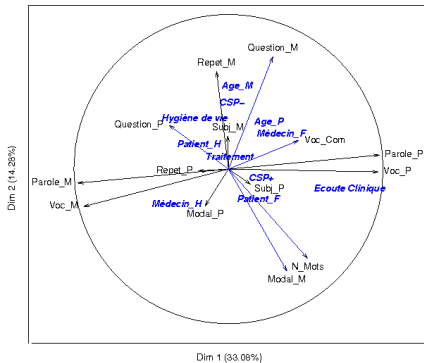
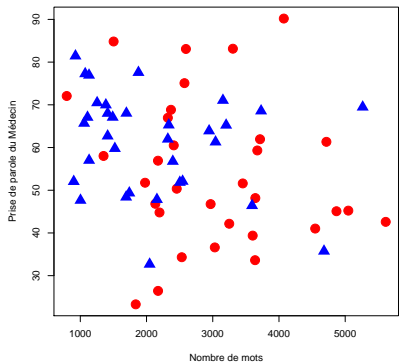
Quelques exemples de visualisations



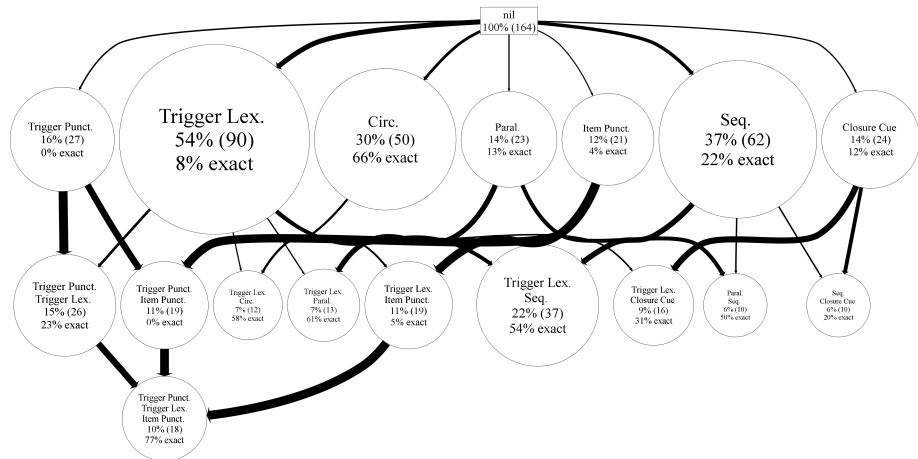
Usages de la visualisation (1) : mode d'interprétation premier avec un référentiel explicite



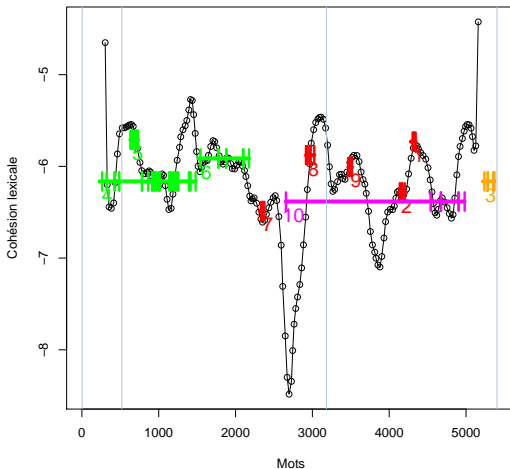
Usage de la visualisation (2) : avant et après l'analyse statistique



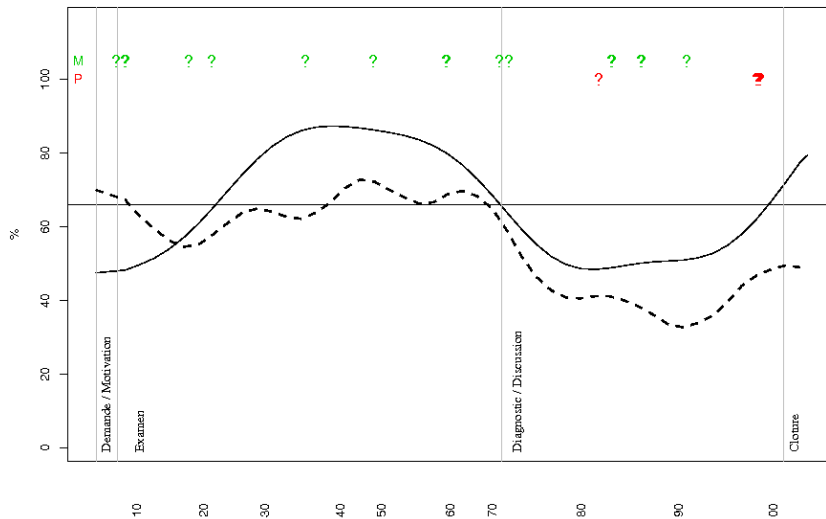
Usages de la visualisation (3) : avoir un point de vue global sur des relations isolées



Usages de la visualisation (4) : croiser des données de natures différentes



Usage de la visualisation (5) : visualiser la disposition dans les textes



Petit bilan

Du point de vue du concepteur

- De gros efforts techniques de conception pour un résultat indéfiniment perfectible
- Des représentations générées à l'intuition, pas toujours adaptées et parfois sources d'erreurs d'interprétation
- Chaque nouveau cas appelle un nouveau mode de visualisation

Du point de vue des utilisateurs

- Des prises en main plus ou moins faciles
- Un bon moyen de voir apparaître des phénomènes globaux inattendus
- Au final de bons supports d'interprétation et de communication

Vers la visualisation des données langagières

Des exemples à suivre

- La visualisation de l'information (*InfoVis*) est désormais un champ de recherche à part entière, intéressant toutes les disciplines
- On y associe la question de l'articulation avec l'investigation quantitative des données (*Visual Analytics*)
- Les coûts élevés en développement d'outils doivent être mutualisés, notamment en ce qui concerne les représentations interactives
- Les avancées se font au sein de communautés partageant des objets d'étude

Diffuser les pratiques

- Le meilleur moyen d'aborder la complexité croissante des données
- Multiplier les expériences pour faire émerger des méthodes productives

Limites

Rester prudent face aux difficultés

- Des dangers de mésinterprétation
- Garder le lien avec les données originelles (minimiser les traitements avant la visualisation)
- Ne pas sous-estimer les connaissances nécessaires :
 - Connaissance des données (initiales et annotations)
 - Expérience des procédures de visualisation
 - Compétences en traitement des données

Un travail intermédiaire et rarement une finalisation

- Un travail qui reste souvent dans l'ombre
- Impossibilité de démontrer un fait par une seule visualisation
- Nécessité de passer aux méthodes statistiques

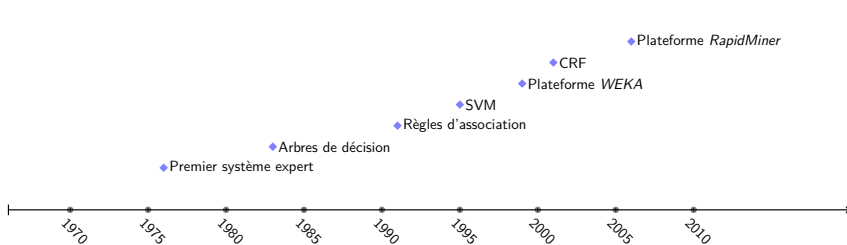
- 1 Vue d'ensemble
- 2 Outillage de l'exploitation des corpus
- 3 Utilisation du Web comme corpus
- 4 Visualisation et interprétation des données linguistiques
- 5 Apprentissage artificiel et articulation avec la linguistique**
- 6 Discussion sur le TAL en linguistique
 - Bienfaits et méfaits
 - Enseignement du TAL en sciences du langage

Frise chronologique des méthodes quantitatives en TAL



Frise chronologique des méthodes quantitatives en TAL

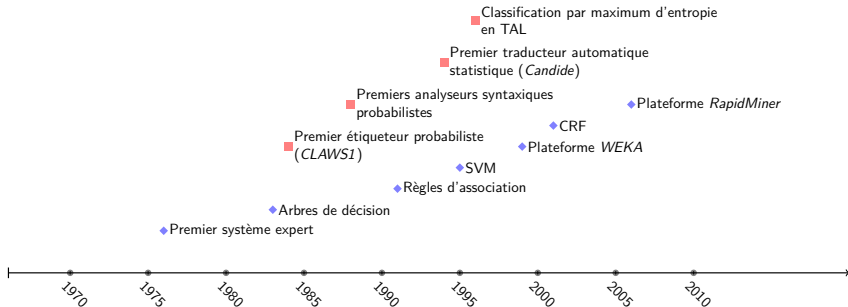
◆ Outils génériques



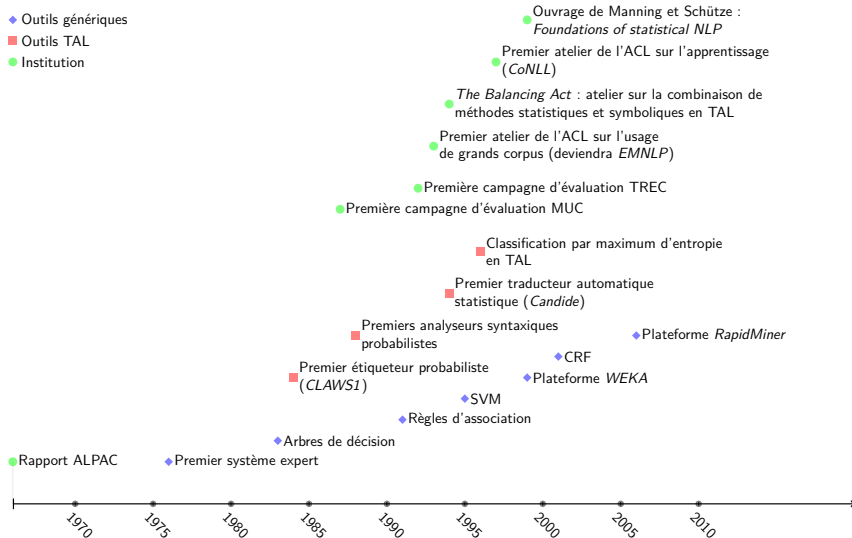
Frise chronologique des méthodes quantitatives en TAL

◆ Outils génériques

■ Outils TAL



Frise chronologique des méthodes quantitatives en TAL



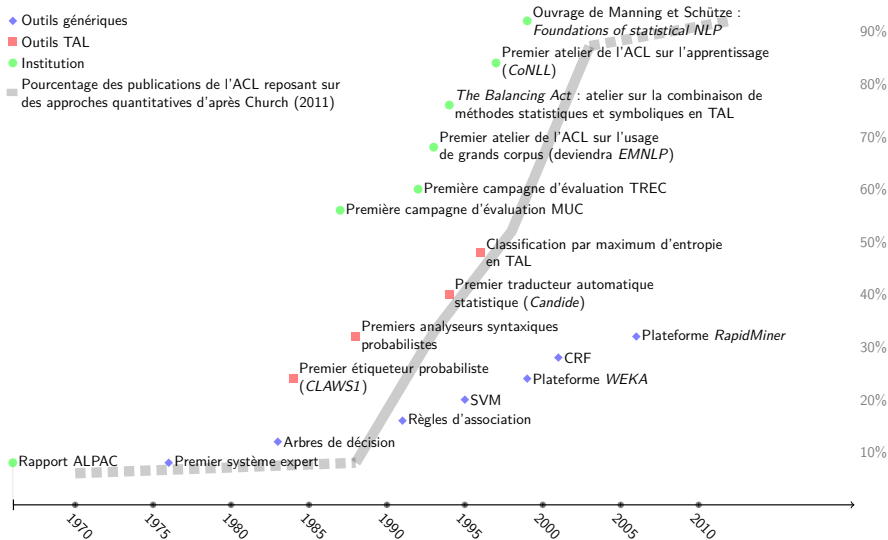
Frise chronologique des méthodes quantitatives en TAL

◆ Outils génériques

■ Outils TAL

● Institution

■ Pourcentage des publications de l'ACL reposant sur des approches quantitatives d'après Church (2011)



Un changement radical en TAL (1)

Exemple de problème : l'analyse de références bibliographiques

Tanguy, Ludovic (2015). Overbearingly Supervised Techniques for Very Small Corpora. Journal of Useless Linguistics, vol 1, pp 10-20.

⇒

```
Nom_auteur : Tanguy
Prénom_auteur : Ludovic
Titre : Overbearingly Supervised Techniques for Very Small Corpora
Journal : Journal of Useless Linguistics
Année : 2015
```

Ancienne méthode (XX^e siècle)

Définir (à la main) des patrons correspondant aux références :

{ $\$$ NOM}, { $\$$ PRENOM} ({ $\$$ ANNEE}). { $\$$ TITRE}. { $\$$ JOURNAL}, vol { $\$$ VOLUME}, pp { $\$$ PAGES}.

Et des sous-patrons pour définir le contenu possible de chaque champ.

Système Paracite : environ 400 patrons de ce type.

Un changement radical en TAL (2)

Nouvelle méthode (XXI^e siècle)

- Étiqueter manuellement quelques centaines de références, champ par champ
- Définir les principaux descripteurs pertinents des segments de texte (ponctuation, majuscules, quelques mots-clés)
- Entraîner un système d'apprentissage automatique (CRF) produisant un modèle probabiliste (boîte noire)
- Le modèle est directement projetable sur de nouvelles données
- Exemple : système *Bilbo* (Kim et al., 2012)

Un changement radical en TAL (3)

Bilan de la transition

- Un temps et un coût de développement réduits
- Des systèmes au moins aussi performants, et souvent plus robustes
- Plus de modèle explicite des données ou de la tâche
- Une décorrélation entre le travail d'expertise sur les données et le développement de l'outil
- Des besoins de post-traitements et/ou d'injection de connaissances pour améliorer les performances

Un phénomène bien installé

- Concerne tous les niveaux du TAL
- Des méthodes éprouvées (CRF, SVM, Entropie Maximale, etc.)
- Des outils facilement disponibles (Weka, RapidMiner, OpenNLP, etc.)
- De nouvelles pratiques de recherche guidées par la tâche (évaluation, comparaison, campagnes...)

Mon point de vue

Raisons de m'y mettre

- Efficacité et coût
- La bonne solution face à la masse et la diversité des données
- Un phénomène incontournable en tant que passeur

Principes

- Me limiter aux méthodes d'apprentissage les plus classiques
- M'en servir comme outil d'exploration des données
- Essayer d'éclairer les boîtes noires et d'observer le comportement face aux données
- En profiter pour valoriser des descripteurs linguistiques

Ma pratique des méthodes par apprentissage

Exploration des données langagières

- Méthodes d'apprentissage symbolique pour identifier les corrélations entre les caractéristiques des données (projets Annodis et Intermede)
- En complément des approches statistiques (bi- et multi-variées)

Classification de documents

- Identification des segments obsolètes : thèse de M. Laignelet
- Attribution d'auteurs : compétitions PAN
- Analyse des rapports d'incidents aéronautiques : travaux avec CFH et thèse de N. Tulechki

Développement d'applications

- Analyse syntaxique probabiliste : *Talismane*, thèse d'A. Urieli
- Typage des requêtes soumises à un moteur de recherche : projet CAAS, thèse de S. Leva

Place des connaissances linguistiques

Un rôle subalterne pour les linguistes

- Dans les scénarios d'utilisation de ces techniques, limité à l'annotation des données d'entraînement et/ou à l'évaluation des résultats
- L'opacité des systèmes est un obstacle aux propositions d'intégration de mécanismes d'analyse basés sur des connaissances des phénomènes

Une tendance à l'appauvrissement des descripteurs

- Le développement de techniques d'apprentissage efficaces et l'augmentation du volume des données a permis l'utilisation de traits pauvres (unités lexicales non traitées, trigrammes de caractères)
- Au détriment des descripteurs plus riches (issus d'une annotation linguistique des données, ou définis à partir d'une observation des données)

Militance pour une meilleure coopération

Des outils d'investigation prêts à l'emploi

- Permettant la prise en compte de nombreux descripteurs hétérogènes
- Plus faciles à l'emploi et à l'interprétation que les techniques statistiques classiques

Des descripteurs plus riches pour les applications

- De nombreuses connaissances produites par la linguistique et ses outils sont aisément intégrables comme descripteurs
- Il faut démontrer leur utilité pour rééquilibrer la collaboration

De premiers résultats encourageants

- Traits linguistiques pertinents pour la prédiction de la difficulté des requêtes en recherche d'information
- Gain significatif apportés par les traits linguistiques pour l'attribution d'auteur

Aller plus loin

D'autres modes d'interaction

- Remettre le linguiste dans la boucle : apprentissage actif, ajout de règles
- Analyser plus finement les données en amont pour sélectionner et définir des descripteurs pertinents
- Proposer de nouveaux descripteurs linguistiques, et en profiter pour étudier leurs interactions

- 1 Vue d'ensemble
- 2 Outillage de l'exploitation des corpus
- 3 Utilisation du Web comme corpus
- 4 Visualisation et interprétation des données linguistiques
- 5 Apprentissage artificiel et articulation avec la linguistique
- 6 Discussion sur le TAL en linguistique**
 - Bienfaits et méfaits
 - Enseignement du TAL en sciences du langage

Un outillage utile

Pour gérer la masse

Rassembler des données, y accéder, les annoter, les fouiller

Pour formaliser

Rechercher un phénomène, ou le faire annoter automatiquement est un des moyens d'arriver à une description fiable (mais partielle) de celui-ci

Pour expérimenter

Vers une linguistique expérimentale : opérationnalisation rapide d'hypothèses à grande échelle

Pour collaborer

Les outils forment un bon lieu de discussion (visualisation, analyses)
Nouvelles applications du TAL ⇒ nouvelles données langagières ⇒ nouvelles questions linguistiques

Des outils coupants

Des comportements à risque

- Fascination de l'outil
- Éloignement des données
- Le marteau et les clous : ne plus voir dans les données que ce qu'une méthode automatisée permet d'atteindre

Déformation du matériau par l'outil

- À vouloir gérer la masse on perd parfois les individus intéressants
- Multiplication des descripteurs : (trop) facile de jouer avec les méthodes de fouille sans (trop) réfléchir

Conséquences sur l'enseignement : nécessité d'armer les futurs spécialistes de TAL

Pour accéder à l'autonomie

Maîtriser les différentes étapes, de l'acquisition à l'analyse des données

Pour affronter les techniques complexes

Trouver sa place dans le TAL moderne

Pour pouvoir s'adapter

Des savoirs fondamentaux et pas d'enfermement dans des outils spécifiques

Compétences techniques dans une formation de TAL (1)

Manipulation du matériau langagier	
L3	<i>Prévenir, diagnostiquer et résoudre les problèmes de codage des fichiers de textes (codage des caractères et des fins de ligne)</i>
L3	<i>Identifier et savoir convertir les différents formats de documents pour au minimum en extraire le contenu textuel (et si possible, quelques méta-données et/ou informations structurelles)</i>
Utilisation de l'outillage	
L3	<i>Appliquer des annotateurs génériques (étiqueteurs, parseurs) en tenant compte du paramétrage de ceux-ci et des spécificités formelles des textes cibles</i>
L3	<i>Convertir les sorties d'un analyseur aux besoins d'un outil qui les exploite (concordancier, outil d'interrogation, extracteur spécialisé, tableur ou outil d'analyse statistique, base de données, etc.)</i>
M1	<i>Projeter des ressources sur un corpus (lexiques, marqueurs, patrons morphosyntaxiques, etc.)</i>
M1	<i>Assembler des unités de traitement en une chaîne complète, cumuler les informations notamment en utilisant un format XML</i>

Compétences techniques dans une formation de TAL (2)

Exploitation	
L3	<i>Effectuer des recherches</i> et des extractions simples dans des fichiers de texte nu (segmenter, extraire des séquences de formes, calculer les fréquences)
L3	<i>Exploiter des annotations</i> automatiques pour rechercher et extraire des unités plus complexes
M1	<i>Évaluer un traitement</i> , en utilisant des mesures adaptées à la tâche, et en se basant sur un étalon ou sur des compétences de linguiste
M1	<i>Diagnostiquer</i> des fonctionnements erronés d'un système, et identifier un contournement ou une solution
Conception	
M1	<i>Déployer une méthode</i> de traitement raisonnablement complexe à partir d'une description du type de celles que l'on trouve présentées dans un article scientifique
M2	<i>Proposer une méthode</i> de traitement adaptée à un besoin nouveau, en se basant sur des solutions connues, mais nécessitant une articulation spécifique
M2	<i>Savoir dialoguer avec un informaticien</i> développeur, dans un travail d'équipe, en participant aux prises de décision nécessaires au développement d'un prototype

Des vertus de l'apprentissage technique

La programmation

- Le prix de l'autonomie
- Un mode de réflexion face à de nouvelles questions
- Une école de la rigueur
- Mieux vaut une trousse qu'un outil supposé universel

La pédagogie par projet

- Savoir s'impliquer dans une équipe
- Avoir la satisfaction d'une réalisation aboutie et en retirer de la confiance en ses capacités

Merci à mes 66 co-auteurs !

Par ordre alphabétique

Afantenos, Armstrong, Asher, Aussenac-Gilles, Benamara, Bouillon, Bourigault, Bras, Calderone, Caudy, Chrisment, Condamines, Dal, Dkaki, Enjalbert, Fabre, Ferrari, François, Génolini, Hathout, Hermann, Ho-Dac, Josselin-Leray, Kanellos, Kiss, Kompaore, Laignelet, Lalleman, Lang, Le Draoulec, Le Dû, Legras, Lehmann, Leva, Lewin, Lignon, Marchand, Mathet, Membrado, Milward, Montermini, Mothe, Muller, Namer, Narquet, Péry-Woodley, Pimm, Plénat, Poulain, Prévot, Raynal, Rebeyrolle, Sajous, Schieber, Schmouchkovitch, Séguéla, Serna, Thlivitis, Tulechki, Urieli, Vergely, Vergez-Couret, Vieu, Walker, Widlöcher, Zaldivar-Carillo