

# Artex est un Autre Résuméur TEXTuel

**Juan-Manuel Torres**  
**2024**

Université d'Avignon

# Objectifs

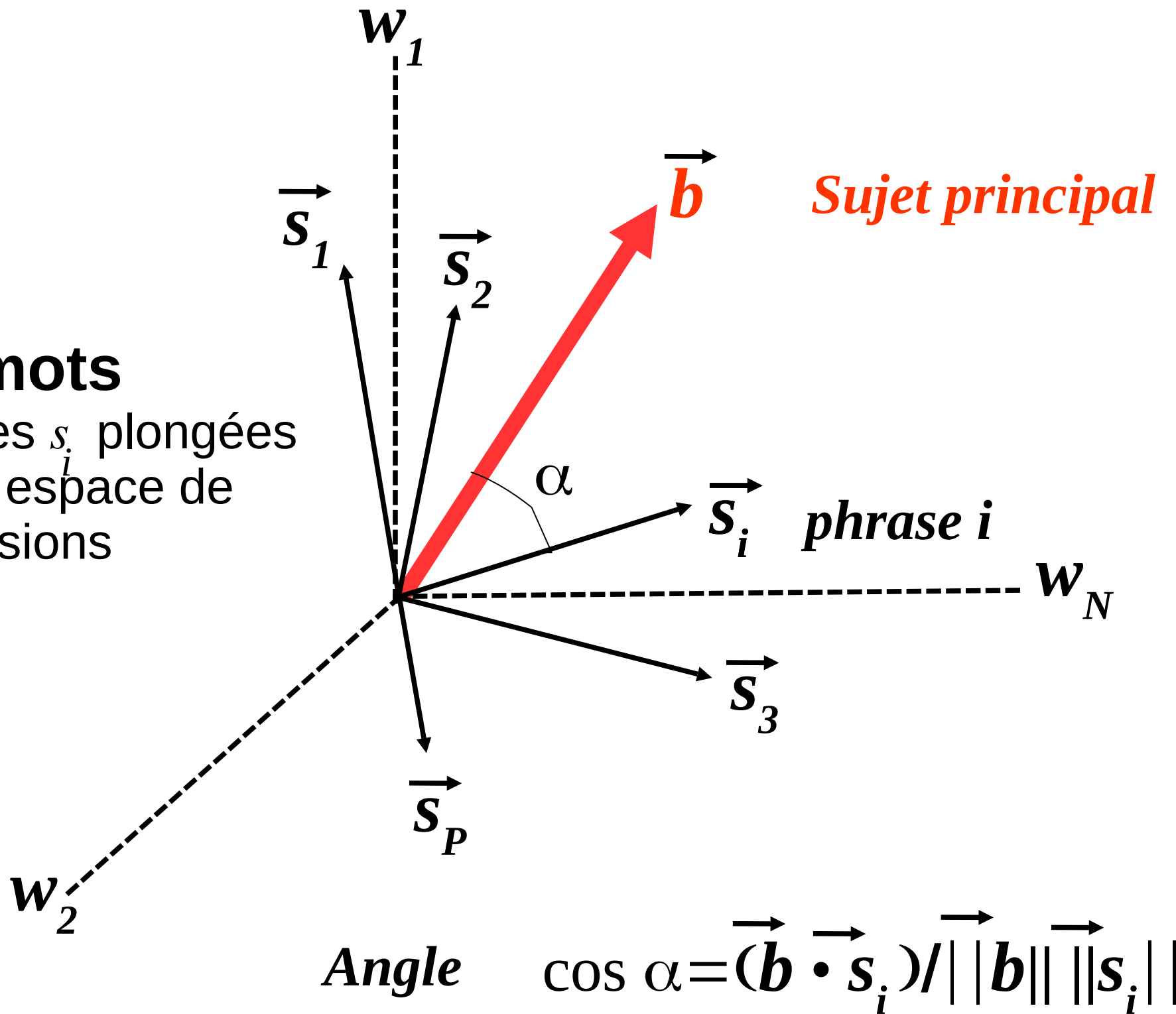
- Poser la tâche de résumé comme un problème géométrique
- Indépendance de la langue
- Indépendance du domaine

# ARTEX

- **Pondération et tri** de phrases
- **Extraction** des phrases importantes
  - **Pre-traitement** classique
    - Segmentation
    - Filtrage
    - Lemmatisation/Stemming/Ultra-stemming
  - Calcul du **vecteur sujet principal a**
  - Calcul du **vecteur poids lexical b**
  - **Produit scalaire a.b**
  - **Post-traitement** de surface

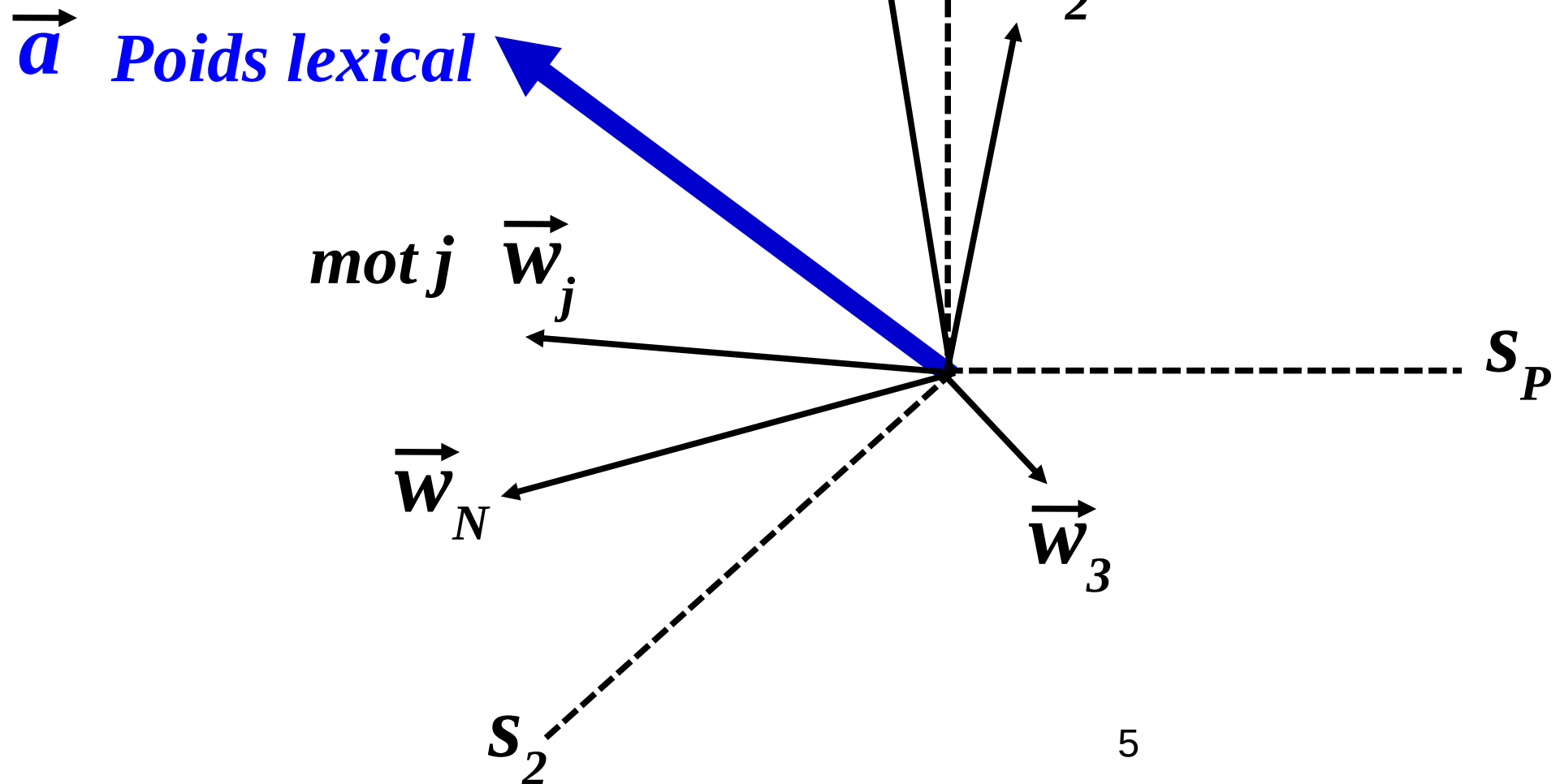
## VSM mots

$P$  phrases  $s_i$  plongées  
dans un espace de  
 $N$  dimensions



# VSM phrases

$N$  mots  $w_j$  plongés  
dans un espace de  
 $P$  dimensions



# Algorithme

$$(1) \quad a_i = \frac{1}{N} \sum_j s_{i,j} \quad \text{Poids lexicale}$$

$$(2) \quad b_j = \frac{1}{P} \sum_i s_{i,j} \quad \text{Sujet principal}$$

$$(3) \quad \text{score}(s_i) = \left( \vec{s} \times \vec{b} \right) \times \vec{a} = \frac{1}{NP} \left( \sum_j s_{i,j} \times b_j \right) \times a_i ;$$

# Artex

P = 3 phrases

N = 5 mots

	avignon	pont	rhone	palais	rempart
Phrase 1	0	1	2	0	1
Phrase 2	1	0	0	1	1
Phrase 3	0	0	0	0	1

$$a = 1/5 \times \begin{pmatrix} 4 \\ 3 \\ 1 \end{pmatrix} = \begin{pmatrix} 4/5 \\ 3/5 \\ 1/5 \end{pmatrix}$$

$$b = 1/3 \times \begin{pmatrix} 1 & 1 & 2 & 1 & 3 \end{pmatrix} = \begin{pmatrix} 1/3 & 1/3 & 2/3 & 1/3 & 1 \end{pmatrix}$$

$$\begin{aligned} \text{Score}(1) &= 1/(N \times P) \begin{pmatrix} 0 & 1 & 2 & 0 & 1 \end{pmatrix} \times \begin{pmatrix} 1/3 & 1/3 & 2/3 & 1/3 & 1 \end{pmatrix} \times 4/5 \\ &= 1/15 \times \begin{pmatrix} 0 + 1/3 + 4/3 + 0 + 1 \end{pmatrix} \times 4/5 \end{aligned}$$

$$\text{Score}(1) = 0,06 \times 2,67 \times 0,8 = 0,143$$

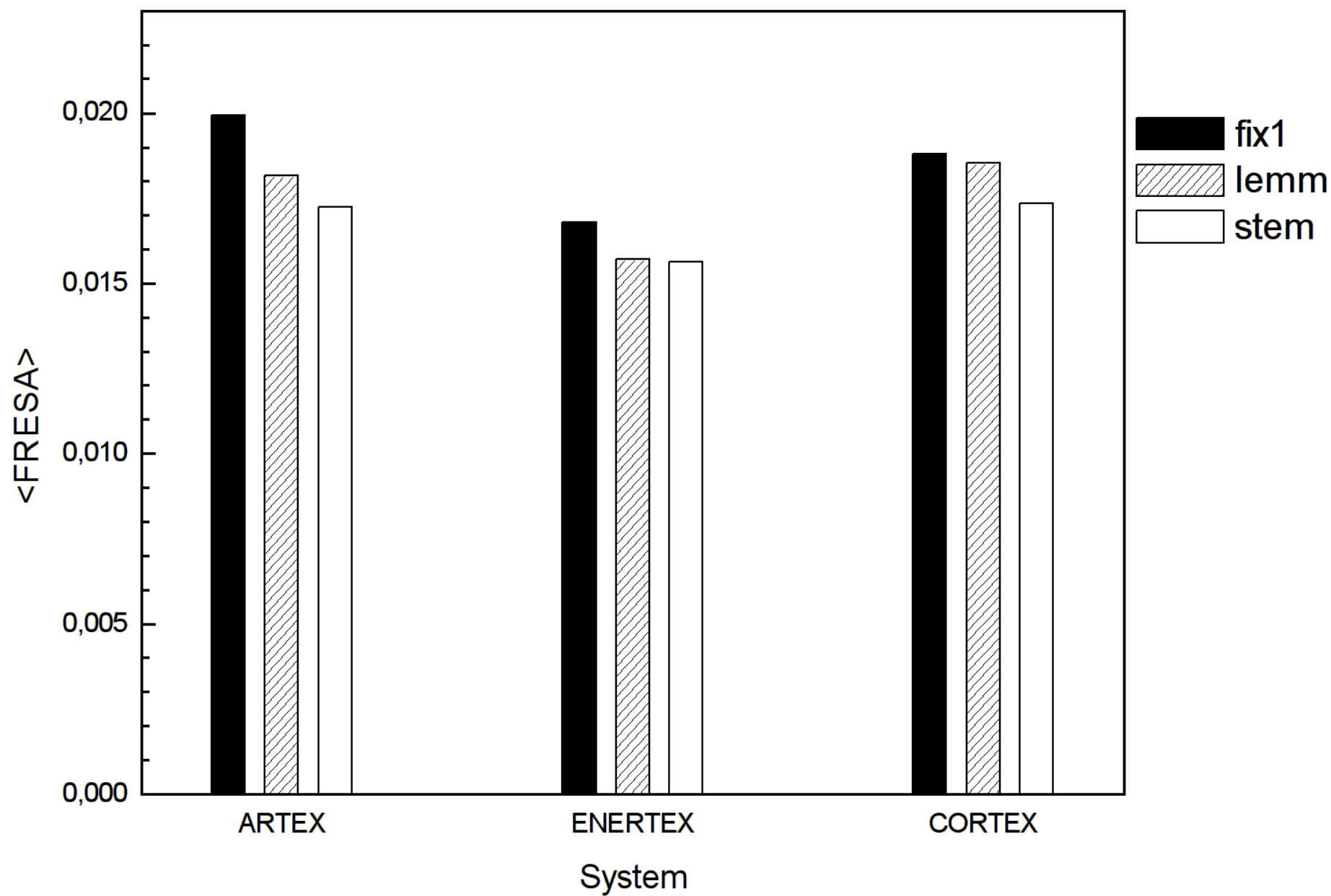
$$\begin{aligned} \text{Score}(2) &= 1/15 \times \begin{pmatrix} 1 & 0 & 0 & 1 & 1 \end{pmatrix} \times \begin{pmatrix} 1/3 & 1/3 & 2/3 & 1/3 & 1 \end{pmatrix} \times 3/5 \\ &= 1/15 \times \begin{pmatrix} 1/3 + 0 + 0 + 2/3 + 1 \end{pmatrix} \times 3/5 \end{aligned}$$

$$\text{Score}(2) = 0,06 \times 2 \times 0,6 = 0,08$$

$$\text{Score}(3) = 1/15 \times \begin{pmatrix} 0 & 0 & 0 & 0 & 1 \end{pmatrix} \times \begin{pmatrix} 1/3 & 1/3 & 2/3 & 1/3 & 1 \end{pmatrix} \times 1/5$$

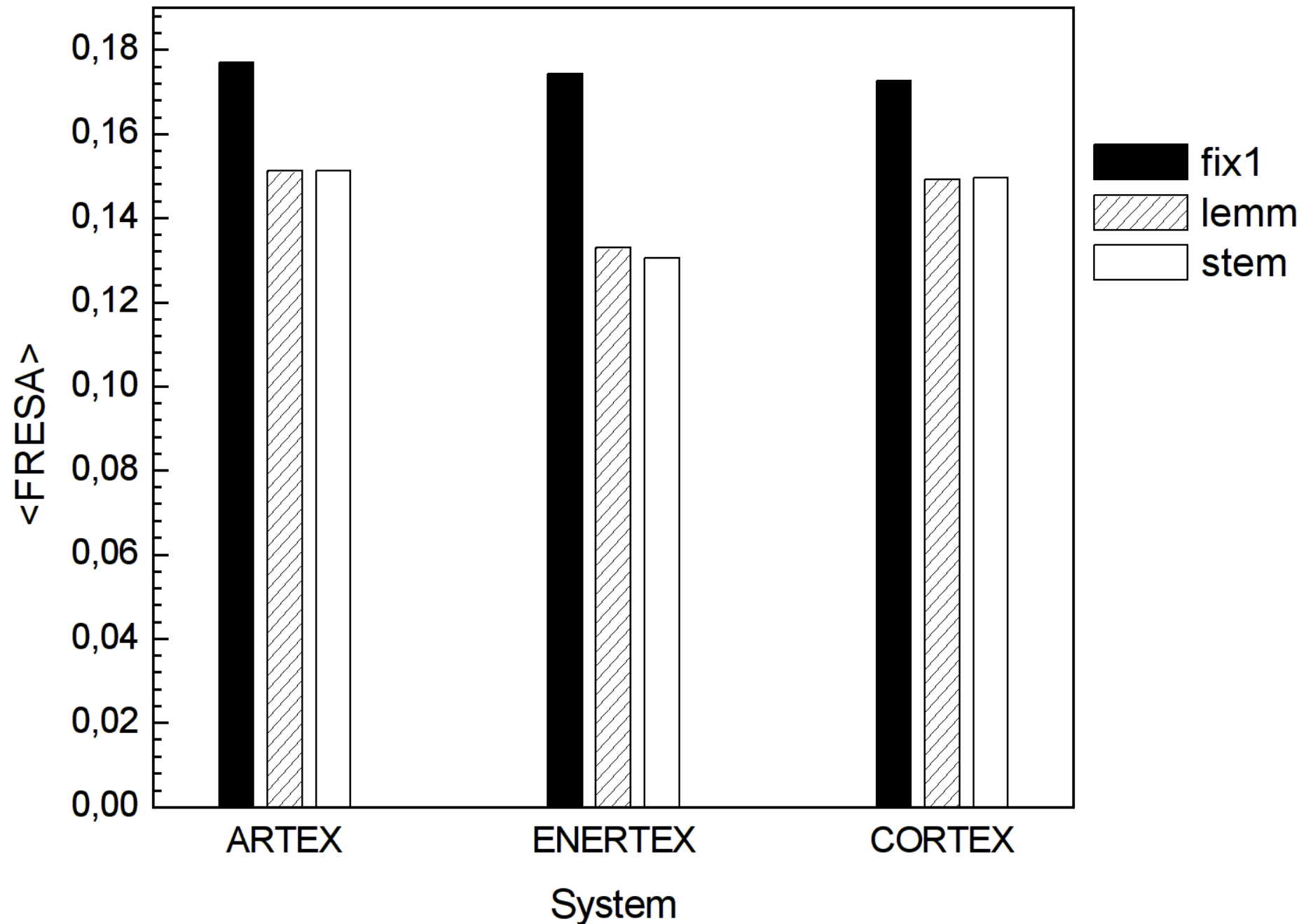
$$\text{Score}(3) = 0,06 \times 1 \times 0,2 = 0,01$$

**CORPUS: DUC'04**  
**(English - 50 Clusters Task 2, Generic summarization)**

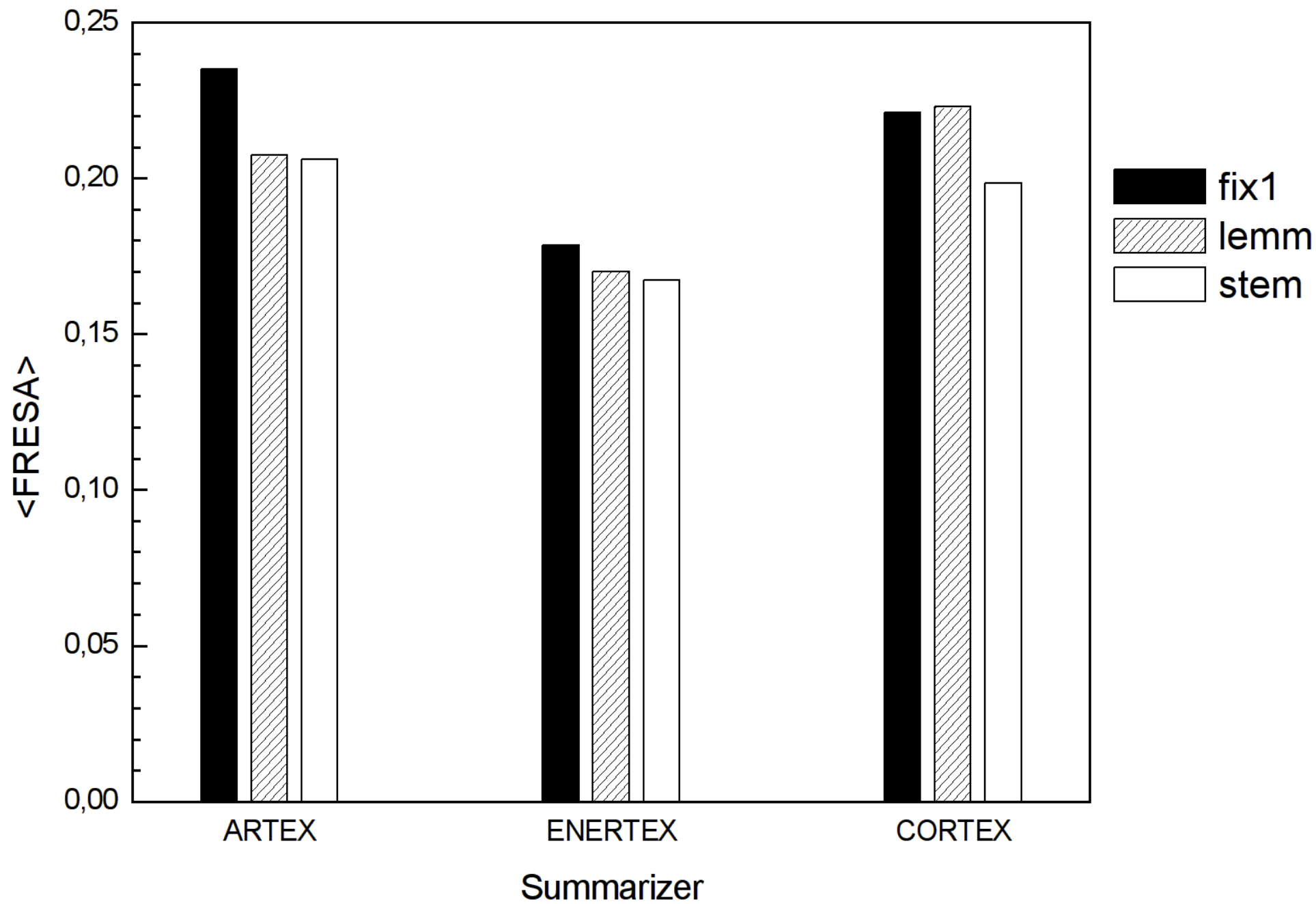




**Corpus: MEDICINA CLINICA**  
**(Spanish - 50 docs - Guided summarization)**



Corpus PISTES  
(French - 50 docs - Generic summarization)



# Artex avec plogements de mots?

Des mots aux embeddings des mots...

# Des mots et des textes...

## Corpus

- *Ensemble de documents (textuels) ayant des caractéristiques intéressantes (taille, domaine, format,...)*
- *Corpus vastes → Apprentissage automatique (modèles)*
- *Corpus spécialisées: évaluation, adaptation des modèles*

## Corpus brut

- Variations de la source: PDF (horrible!), texte (encodage), etc
- Disponibilité
- Propriété

En IA il faut pré-traiter les documents → Corpus pré-traité

# Pre-traitement de textes

## Corpus

- *Avignon, la belle capitale du Vaucluse, jouit d'un climat vraiment bon*
- *De ce fait, les courageux étudiants de l'Université, sont ainsi récompensés dans leurs énormes efforts.*
- *Espérons alors qu'ils réussiront ce cours très intéressant !*

## Filtrage/Lemmatisation (mots porteurs de sens)

- avignon capitale vaucluse jouir climat
- courage étudiant université récompense effort
- réussir cours

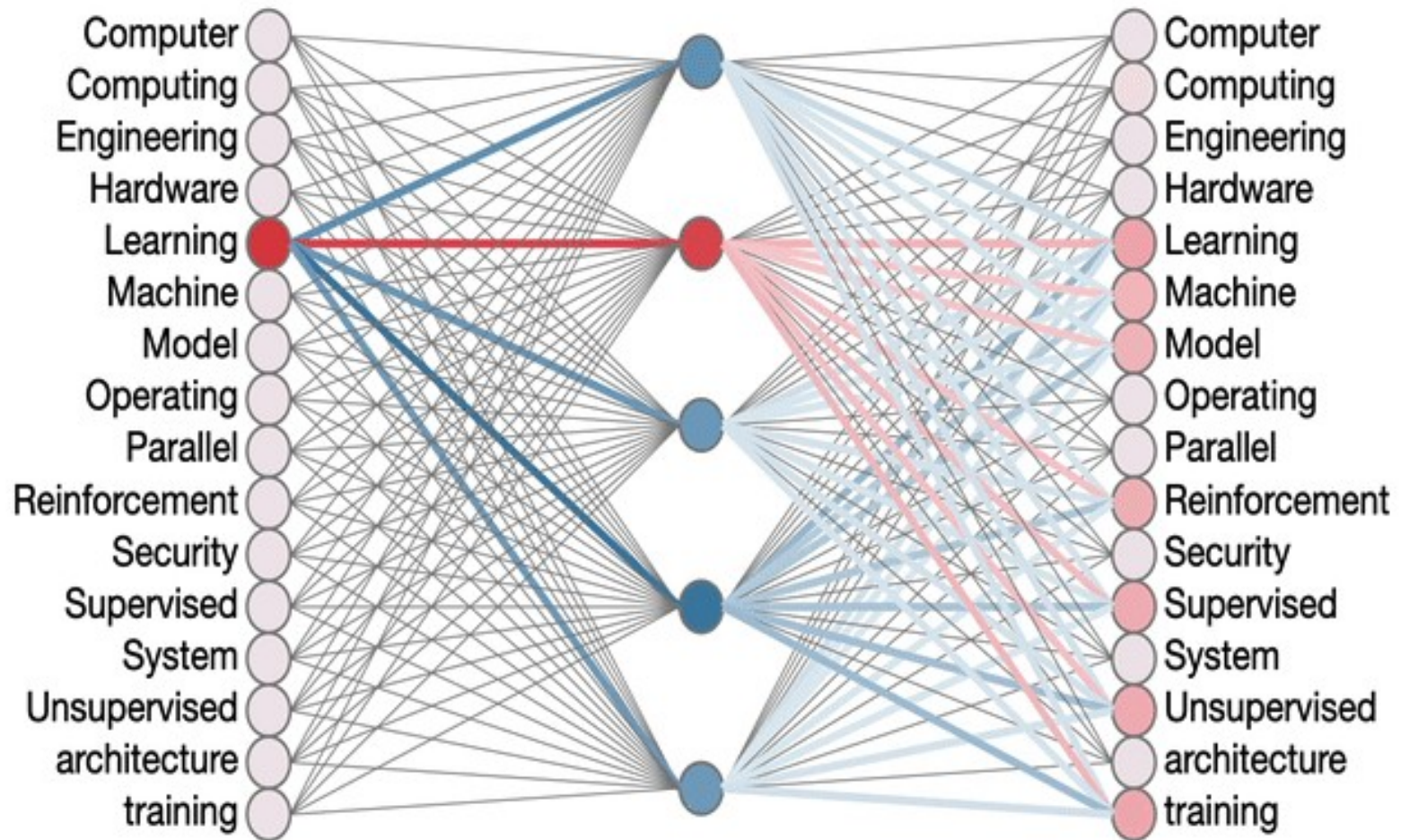
# Representation de textes (vectoriel)

## Corpus pré-traité:

avignon capitale vaucluse jouir climat  
courage étudiant université récompense effort  
espérer réussir cours

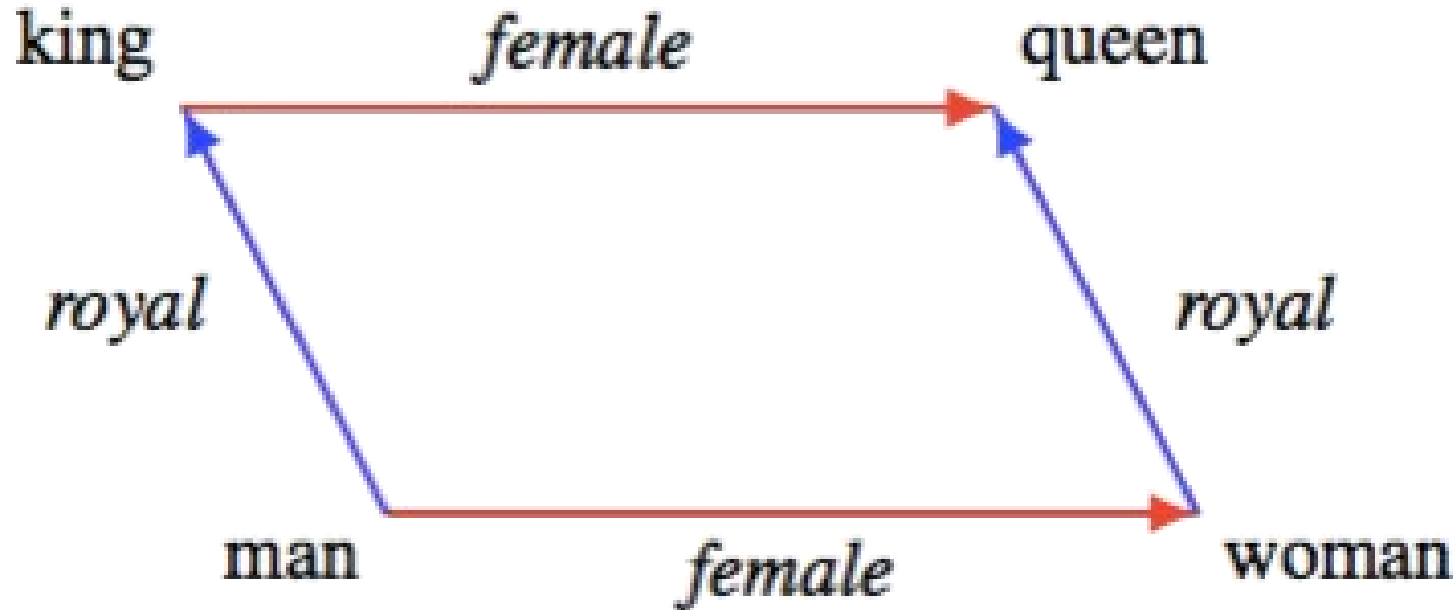
[illegible]

# Représentation embeddings



<https://medium.com/@evertongomede/embeddings-in-machine-learning-unleashing-the-power-of-representation-2402bab526fe>

# Word Embedding Analogies: Understanding King - Man + Woman = Queen





# Artex + IA

- Construire Artex avec embeddings
- Focaliser sur des endroits “intéressants” des textes
- Zones intéressantes dans le cadre du défi = zones avec des EN (Noms des personnes)
- Explorer si vaut la peine de chercher sur tout le texte ou sur un résumé du texte...