

FOCUS 3

Discours : ressources et systèmes par apprentissage

Jean-Yves Antoine (LIFAT), Agata Savary (LIFAT), Denis Maurel (LIFAT), Anne-Lyse Minard (LLL), Emmanuel Schang (LLL), Lotfi Abouda (LLL), Flora Badin (LLL), Guillaume Cleuziou (LIFO), Sylvie Billot (LIFO), Gaëtan Caillaut (LIFO), Anais Lefeuvre- Halftermeyer (LIFO), Loïc Grobol (LATTICE)

Plan

Phénomènes

Données

Méthodes

Outils

Phénomènes

- Coréférences
- Temporalité

j'ai quitté Orléans en 1964, j'y suis revenue en 1968.

Phénomènes

- Coréférences
- Temporalité

j'ai quitté Orléans en 1964, j'y suis revenue en 1968.



Phénomènes

- **Coréférences**
- Temporalité



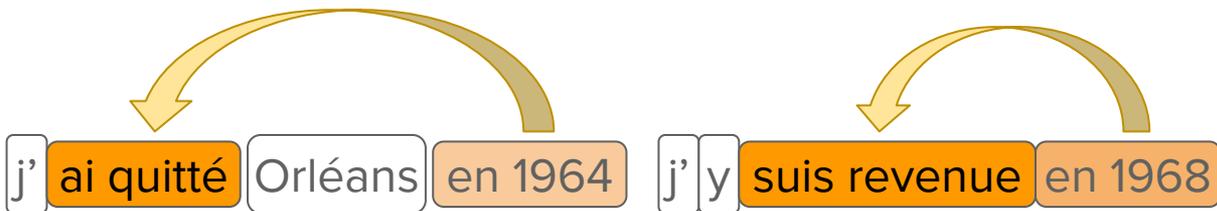
Phénomènes

- Coréférences
- Temporalité

j' ai quitté Orléans en 1964 j' y suis revenue en 1968

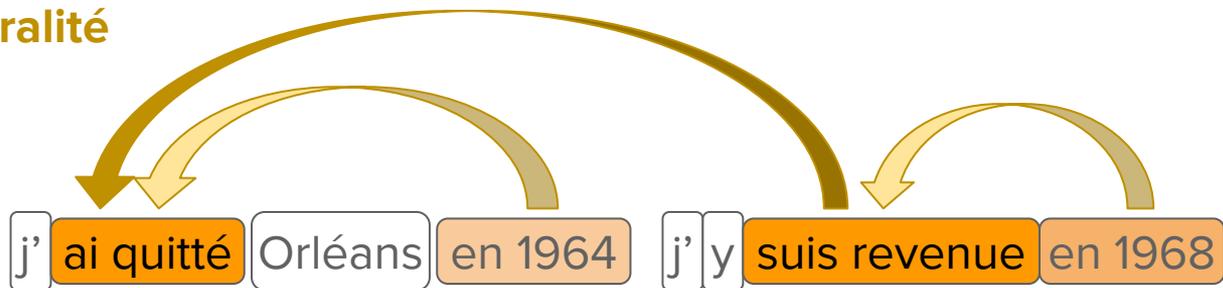
Phénomènes

- Coréférences
- Temporalité



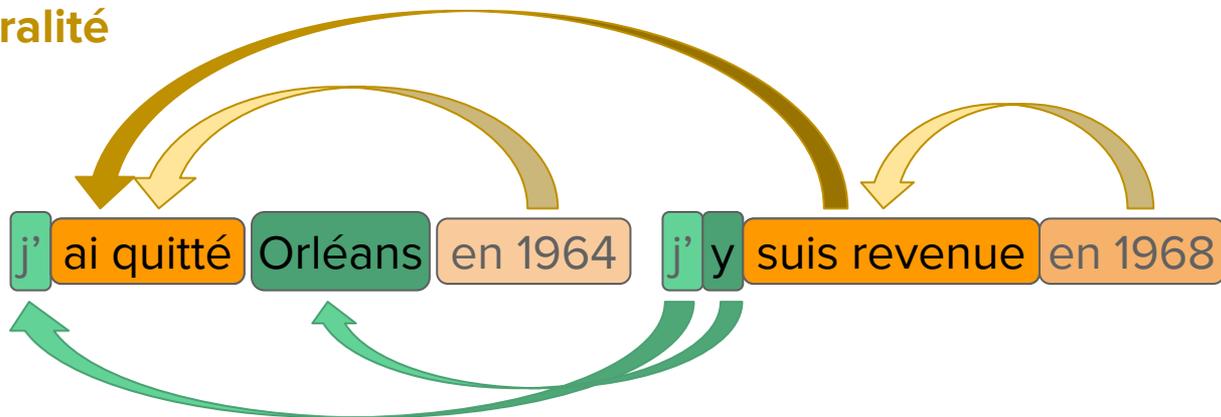
Phénomènes

- Coréférences
- **Temporalité**



Phénomènes

- Coréférences
- Temporalité



Données

- Corpus ANCOR

oral transcrit, spontané, avec interactivité variable

⇒ plus grand corpus de la coréférence à l'oral, plus gros corpus français annoté en coréférence = 488 000 mots (jusqu'à DEMOCRAT, imminent)

- Corpus ODIL

en cours de réalisation, autre couche du corpus ANCOR, annoté en syntaxe automatiquement, temporalité (semi ?) manuellement

Données

- Corpus ANCOR

Annotation manuelle de traits fins :

- genre,
- nombre,
- type d'EN,
- définitude,
- généricité,
- etc.

- Corpus ODIL

Syntaxe automatique

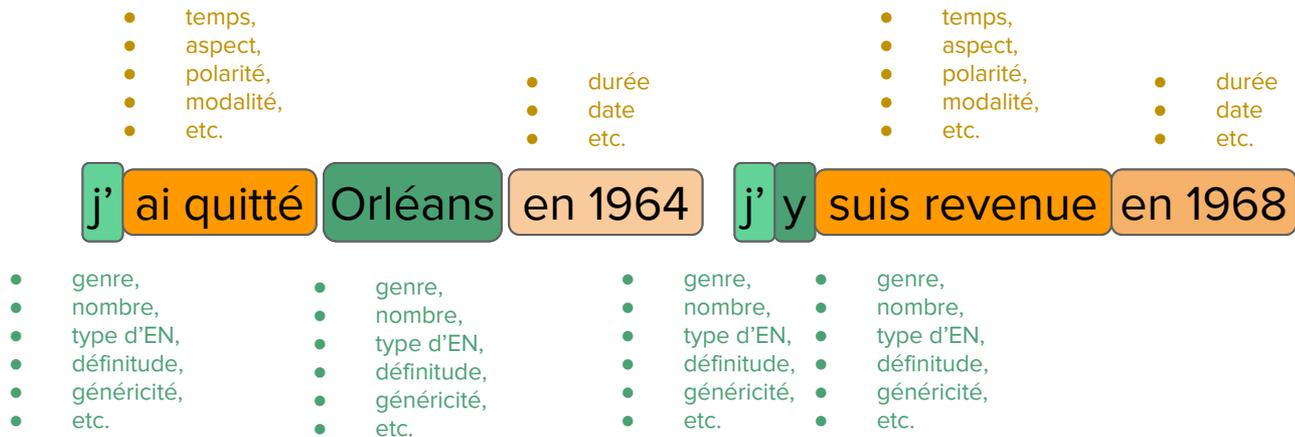
Correction manuelle et annotation manuelle de traits fins :

- temps,
- aspect,
- polarité,
- modalité,
- etc.

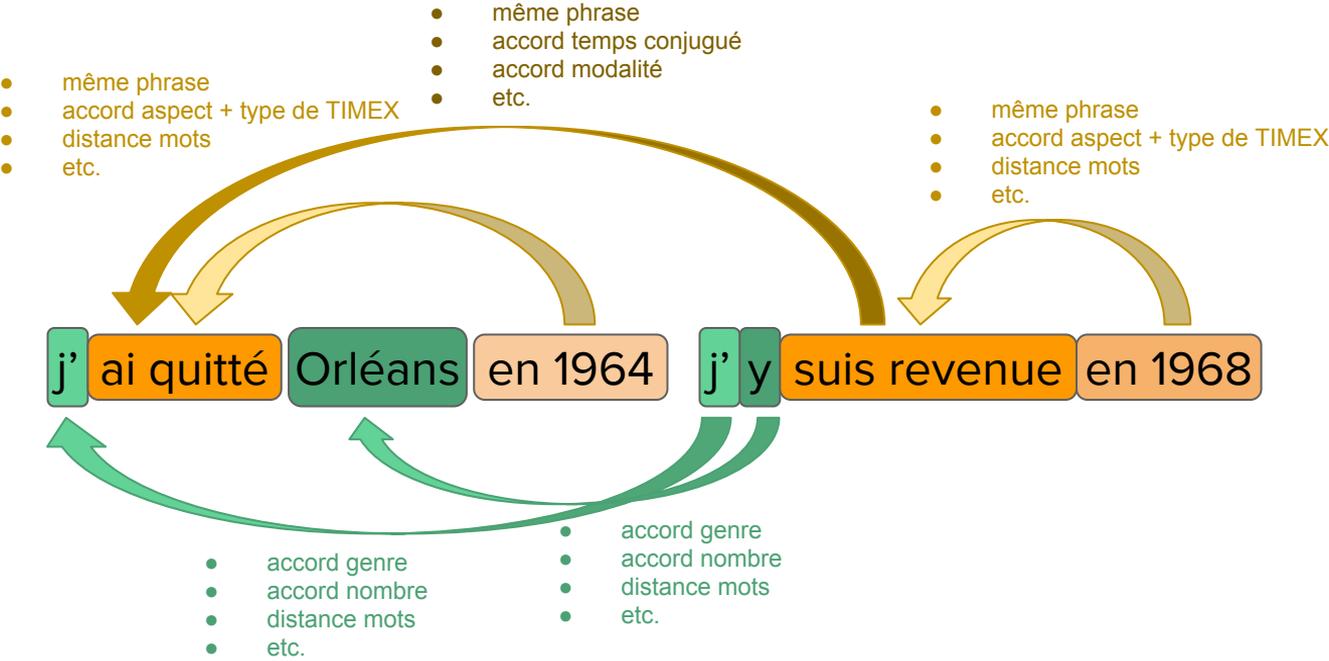
Données

j' ai quitté Orléans en 1964 j' y suis revenue en 1968

Données



Données



Méthodes

- Travail sur la norme ISO-TIMEML
- Production de corpus annotés

- Apprentissage supervisé pour la création de chaînes référentielles :
Naive Bayes, Decision Tree, SVM

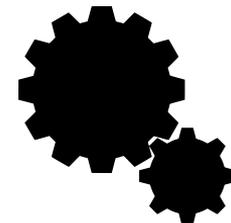
- Apprentissage d'espace prétopologique structurant pour extraire un graphe temporel issu d'un discours (Gaëtan Caillaut LIFO)

Méthodes

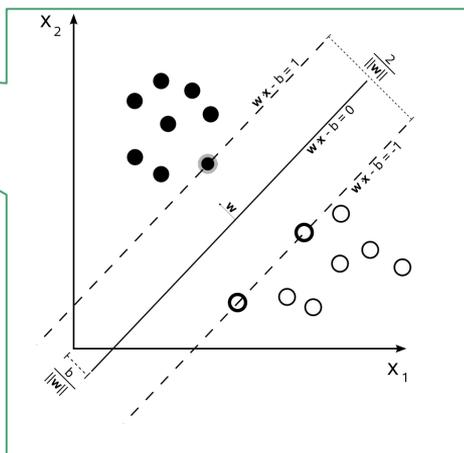


```
==== Classifier model (full training set) ====
348 pruned tree
-----
com_rate <= 0.25
| distance_minion <= 35
| | id_gendre = YES
| | | id_an = YES
| | | | id_mest = YES
| | | | | distance_char <= 1: NOT_COREF (441.0/3.0)
| | | | | distance_char > 1
| | | | | distance_turn <= 81: COREF (109.0/9.0)
| | | | | distance_turn > 81: NOT_COREF (54.0/5.0)
| | | | | id_mest = NO
| | | | | distance_char <= 784
| | | | | distance_char <= 8
| | | | | distance_word <= -7
| | | | | | id_nombre = YES
| | | | | | | m2_type = N; NOT_COREF (6.0/1.0)
| | | | | | | m2_type = PR
| | | | | | | distance_char <= -187: NOT_COREF (9.0/3.0)
| | | | | | | distance_char > -187: COREF (81.0/8.0)
| | | | | | | m2_type = NULL: COREF (0.0)
| | | | | | | id_nombre = NO: NOT_COREF (14.0/1.8)
| | | | | | | id_nombre = UNK: COREF (0.0)
| | | | | | | distance_char > 7
| | | | | | | id_spk = YES
| | | | | | | distance_char <= 1
| | | | | | | | m1_new = YES: NOT_COREF (41.0)
| | | | | | | | m1_new = NO
| | | | | | | | m2_type = N: NOT_COREF (3.0)
| | | | | | | | m2_type = PR: COREF (2.0)
| | | | | | | | m2_type = NULL: NOT_COREF (0.0)
| | | | | | | | m1_new = UNK: NOT_COREF (0.0)
| | | | | | | | m1_new = NULL: NOT_COREF (0.0)
| | | | | | | | distance_char > 21: COREF (15.0)
| | | | | | | | id_spk = NO: NOT_COREF (341.0/2.0)
| | | | | | | | id_spk = NA: NOT_COREF (0.0)
| | | | | | | distance_word > 0
| | | | | | | distance_word <= 26
| | | | | | | | id_nombre = YES: COREF (8144.0/187.0)
```

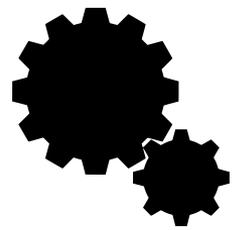
modèle



- Apprentissage supervisé pour la création de chaînes référentielles :
Naive Bayes, **Decision Tree**, **SVM**



modèle



Méthodes

- Apprentissage supervisé pour la **création de chaînes référentielles** :
Naive Bayes, Decision Tree, SVM

J' ... une enseignante ... il ... j'

Méthodes

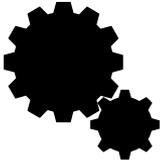
- Apprentissage supervisé pour la **création de chaînes référentielles** :
Naive Bayes, Decision Tree, SVM



Méthodes

- Apprentissage supervisé pour la **création de chaînes référentielles** :
Naive Bayes, Decision Tree, SVM

- genre : U
- nombre : SG
- EN : PERS
- etc.



- genre : F
- nombre : SG
- EN : PERS
- etc.

J'

une enseignante

...

...

il

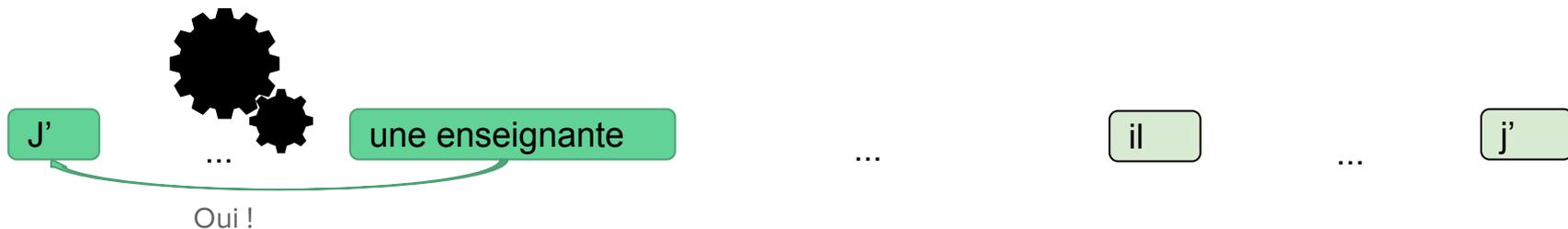
...

j'

- accord genre : no
- accord nombre : yes
- accord EN : yes
- distance = 6 mots
- etc.

Méthodes

- Apprentissage supervisé pour la **création de chaînes référentielles** :
Naive Bayes, Decision Tree, SVM



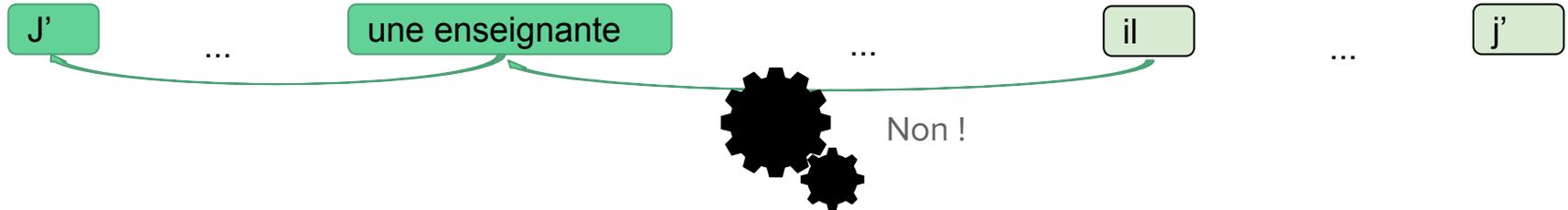
Méthodes

- Apprentissage supervisé pour la **création de chaînes référentielles** :
Naive Bayes, Decision Tree, SVM



Méthodes

- Apprentissage supervisé pour la **création de chaînes référentielles** :
Naive Bayes, Decision Tree, SVM



Méthodes

- Apprentissage supervisé pour la **création de chaînes référentielles** :
Naive Bayes, Decision Tree, SVM

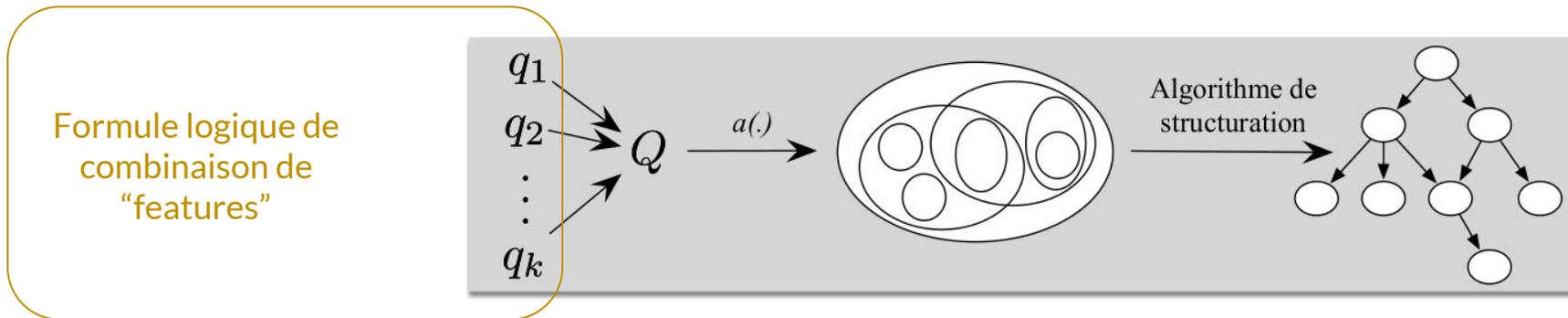


Méthodes

- Apprentissage supervisé pour la **création de chaînes référentielles** :
Naive Bayes, Decision Tree, SVM



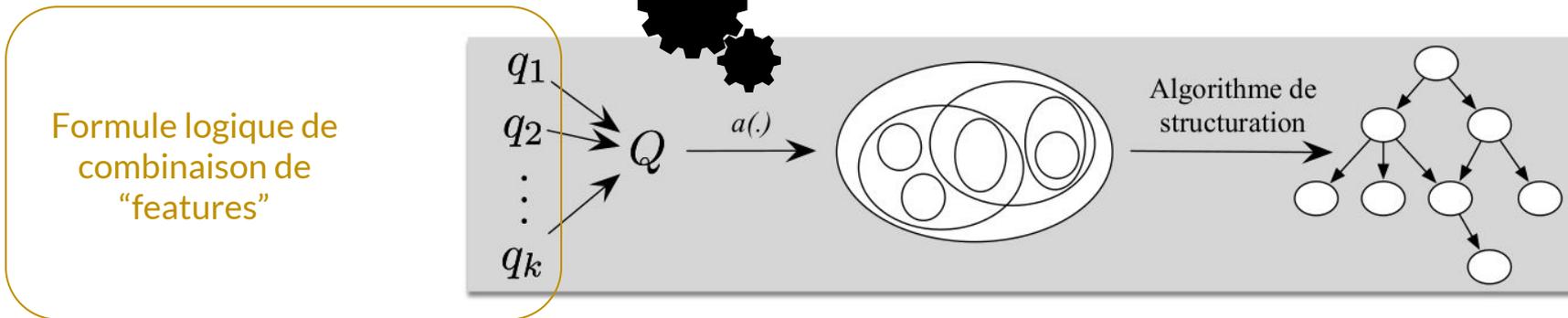
Méthodes



- Apprentissage d'espace prétopologique structurant pour extraire un graphe temporel issu d'un discours (Gaëtan Caillaut LIFO)

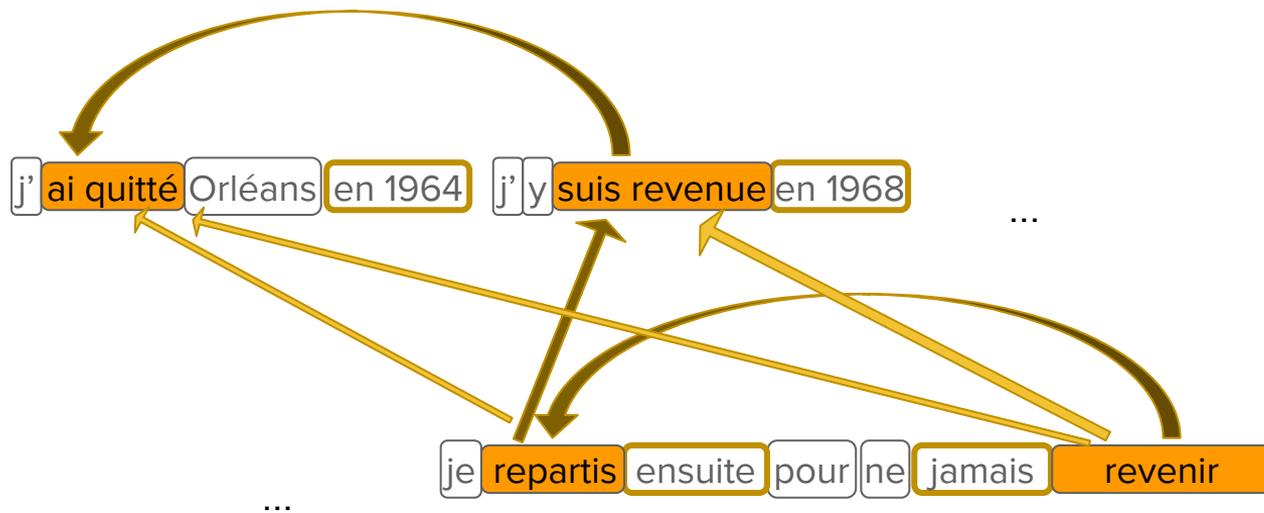
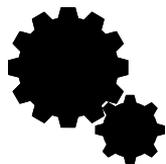
Méthodes

Fonction d'adhérence

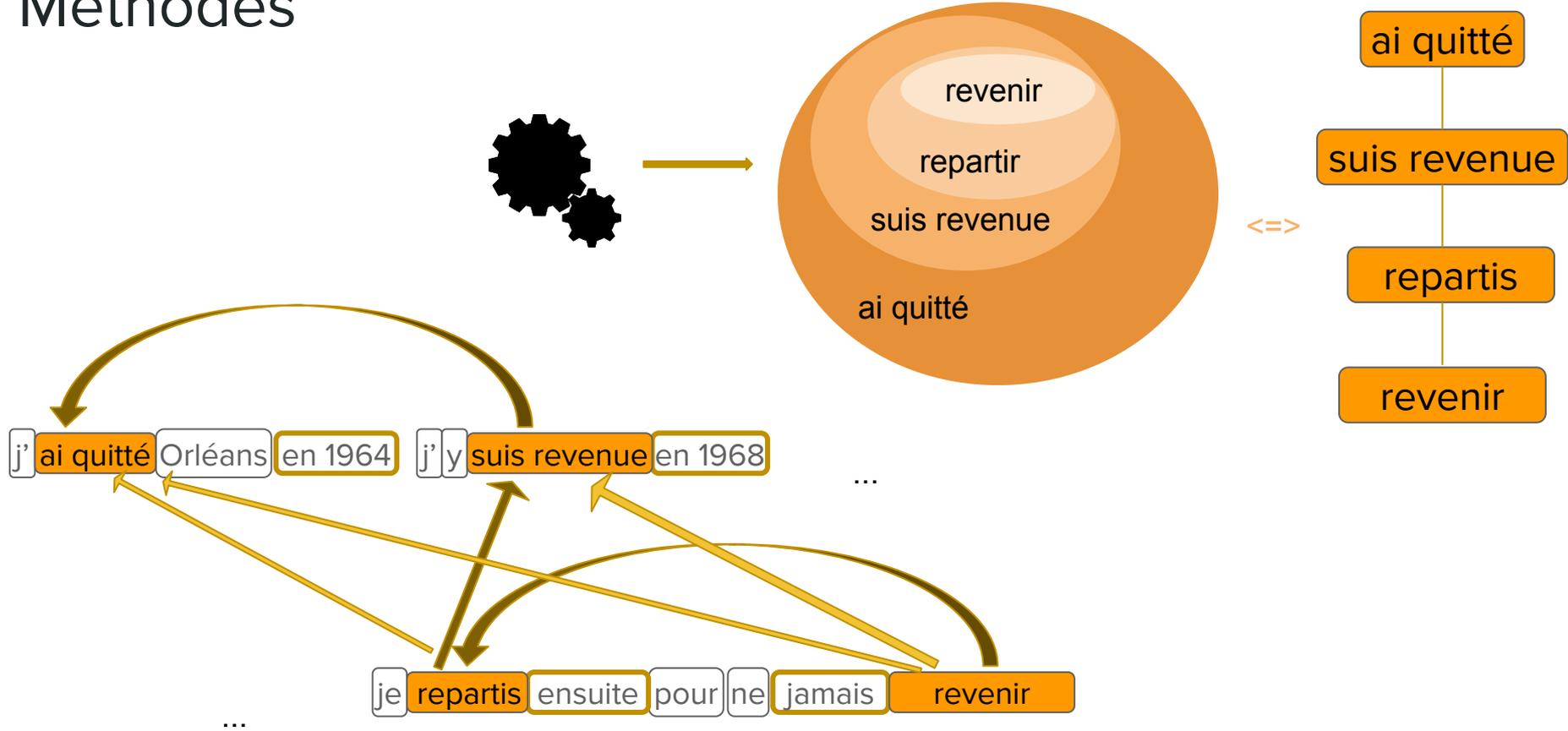


- Apprentissage d'espace prétopologique structurant pour extraire un graphe temporel issu d'un discours (Gaëtan Caillaut LIFO)

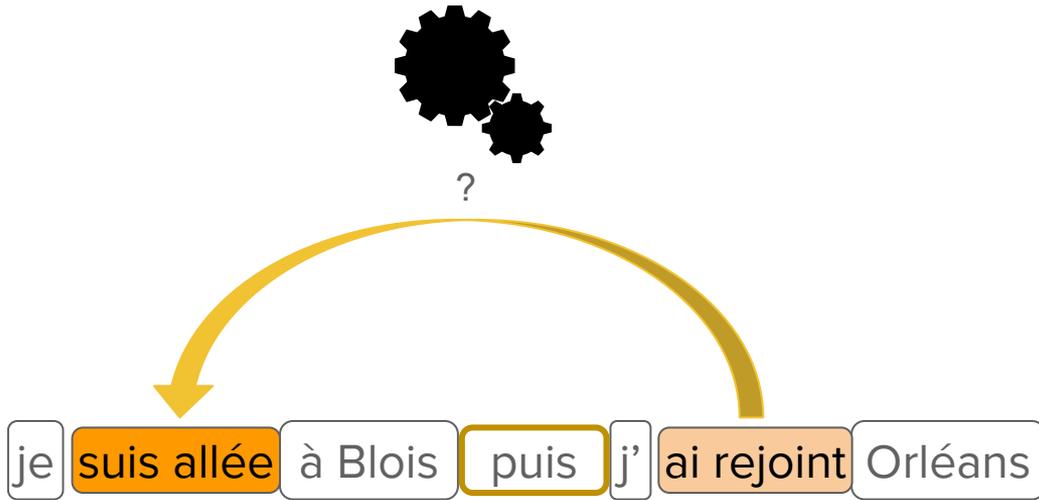
Méthodes



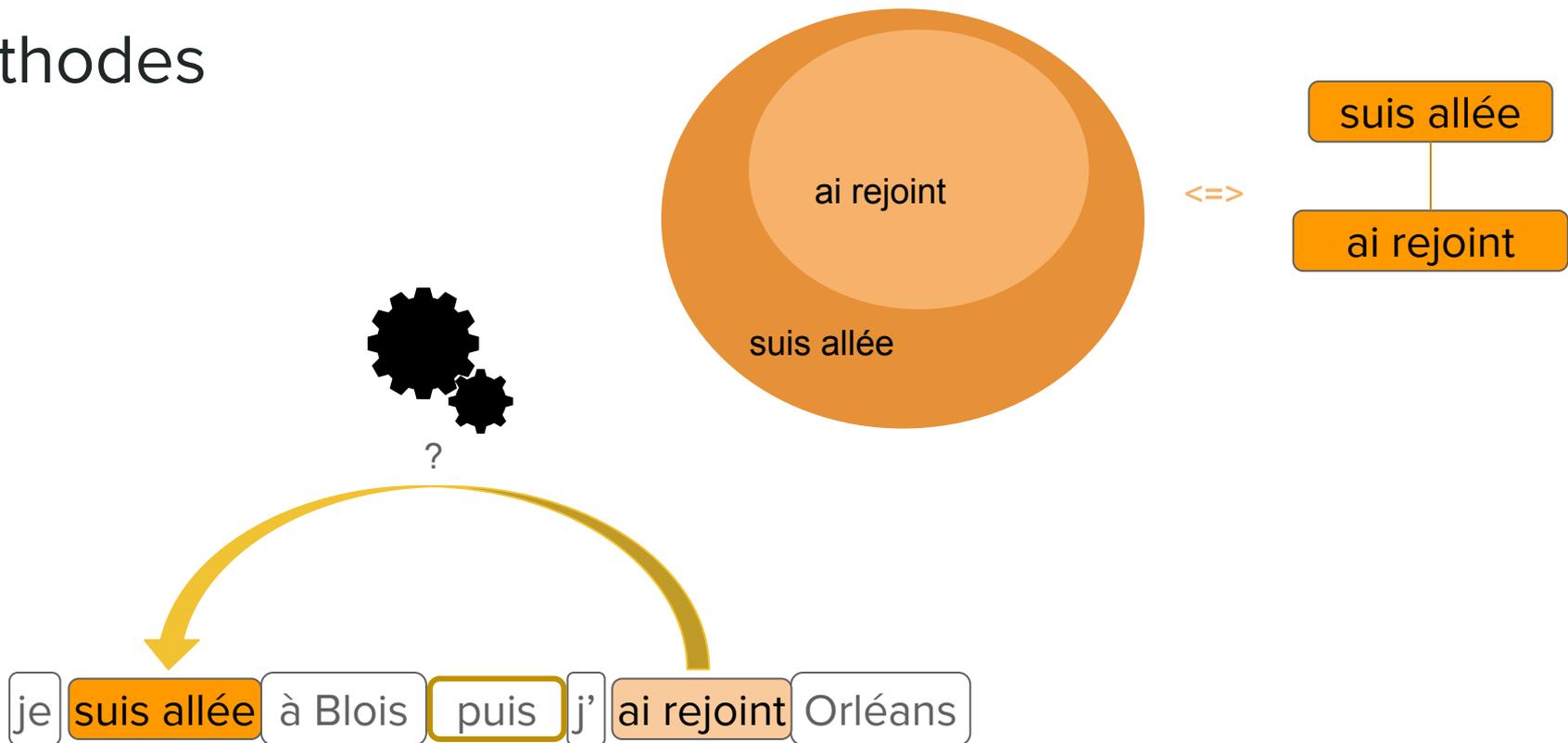
Méthodes



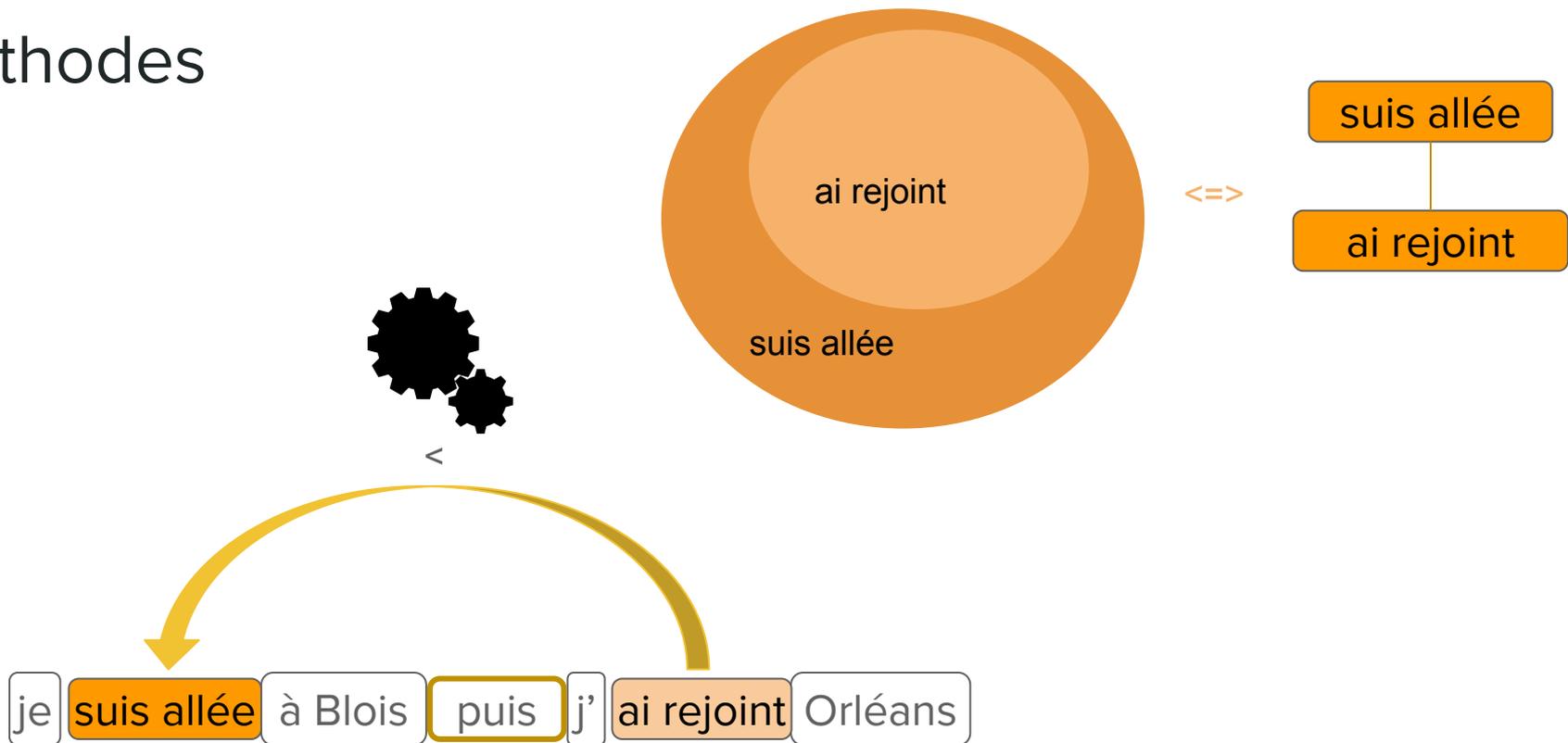
Méthodes



Méthodes



Méthodes



Outils

- **RÉSOLUTION** : CROC (LATTICE)... et ses petits ! (collaboration Loïc Grobol LATTICE)
 - rendre la chaîne end-to-end
 - intégration de la détection de mention

- **RÉSOLUTION** : PrétopoTal
 - adaptation de la création d'ontologies à la création de graphes temporels de discours

- **ANNOTATION** : CONTEMPLATA
 - refonte annotation

Autres collaborations

RAVIOLI

Métriques

PREDICT4ALL

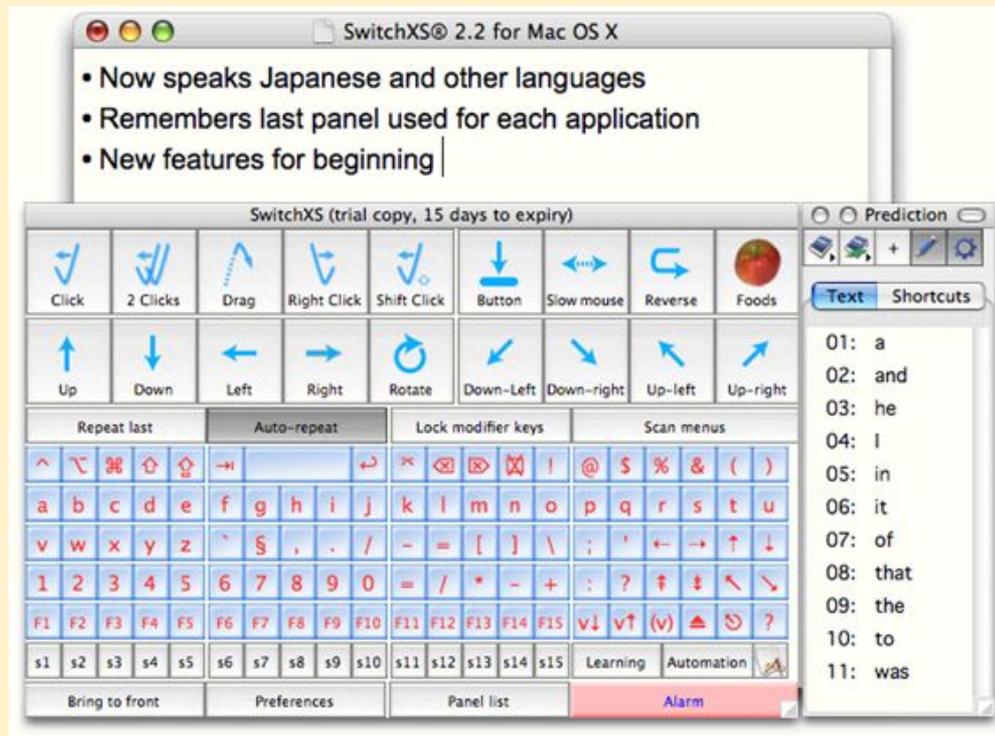
PREDICT4ALL (aide au handicap)

Phénomènes et Données

Méthodes et Outils

Phénomènes et données

- Saisie de texte pour personnes souffrant de handicaps moteurs sévères (locked-in syndrom, paralysie cérébrale, maladie de Charcot, tétraplégie...)
- **Clavier virtuel**
- **Enjeu TAL** : prédiction de texte pour saisie plus rapide



Méthodes

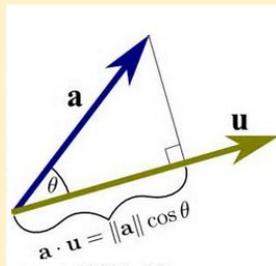
- Prédiction de texte combinant un modèle de langage statistique (syntaxe) et un modèle d'adaptation sémantique (thématique du discours)

$$P'(w_i) = \frac{P_b(w_i)^{\lambda_1} \cdot P_s(w_i)^{(1-\lambda_1)}}{\sum_{j=1}^n P_b(w_j)^{\lambda_1} \cdot P_s(w_j)^{(1-\lambda_1)}}$$

$$P(w_i) = \lambda_1 \times P_{\text{Base}}(w_i | w_{i-1} w_{i-2} w_{i-3}) + (1-\lambda_1) \times P_{\text{DUM}}(w_i | w_{i-1} w_{i-2})$$

N-gram général +
N-gram utilisateur

$$P_{\text{LSA}}(w_i | h) = \frac{(\cos(\vec{w}_i, \vec{h}) - \cos_{\min}(\vec{h}))^\gamma}{\sum_k (\cos(\vec{w}_k, \vec{h}) - \cos_{\min}(\vec{h}))^\gamma}$$



Analyse sémantique latente
(équivalence *word embeddings*)

=> espace vectoriel sémantique

Outils

- Système Sibylle vK : diffusion libre (*open source* en 2020)

grâce à Sibylle, je peux communiquer avec mon père et ma

mots	pictos	l	j	t	d	p	a	<	mère					
e	c	s	m	r	i	f	u	q	filie					
n	à	o	é	v	b	g	h	y	famille					
ê	k	w	z	ù	x	ç	_	è	grand-mère					
									soeur					
,	.	-	?	!	.	@	;		vie					
HP	abc	pictos	clavier	<	<<				tante					
									abc				!	123

- **Predict4All** adaptation troubles dyslexiques