

# Systemes de spin appliqués au TALN

---

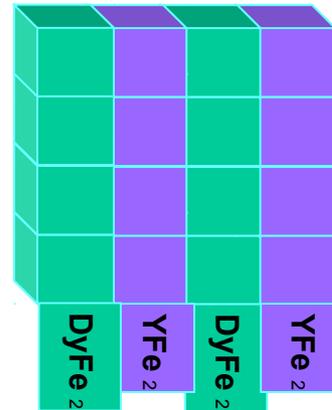


# Nouveaux matériaux magnétiques

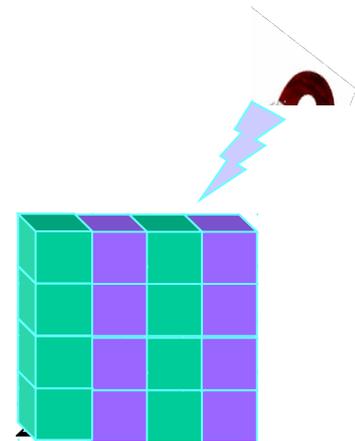
1)

Métaux	M	<b>T</b> (300K)
Terres rares	<b>M</b>	T (4K)

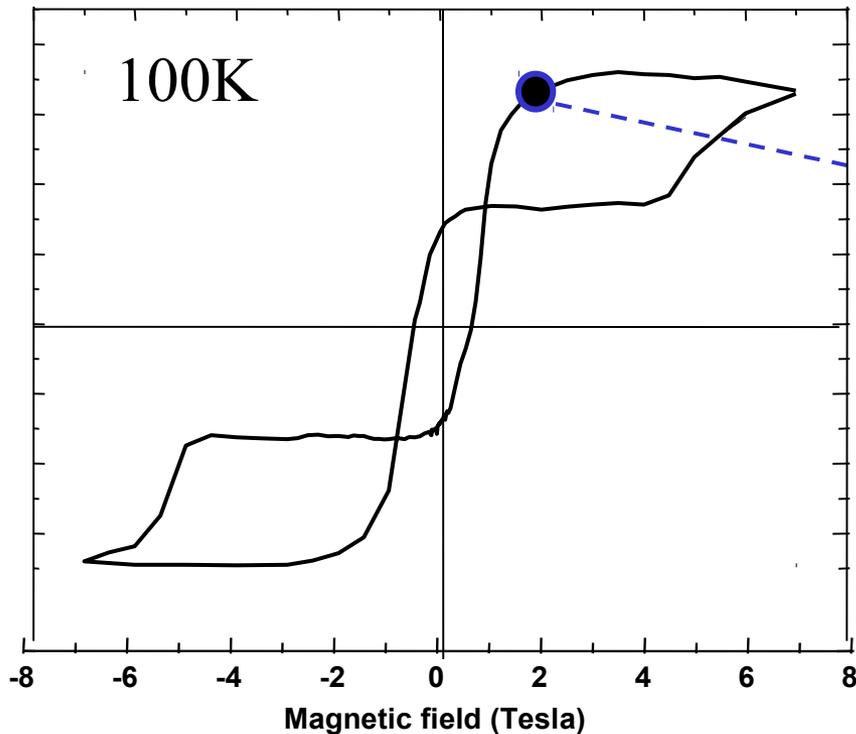
2) Epitaxie par jet moléculaire



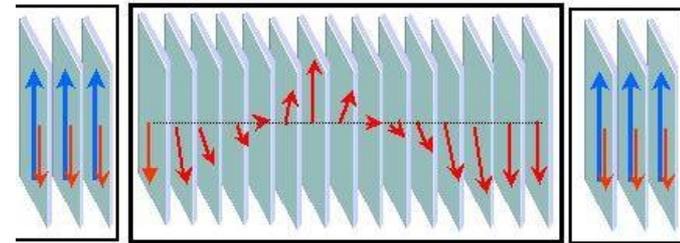
3) Mesures magnétiques



# Mesure magnétique et configuration de spin



**Spin** : représentation de chaque atome comme un petit aimant



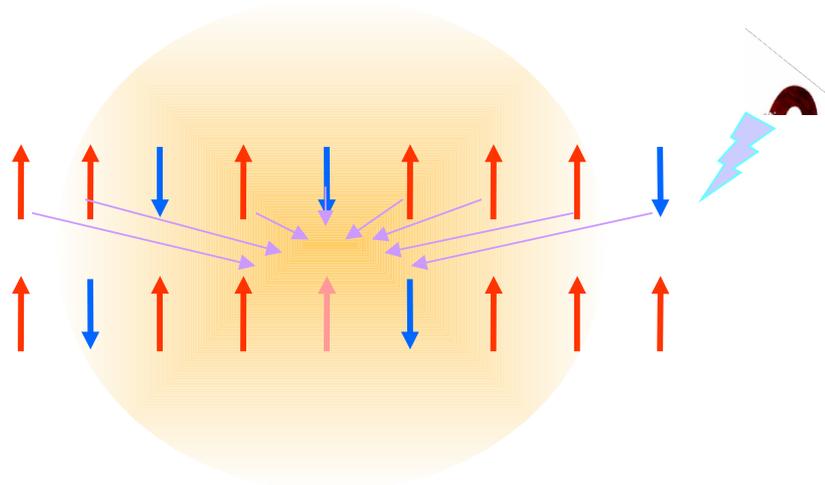
Modèles théoriques de la Physique Statistique:

Modèle d'Ising: deux orientations possibles  $\downarrow \uparrow$

# Energie du système

$$E = E(\text{interactions}) + E(\text{champ})$$

$$E_{ij} = J_{ij} s_i s_j \quad + \quad E_i = H s_i$$
$$J_{ij} = J_{ji}$$

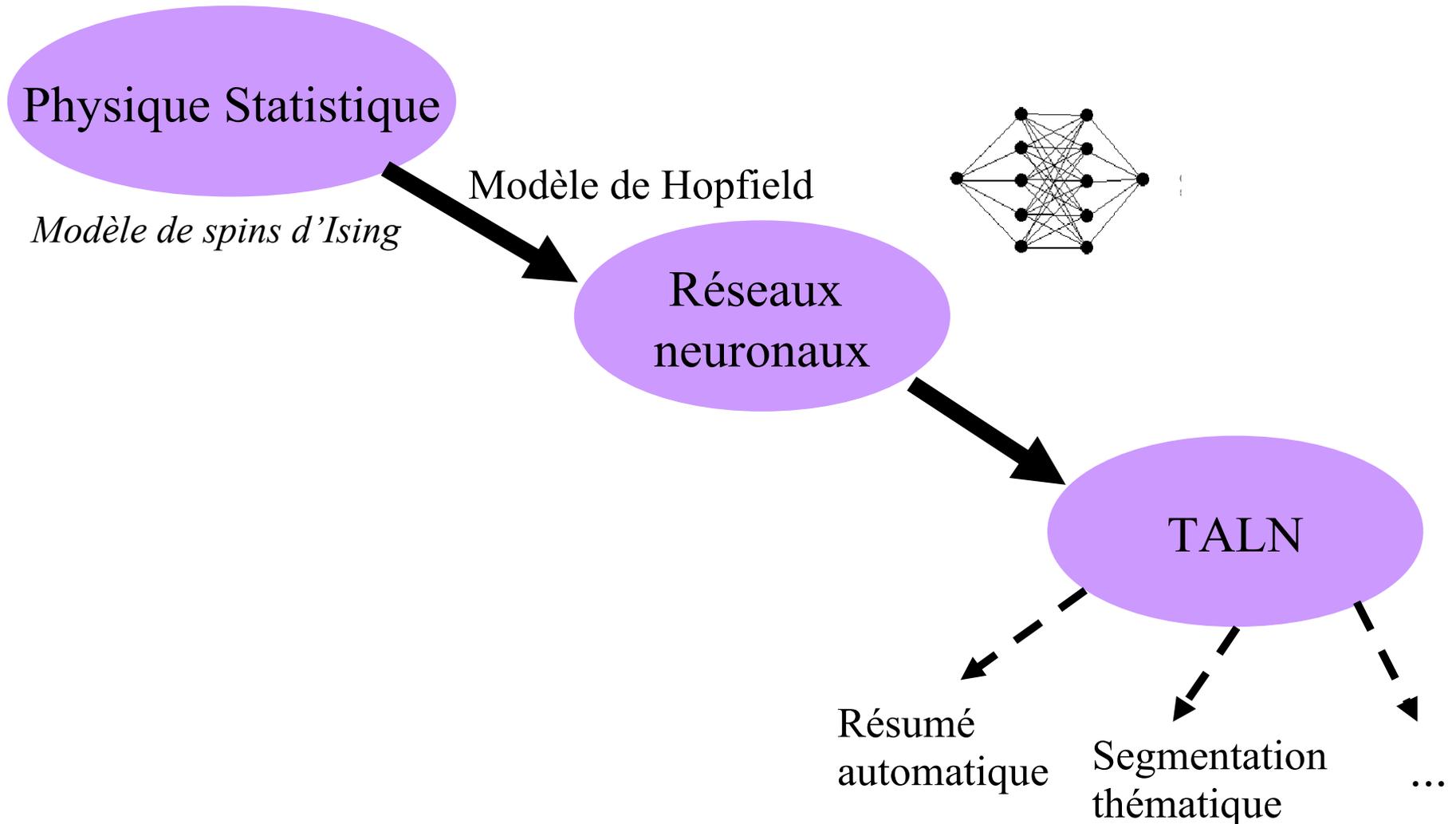


Configuration de spin final : minimisation de  $E$

Prob (état du système) =  $f(E, T, Z)$  ;  $Z$  = fonction de partition

**Mais... où entre le TALN  
dans toute cette histoire ?**

# Applications *exotiques* de la physique statistique



# Mémoire associative

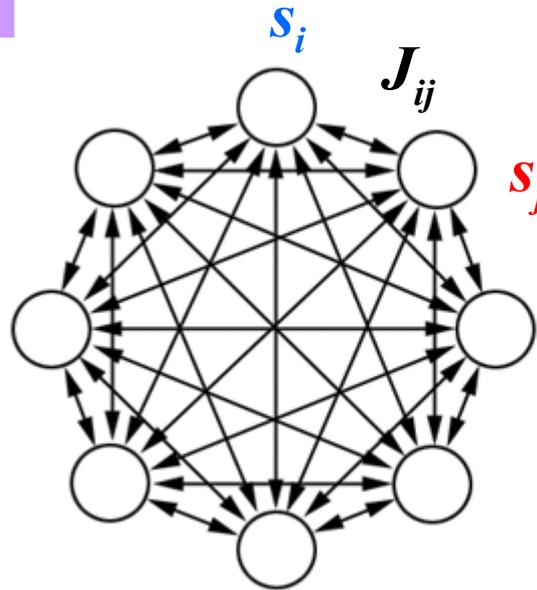
(Hopfield, 1982)

## Modèle de spins d'Ising

neurone = spin  $\downarrow \uparrow$

$$J_{ij} = J_{ji}$$

$$E_{ij} = J_{ij} s_i s_j$$



## Réseaux neuronaux

Règle d'Hebb

$$J_{ij} = s_i s_j$$

Apprentissage

Récupération: minimisation de  $E$

Limitations :

- Patrons corrélés  $\rightarrow$  erreur de récupération
- Capacité  $\approx 0,14 N$

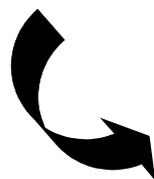
Les maisons bleues de ma tante.  
 Ma tante s'appelle Lulu.  
 J'adore sa maison.  
 Le bleu c'est ma couleur préférée.  
 J'ai des chaussures bleues toutes  
 neuves.

# Codage des documents comme un système de spins!

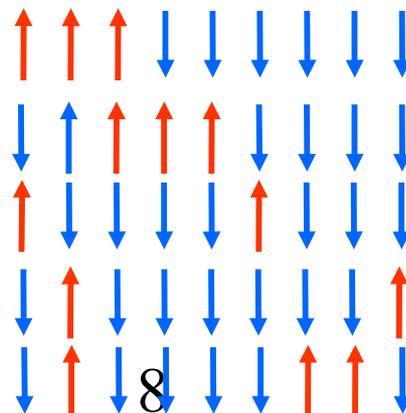


	maison	bleu	tante	appeler	lulu	adorer	neuf	chaussure	couleur
maison	TF	TF	TF	0	0	0	0	0	0
bleu	0	0	TF	TF	TF	0	0	0	0
tante	TF	0	0	0	0	TF	0	0	0
appeler	0	TF	0	0	0	0	0	0	TF
lulu	0	TF	0	0	0	0	TF	TF	0

Corrélés!!



- Modèle vectoriel (sac-à-mots)
- Mots filtrés, normalisés et lemmatisés  
*(Porter, 1980; Manning & Schutze, 2000)*



# Interaction entre spins

mot  $\sim$  neurone  $\sim$  spin  $s_i$

$$[\text{TF TF } \mathbf{0} \dots \mathbf{0}] = [s_0 \ s_1 \ s_2 \dots s_N]$$

Phrase  $\sim$  chaîne de spins



$$S = \begin{pmatrix} s_1^1 & s_2^1 & \dots & s_N^1 \\ s_1^2 & s_2^2 & \dots & s_N^2 \\ \vdots & \vdots & \ddots & \vdots \\ s_1^P & s_2^P & \dots & s_N^P \end{pmatrix}$$

Phrases x mots

$$J^\mu = \begin{pmatrix} s_1^\mu \\ \vdots \\ s_i^\mu \\ \vdots \\ s_N^\mu \end{pmatrix} \times (s_1^\mu \ \dots \ s_i^\mu \ \dots \ s_N^\mu) = \begin{pmatrix} j_{1,1}^\mu & j_{1,j}^\mu & \dots & j_{1,N}^\mu \\ \vdots & \vdots & \ddots & \vdots \\ j_{i,1}^\mu & j_{i,j}^\mu & \dots & j_{i,N}^\mu \\ \vdots & \vdots & \ddots & \vdots \\ j_{N,1}^\mu & j_{N,j}^\mu & \dots & j_{N,N}^\mu \end{pmatrix}$$

$J = \sum J^\mu = (S^T \times S)$  : c'est la mémoire d'Hopfield !

*L'énergie n'est pas utilisée*

# Energie textuelle

$$E = - \begin{pmatrix} s_1^1 & s_2^1 & \cdots & s_N^1 \\ s_1^2 & s_2^2 & \cdots & s_N^2 \\ \vdots & \vdots & \ddots & \vdots \\ s_1^P & s_2^P & \cdots & s_N^P \end{pmatrix} \times J \times \begin{pmatrix} s_1^1 & s_1^2 & \cdots & s_1^P \\ s_2^1 & s_2^2 & \cdots & s_2^P \\ \vdots & \vdots & \ddots & \vdots \\ s_N^1 & s_N^2 & \cdots & s_N^P \end{pmatrix} = - S \times (S^T \times S) \times S^T$$

$$= -(S \times S^T) \times (S \times S^T)$$

$$= -(S \times S^T)^2$$

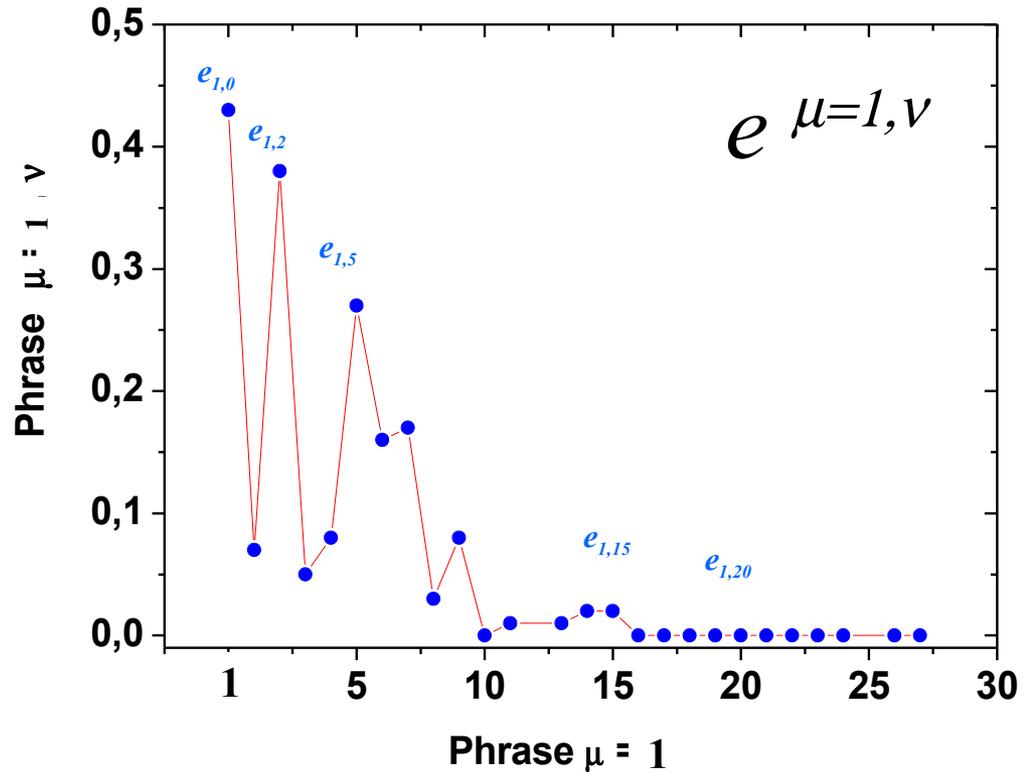
$$E = \begin{pmatrix} e^{1,1} & \cdots & e^{1,P} \\ \vdots & & \vdots \\ e^{\mu,1} & \cdots & e^{\mu,P} \\ \vdots & & \vdots \\ e^{P,1} & \cdots & e^{P,P} \end{pmatrix}$$

$e^{\mu,\nu}$  = énergie entre la phrase  $\mu$  et la phrase  $\nu$

# Energie textuelle

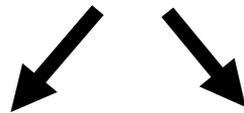
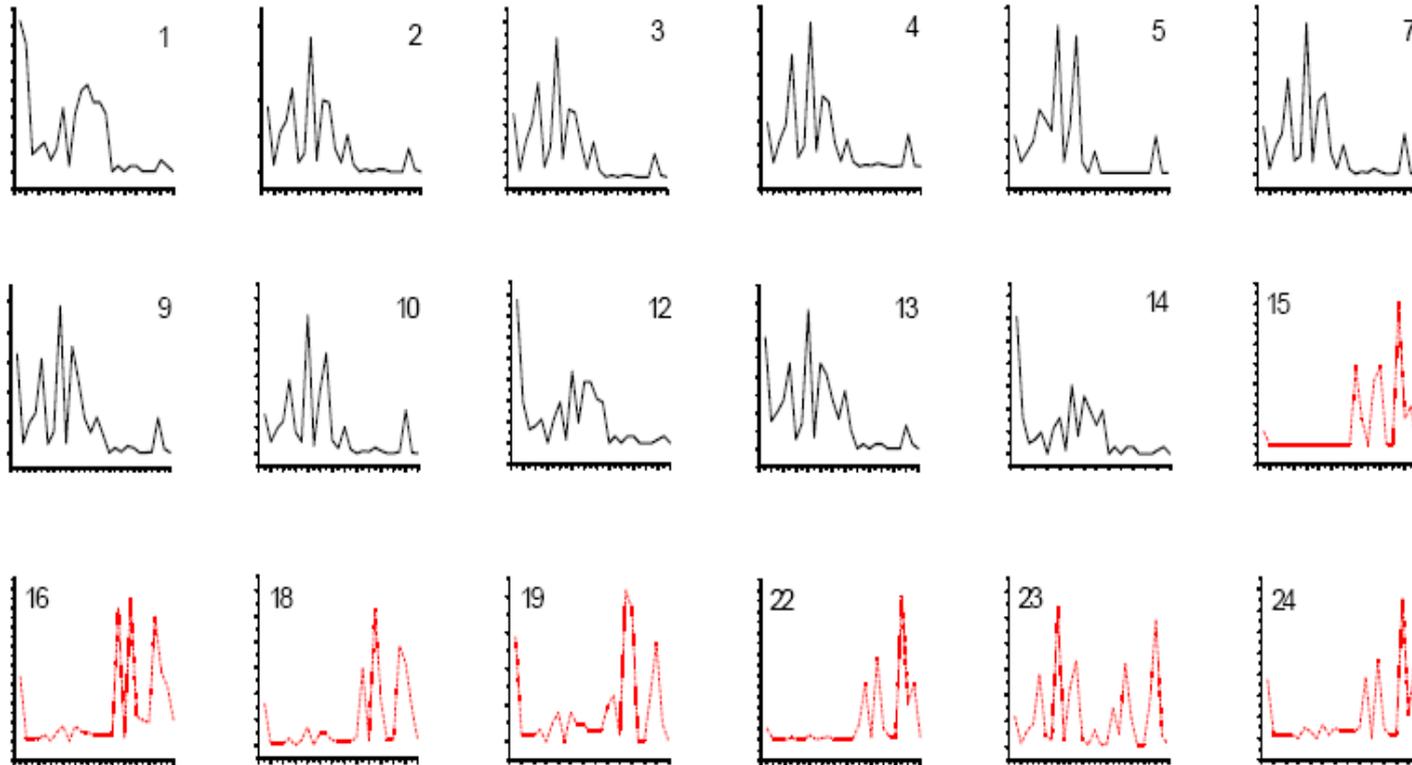
$$E = \begin{pmatrix} e^{1,1} & \dots & e^{1,P} \\ \vdots & & \vdots \\ e^{\mu,1} & \dots & e^{\mu,P} \\ \vdots & & \vdots \\ e^{P,1} & \dots & e^{P,P} \end{pmatrix}$$

$$E^1 = - \sum_v^p e^{v,1}$$



Energie totale de la phrase  $\mu=1$

# Energie textuelle du document



$|E^\mu|$  des phrases  
 $\Rightarrow$  résumé automatique

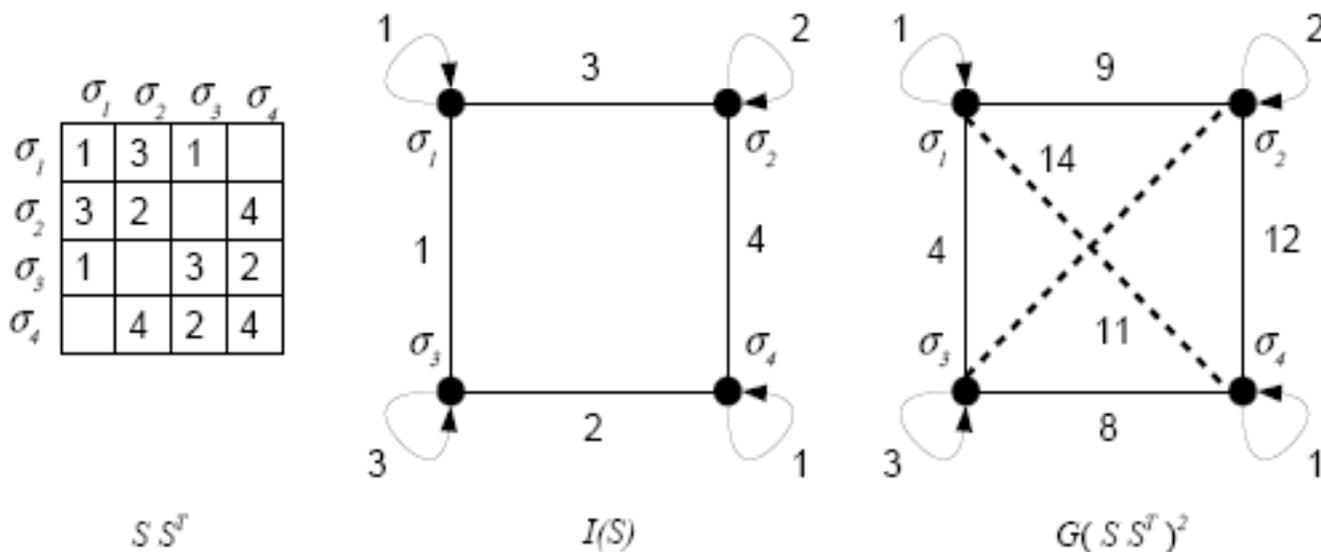
Concordance entre courbes  
 $\Rightarrow$  segmentation thématique

# Interprétation (théorie de graphes)

La matrice d'énergie textuelle s'écrit :

$$E = - S \times (S^T \times S) \times S^T = E = -(S \times S^T)^2$$

Exemple



$\sigma_1 \cap \sigma_4 = \emptyset$  mais  $\sigma_3 \cap \sigma_1 \neq \emptyset$  et  $\sigma_3 \cap \sigma_4 \neq \emptyset$

$\Rightarrow$  l'énergie entre  $\sigma_1$  et  $\sigma_4$  n'est pas nulle

# Somme de trajets de longueur 2 dans le graphe

$$S \times S^T$$



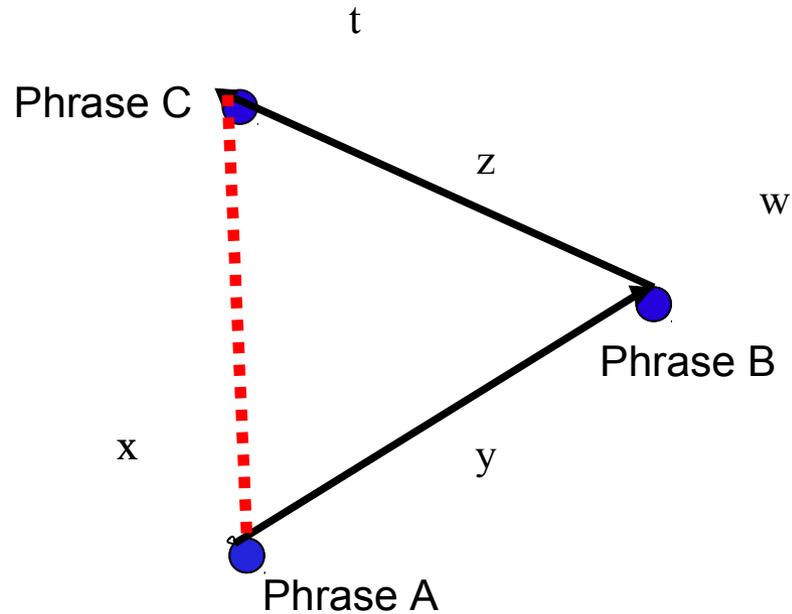
Interactions entre phrases (A,B) et (B,C) ayant des mots en commun

$$(S \times S^T)^2$$



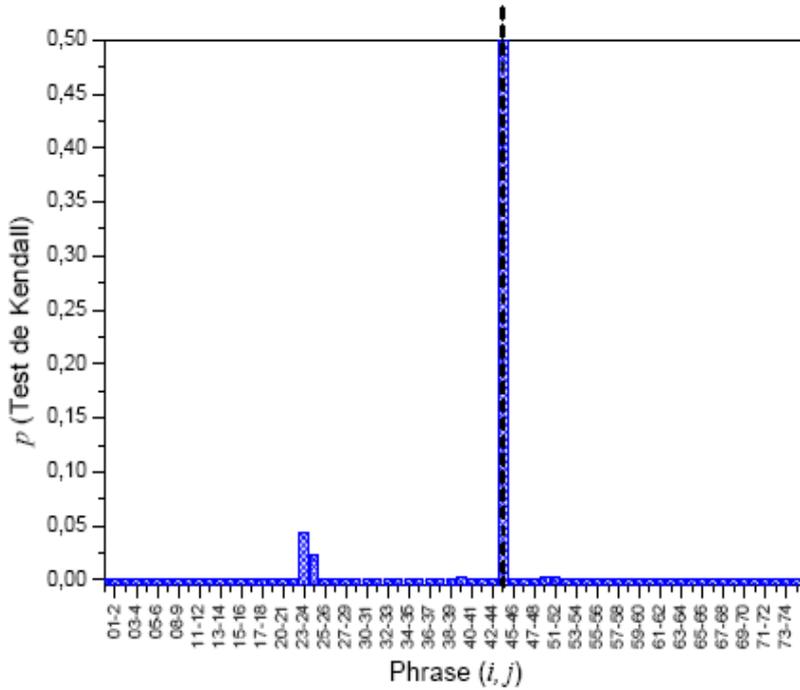
Interactions entre phrases ne partageant pas des mots (A,C) mais ayant des mots en commun avec des *phrases voisines* (B)

$$\text{Coût (A, C)} = y \times z + w \times z + z \times t$$



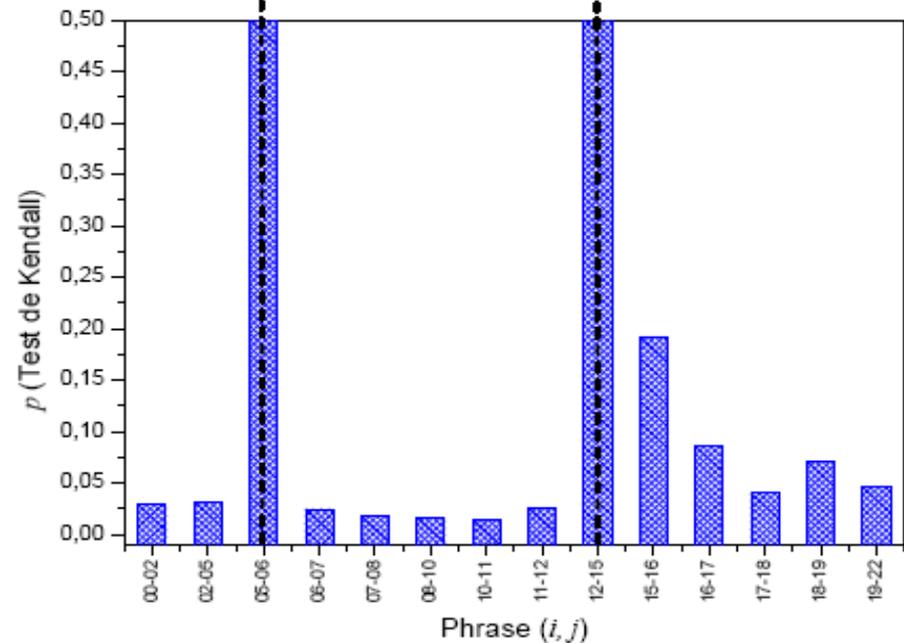
# Résultats

# Frontières thématiques



2 thématiques (en)

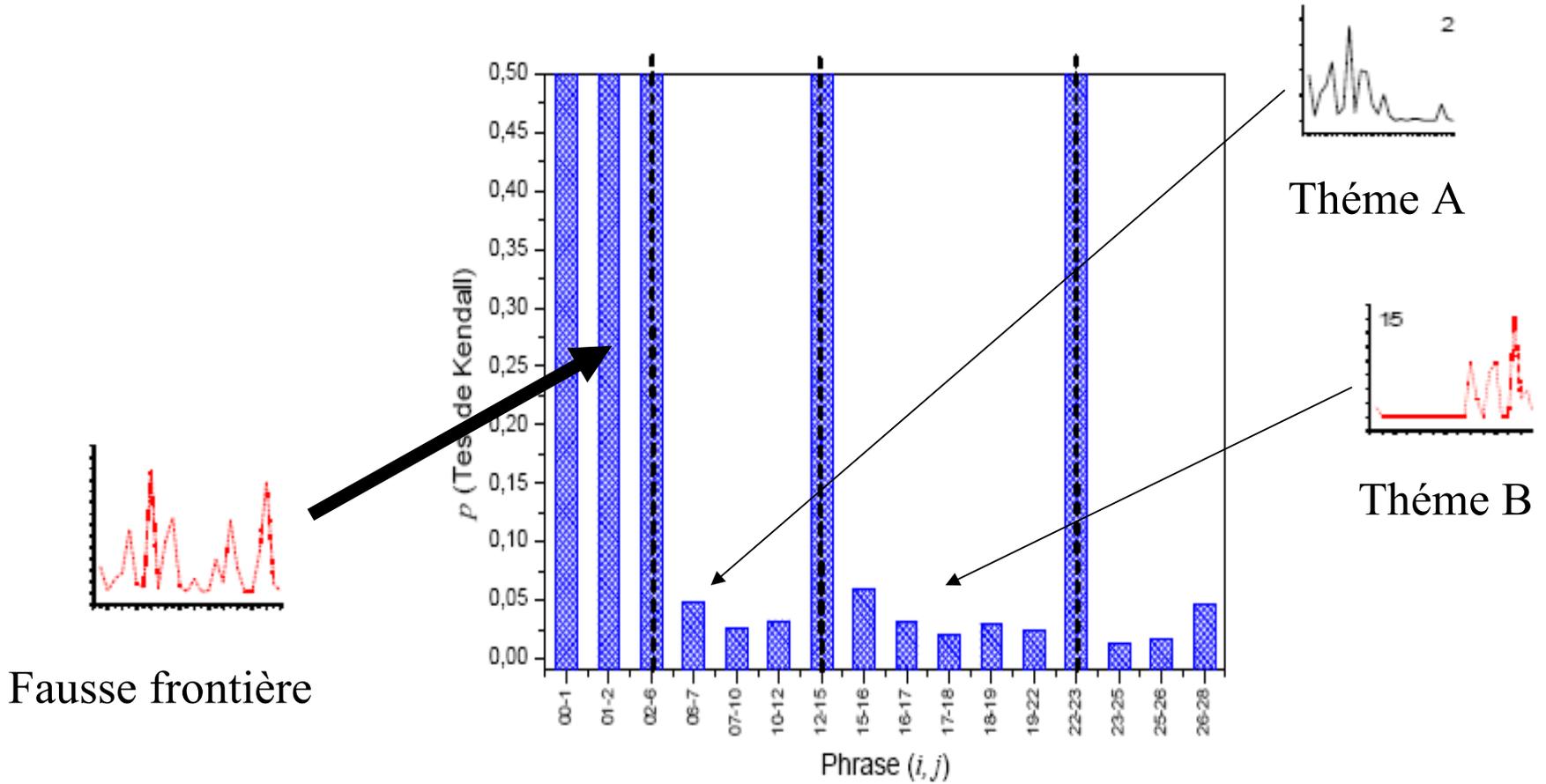
3 thématiques (fr)



Coefficient de concordance  $W$   
de Kendall et sa probabilité  $p$

Test non-paramétrique  
(Siege & Castellan 1988)

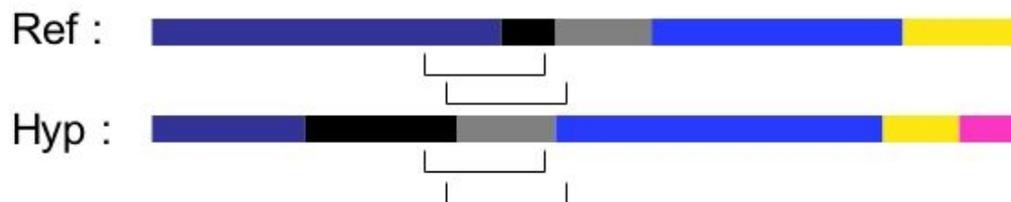
# Erreurs en frontières



Texte avec 4 thématiques

# Frontières thématiques (français)

Evaluation WinDiff (*Pevzner&Hearst, 2002*)



$$\text{WindowDiff}(\text{ref}, \text{hyp}) = \frac{1}{N - k} \sum (|b(\text{ref}_i, \text{ref}_{i+k}) - b(\text{hyp}_i, \text{hyp}_{i+k})|)$$

Corpus (LIA) basé sur la méthode de (*Choi, 2000*) :

200 documents  
10 thématiques par document, *Le Monde 2001*

Taille du segmente	LIA_Seg (Sitbon et Bellot, 2005)	Energie
9-11	<b>0,319</b>	0,449
3-11	<b>0,369</b>	0,409

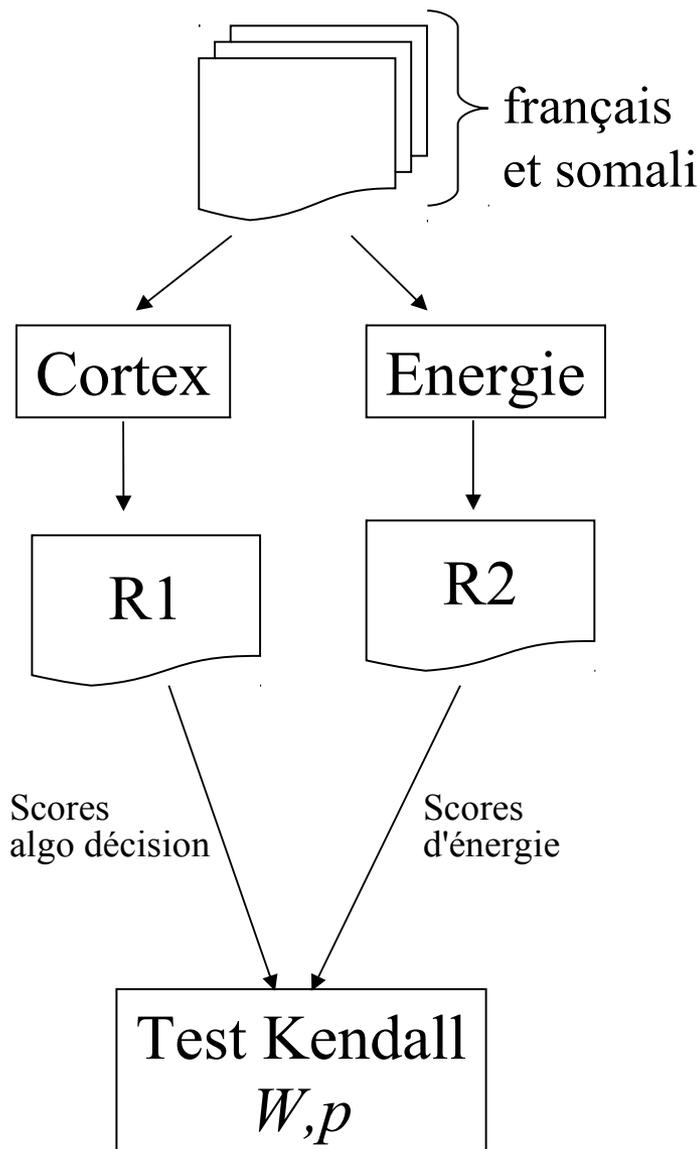
# Résumé générique

*F-score* – Rouge-SU4 normalisé (Lin, 2004)

<b>Document</b>	<b>Energie</b>	<b>Cortex*</b>	<b>Baseline</b>
<b>3-mélanges (web, Fr)</b> 27 phr, 826 mots, 25%, 8 réf	<b>0,47150</b>	0,43068	0,32936
<b>puces (web, Fr)</b> 29 phr, 653 mots, 25%, 8 réf	0,53574	<b>0,55628</b>	0,32723
<b>J'accuse (E. Zola, Fr)</b> 206 phr, 4936 mots, 12%, 6 réf	0,58479	<b>0,60037</b>	0,26152
<b>Lewinsky (Wikipedia, En)</b> 30 phr, 816 mots, 20%, 7 réf	0,47757	<b>0,51076</b>	0,29248
<b>Québec (Wikipedia, En)</b> 44 phr, 1190 mots, 25%, 8 réf	0,51179	<b>0,55656</b>	0,35244

\*Torres et al. 2002

# Résumé : concordance entre Cortex et Energie



Texte	$W$ Kendall	$p$ -value
parlement-fr	0,96847	NA
baarlamanka-som	1,00000	NA
ONU-fr	0,95105	0,02506
ONU-somali	0,91958	0,04227
prince-1-fr	0,88879	0,01087
prince-1-somali	0,97832	0,00339

NA  $\approx$  0

# Conclusion

- Pont entre la Physique Statistique et le TALN
- Notion d'énergie textuelle
- Deux applications par le même prix !
  - Résumé générique
    - comparable au système Cortex (générique) en termes de précision, rappel et  $F$ -score
  - Frontières thématiques
    - Combinaison de l'énergie avec une méthode non paramétrique (test de Kendall)
    - Moins performante que *LIA\_Seg*

# Perspectives

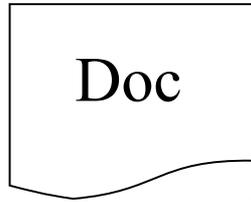
- Résumés guidés par des requêtes (tests sur DUC...)
- Multilingue : anglais, français, somali, espagnol, maya
- Améliorer la détection des frontières

# Perspectives *exotiques*

## Application d'un champ externe

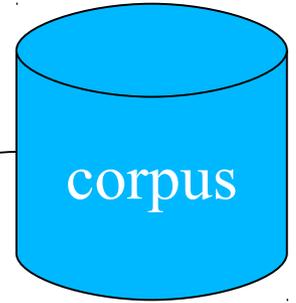
Catégorisation

Résumé  
personnalisé



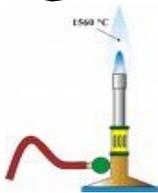
*Champ externe*

$$H = (w_1, w_2, \dots, w_P)$$



Requête

## Introduction d'une température



$T$

$$\text{Prob}(S_i = \pm 1) = 1/(1 + e^{(2\beta h_i)})$$

$\beta = 1/k_B T$  = température inverse

$k_B$  = cte de Boltzmann

$\Rightarrow$  Modifier le paysage d'énergie