

# CHAPITRE V

## Traitement automatique du langage naturel

Natural Language Processing



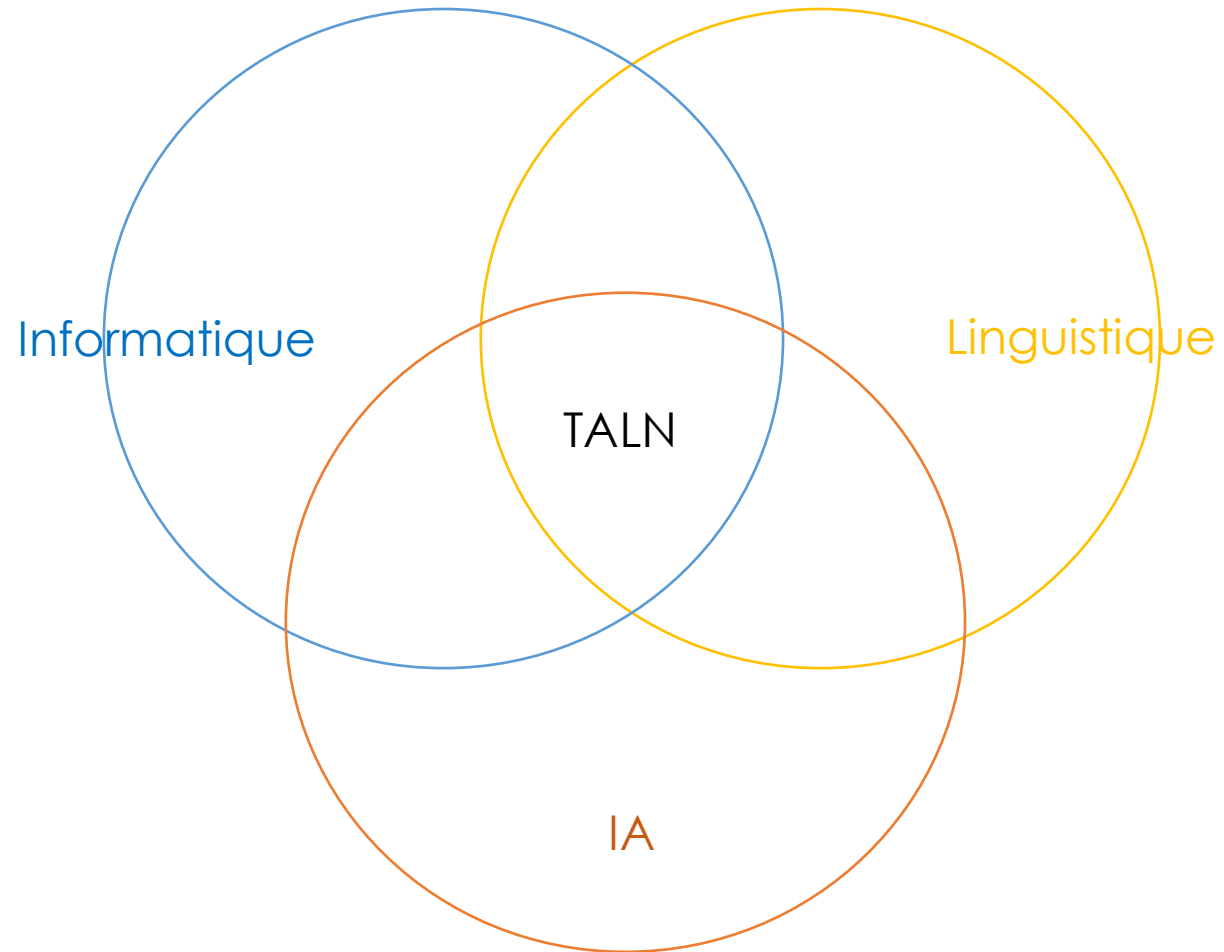
# Histoire

- L'histoire du TALN commence dans les années 50s.
- Les premiers travaux concernaient la traduction automatique qui en fut l'une des premières applications informatiques
- Traitement automatique des conversations
- Tester l'intelligence des machines (test de Turing)
- 1954, traduction automatique des phrases russes (politique) vers l'anglais
- SHRDLU (1960)
- ELIZA (1964-1966)

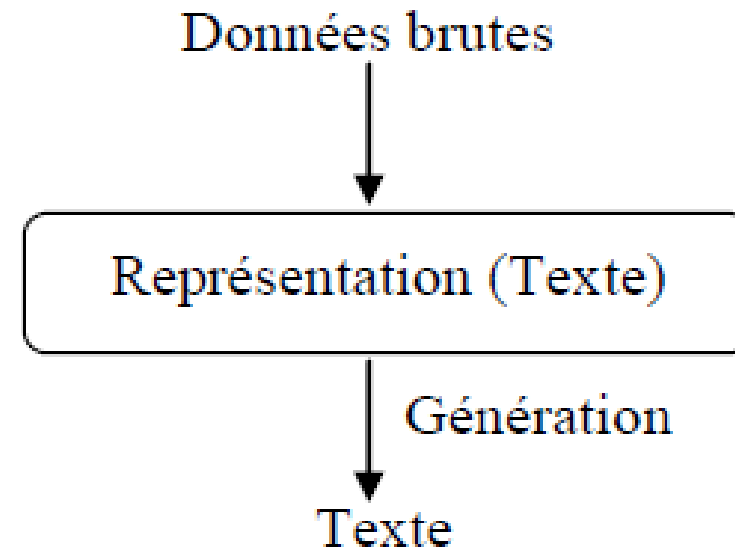
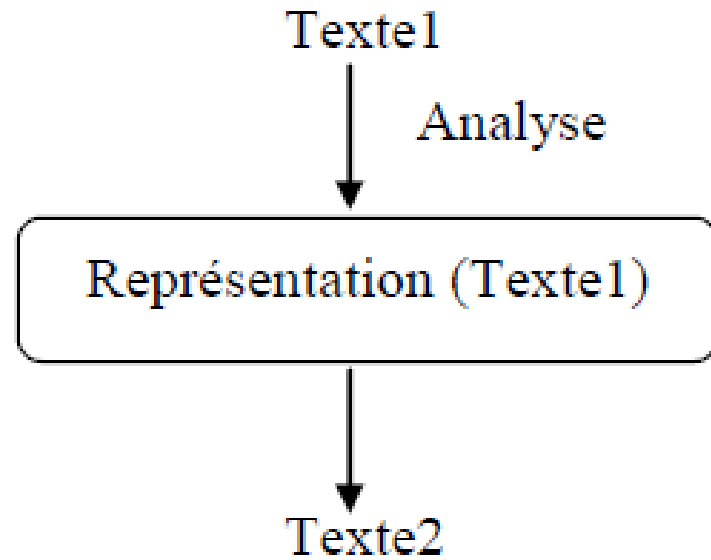
## Définition

- Le TALN à pour objet la création de programmes informatiques capables de traiter automatiquement les langues naturelles
- La langue naturelle désigne la langue parlée ou écrites par les êtres humains
- Traitement des données linguistiques :
  - Textes écrits (paragrapes, phrases, mots,..)
  - Données orales (phonèmes, prosodies,..)

# TALN : Domaine multidisciplinaire



# Analyse vs Génération



# Domaines d'application

	<b>Texte</b>	<b>Parole</b>
<b>Analyse</b>	<ul style="list-style-type: none"><li>▪ Correction orthographique</li><li>▪ Aide à la reformulation</li><li>▪ Recherche d'information fouille textuelle</li><li>▪ Reconnaissance d'entités nommées</li><li>▪ Classification et catégorisation de documents</li><li>▪ Reconnaissance de l'écriture manuscrite</li><li>▪ Annotation morpho- syntaxique / sémantique</li></ul>	<ul style="list-style-type: none"><li>▪ Reconnaissance vocale</li><li>▪ Identification du locuteur</li></ul>
<b>Génération</b>	<ul style="list-style-type: none"><li>▪ Génération automatique de textes</li><li>▪ Résumé automatique</li></ul>	<ul style="list-style-type: none"><li>▪ Synthèse de la parole</li></ul>

# Outils TALN

- **Linguistiques** : théories linguistiques qui décrivent les différentes connaissances relatives à la langue
- **Formels** : expriment les connaissances linguistiques dans un formalisme qui convient à un traitement automatique
- **Informatiques** : utilisent la description formelle des connaissances dans une application informatique concrète

# Niveaux de traitement



Traitement phonétique

Prétraitement



Traitement morphologique

Extraction des caractéristiques

Traitement syntaxique

Traitement sémantique

Traitement pragmatique





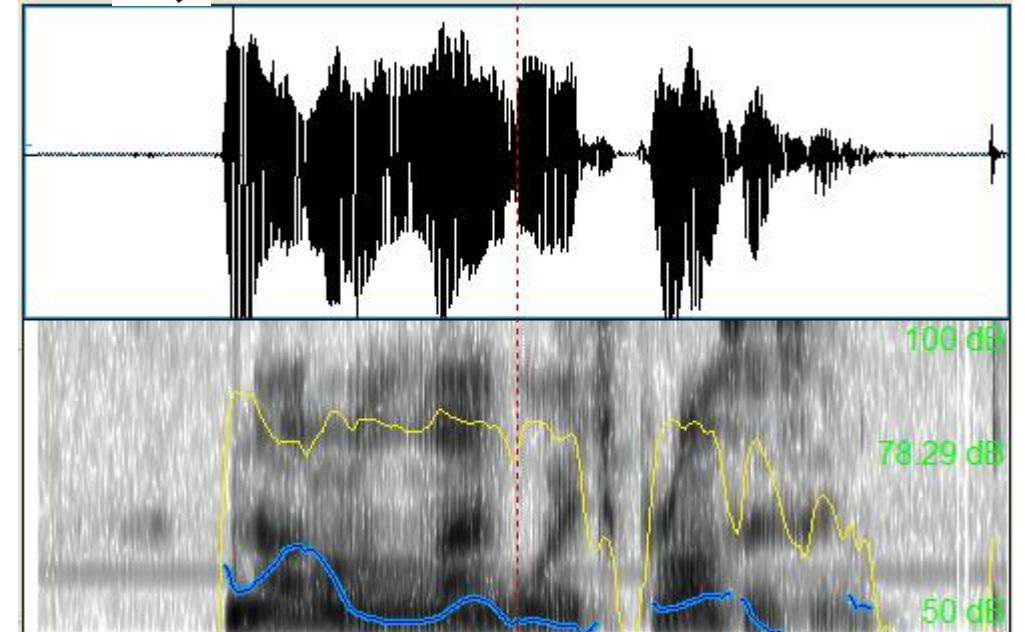
# Niveaux de traitement

## ■ Niveau phonétique (phonologique)

- Traitement de la langue orale
- La machine doit reconnaître des signaux acoustiques et les identifier en temps que mots via une interface vocale
- Identifier les prosodies



Natural | Language | Processing



'nætʃrəl 'læŋgwɪdʒ 'prəʊsesɪŋ

# Niveaux de traitement

## ■ Prétraitement

- Segmentation
- Tokenisation
- Suppression de mots vides
- Normalisation
- Encodage des caractères

يذهب محمد إلى المسجد كل يوم

يذهب، محمد، إلى، المسجد، كل، يوم

يذهب، محمد، المسجد، يوم

# Niveaux de traitement

## ■ Niveau morphologique (morpholexical)

- Etudier la formation des mots et leur variation de formes.
- Flexion vs Dérivation
  - Flexion : modifications que subit un mot dans sa terminaison (ou son début)
    - Ajout d'un affixe qui ne crée pas un nouveau lexème
    - Exemple : cheval / chev**aux**      œuf / œuf**s**
  - Dérivation : formation de mots nouveaux par addition, suppression ou remplacement d'un élément grammatical d'un mot simple
    - création d'un nouveau lexème par l'ajout d'un affixe
    - Exemple : **dé**lavé, **in**tolérable, **ill**isible

# Niveaux de traitement

## ■ Niveau morphologique (morpholexical)

- Racinisation (stemming)
  - petit, petite, petits, petites ⇒ petit
- Lemmatisation
  - Pêche, pêcher, pêcheur ⇒ pêche
  - eu, avions ⇒ avoir

محمد، المسجد

Stemming

حمد، سجد

يسترجعون

Lemmatisation

استرجع

# Niveaux de traitement

## ■ Extraction des caractéristiques

- Text vectorization
  - Bag of words
  - N-grams
- TF (term frequency)
- TF-IDF (term frequency-inverse document frequency)

Variantes de TF

Schéma de pondération	formule du TF
binaire	0, 1
fréquence brute	$f_{t,d}$
normalisation logarithmique	$1 + \log(f_{t,d})$
normalisation « 0.5 » par le max	$0.5 + 0.5 \cdot \frac{f_{t,d}}{\max_{\{t' \in d\}} f_{t',d}}$
normalisation par le max	$K + (1 - K) \frac{f_{t,d}}{\max_{\{t' \in d\}} f_{t',d}}$

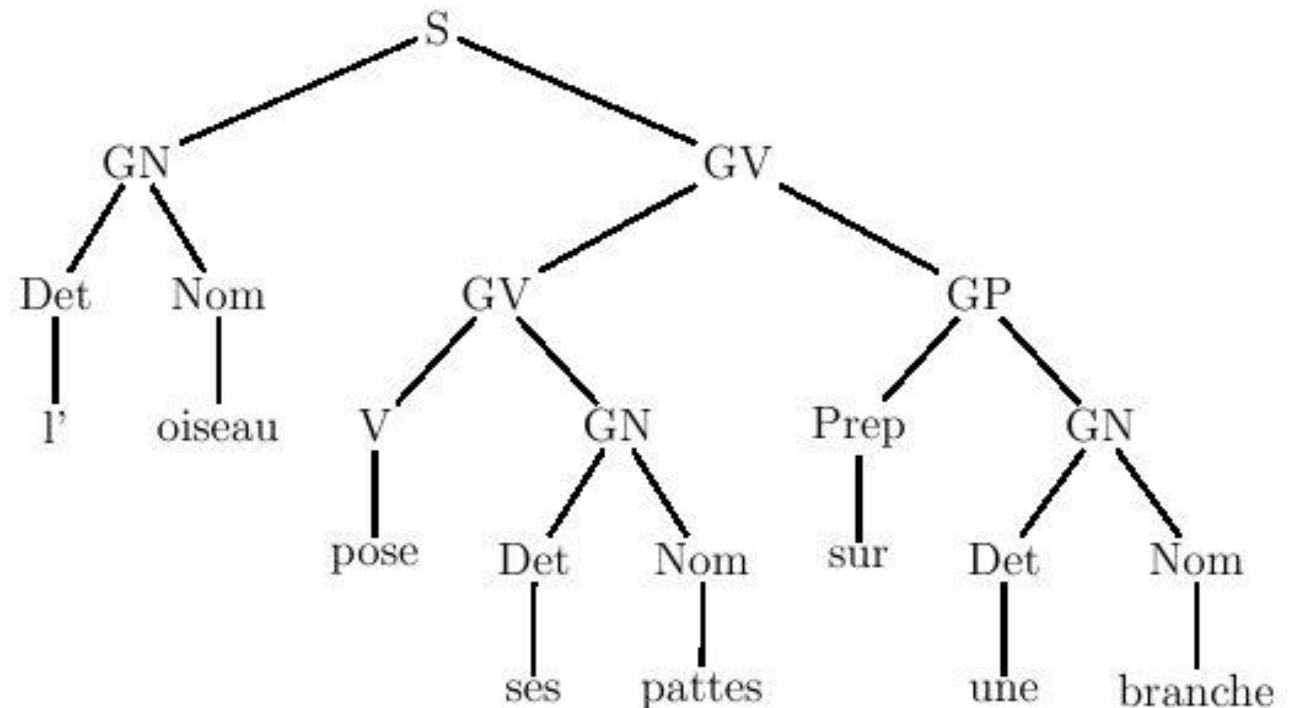
$$\text{idf}_i = \log \frac{|D|}{|\{d_j : t_i \in d_j\}|}$$

$$\text{tfidf}_{i,j} = \text{tf}_{i,j} \cdot \text{idf}_i$$

# Niveaux de traitement

## ■ Niveau syntaxique

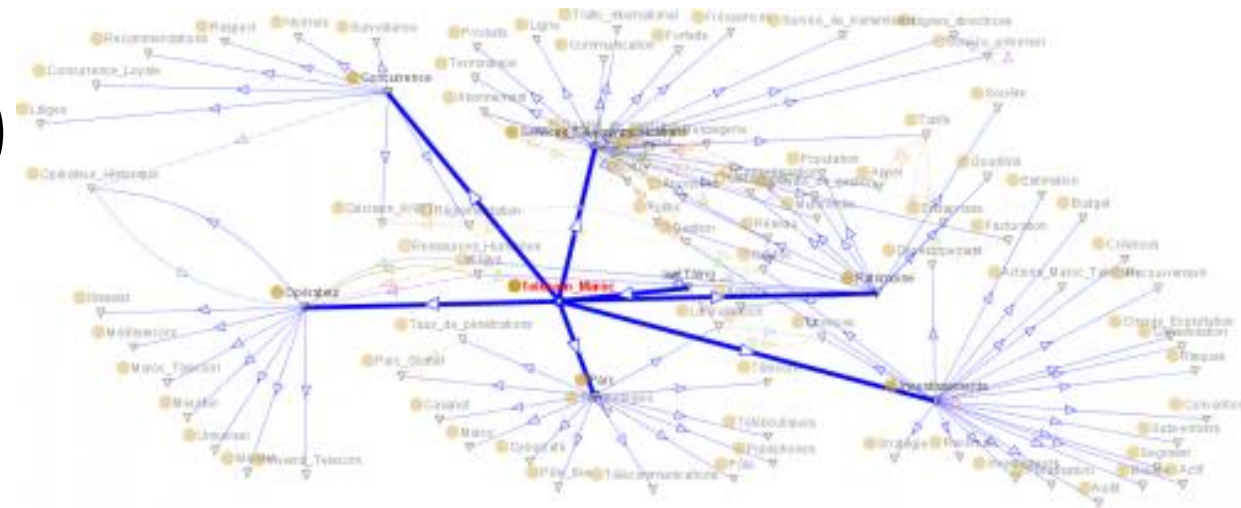
- Catégorisation des mots
- Etudier l'agencement des mots et leurs relations structurelles
- Détection et formation des phrases correctes (en utilisant des grammaires par exemple)
- Lever les ambiguïtés grammaticales



# Niveaux de traitement

## ■ Niveau sémantique

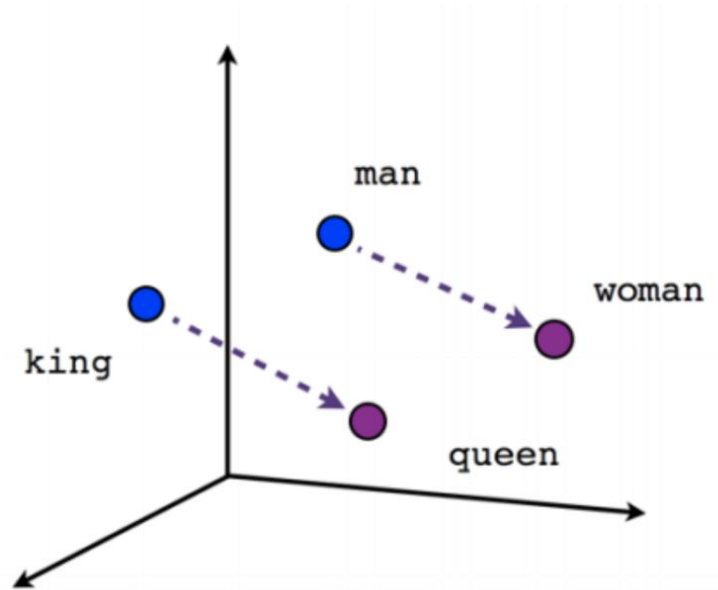
- Associer un sens aux mots, phrases, énoncés
- Détecter les structures n'ayant pas de sens
- Utilisation des ressources lexicales et des ontologies pour :
  - Chercher des synonymes (WordNet)
  - Identifier les entités nommées
  - Voyéllation automatique (Arabe)



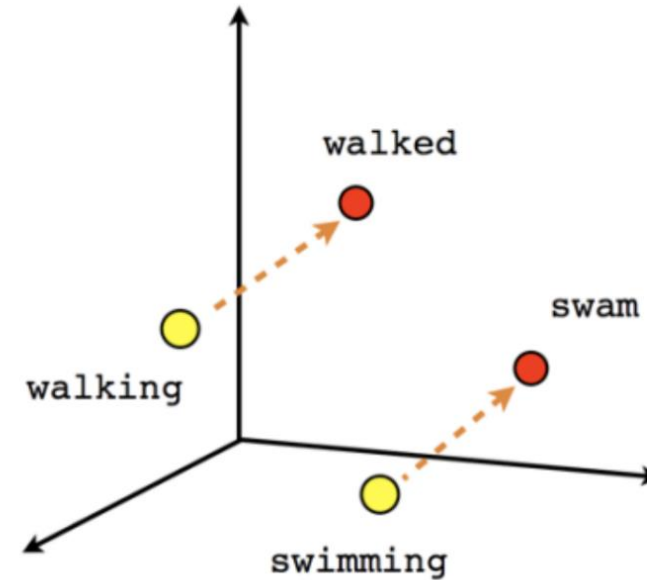
# Niveaux de traitement

- **Niveau pragmatique**

- Situer les mot et les énoncés dans les bons contextes



Male-Female



Verb tense



# Evaluer un système TALN

- Connaitre la sortie souhaitée ?
  - Recherche d'information : connaitre les documents pertinents
  - Résumé automatique : connaitre les phrases correctes
  - Traduction automatique : connaitre les traductions correctes
  - Classification de documents : connaitre les documents bien classés

# Mesures d'évaluation d'un système TALN

- Rappel vs Précision
  - Classification de documents : connaitre les documents bien classés

$$\text{Rappel}_i = \frac{\text{Documents correctement attribués à la classe}_i}{\text{Nombre de documents attribués à la classe}_i}$$

$$\text{Précision}_i = \frac{\text{Documents correctement attribués à la classe}_i}{\text{Nombre de documents appartenant à la classe}_i}$$

# Mesures d'évaluation d'un système TALN

- Rappel vs Précision
  - Classification de documents : connaître les documents bien classés

Sujet réel	Sport	Sport	Autre	Autre	Autre	Sport	Autre	Sport	Autre	Sport
Document	1	2	3	4	5	6	7	8	9	10
Sortie système	Sport	Autre	Sport	Sport	Autre	Sport	Autre	Sport	Sport	Sport

Réponses correctes

- Rappel et Précision pour la classe sport ?

$$\text{Rappel}_i = \frac{\text{Documents correctement attribués à la classe}_i}{\text{Nombre de documents attribués à la classe}_i}$$

$$\text{Précision}_i = \frac{\text{Documents correctement attribués à la classe}_i}{\text{Nombre de documents appartenant à la classe}_i}$$