

UNIVERSITÉ D'AVIGNON  
ET DES PAYS DE VAUCLUSE



# **COMPRESION AUTOMATICA DE FRASES: ¿CÓMO DECIR ALGO EN MENOS PALABRAS Y AÚN ASÍ DECIRLO BIEN?**

**JUAN-MANUEL TORRES**

**ALEJANDRO MOLINA**

**juan-manuel.torres@univ-avignon.fr**

LABORATOIRE INFORMATIQUE D'AVIGNON

UNIVERSITÉ D'AVIGNON ET DES PAYS DE VAUCLUSE

ÉCOLE POLYTECHNIQUE DE MONTRÉAL

# COMPRESIÓN AUTOMÁTICA DE FRASES

## TAREA SUBJETIVA

TRABAJOS PIONEROS DE MARCU (CANAL RUIDOSO) Y ARBOLES SINTACTICOS

– TAREA BIEN PLANTEADA ?

SE DEBEN ELIMINAR PALABRAS ?

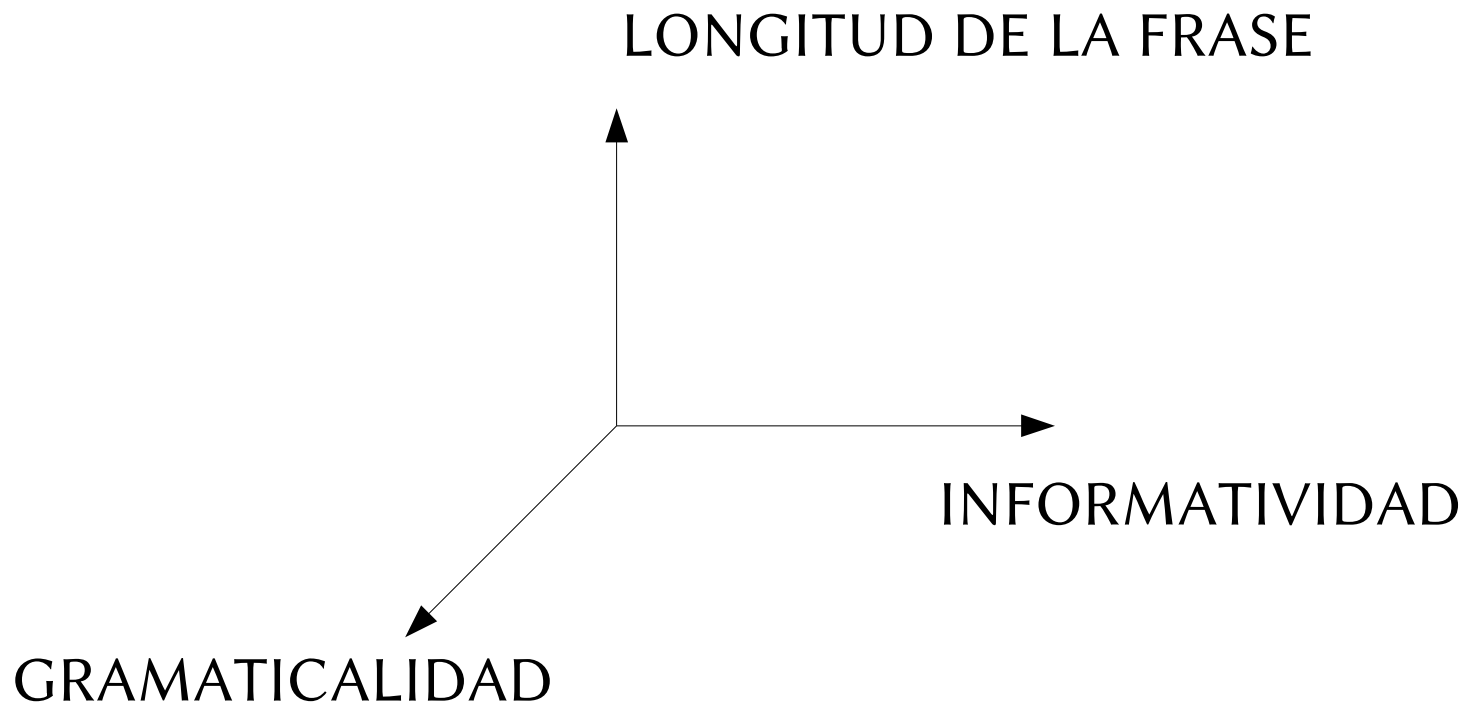
– LAS PALABRAS DEPENDEN DEL CONTEXTO

POR TANTO, ELIMINAR MEJOR SEGMENTOS...

– QUÉ SEGMENTOS ?

# COMPRESIÓN AUTOMÁTICA DE FRASES

PROBLEMA A TRIPLE DIMENSION



# COMPRESIONES GRAMATICALES... PERO QUÉ TAN INFORMATIVAS SON?

$\varphi$  : Juliette prépare un gâteau pour le manger bien qu'elle n'ait pas faim.

$\tilde{\varphi}_1^*$  : Juliette prépare un gâteau bien qu'elle n'ait pas faim.

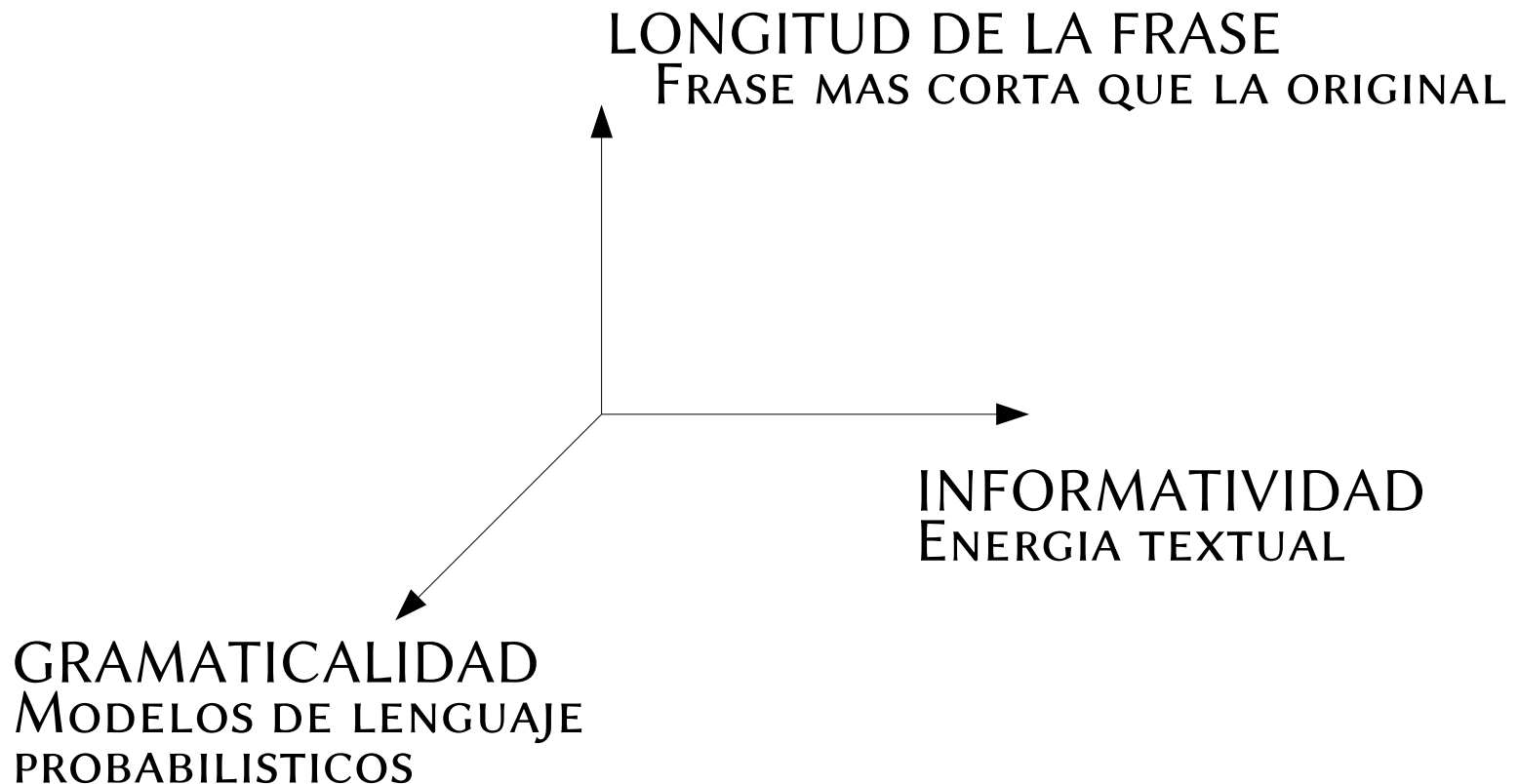
$\tilde{\varphi}_2^*$  : Juliette prépare un gâteau pour le manger.

$\tilde{\varphi}_3^*$  : Juliette prépare un gâteau.

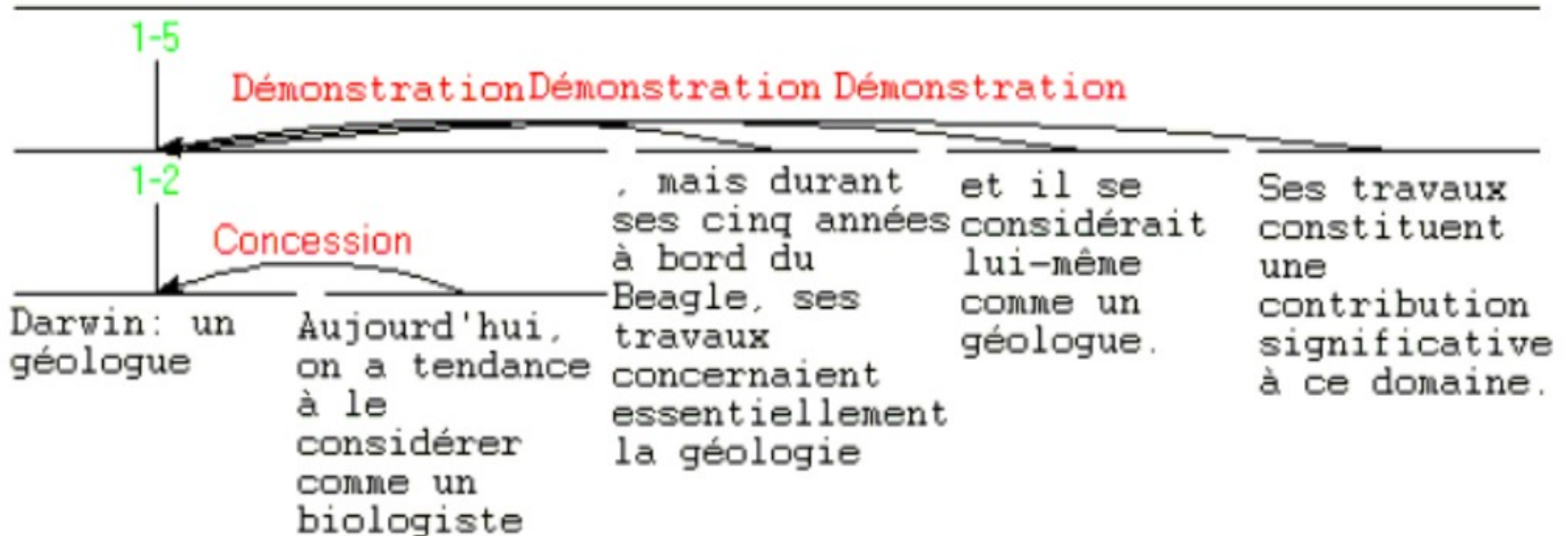
## DIMENSIONES ANTAGONISTAS (ORTOGONALES)

# COMPRESIÓN AUTOMÁTICA DE FRASES

PROBLEMA A TRIPLE DIMENSION



# INFORMATIVIDAD : INFORMACION DISCURSIVA



# **INFORMATIVIDAD :**

## **SEGMENTOS DISCURSIVOS**

**S0.** DARWIN : UN GÉOLOGUE.

**S1.** AUJOURD'HUI ON A TENDANCE A LE CONSIDÉRER COMME UN BIOLOGISTE

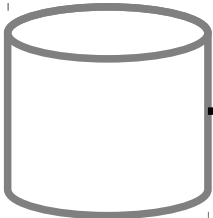
**S2.** MAIS SES 5 ANNÉES À BORD DU BEEGLE, SES TRAVAUX CONCERNANT ESSENTIELLEMENT LA GÉOLOGIE

**S3.** ET IL SE CONSIDÉRerait LUI-MÊME COMME UN GÉOLOGUE.

**S4.** SES TRAVAUX CONSTITUENT UNE CONTRIBUTION SIGNIFICATIVE À CE DOMAINE.

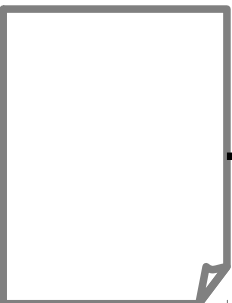
**SEGMENTADOS AUTOMATICAMENTE CON HERRAMIENTAS  
HECHAS EN « CASA » (FRANCIA/MEXICO)**

# SEGMENTACION DISCURSIVA ELEMENTAL



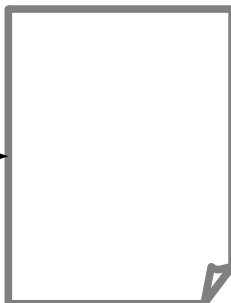
LISTA DE MARCADORES DISCURSIVOS

TEXTO UTF8



SEGMENTADOR FRASES

SEGMENTADOR DISCURSIVO

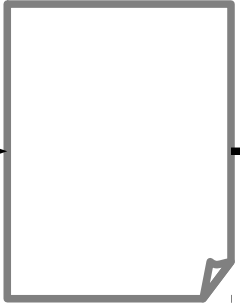


TEXTO SEGMENTADO EDUs

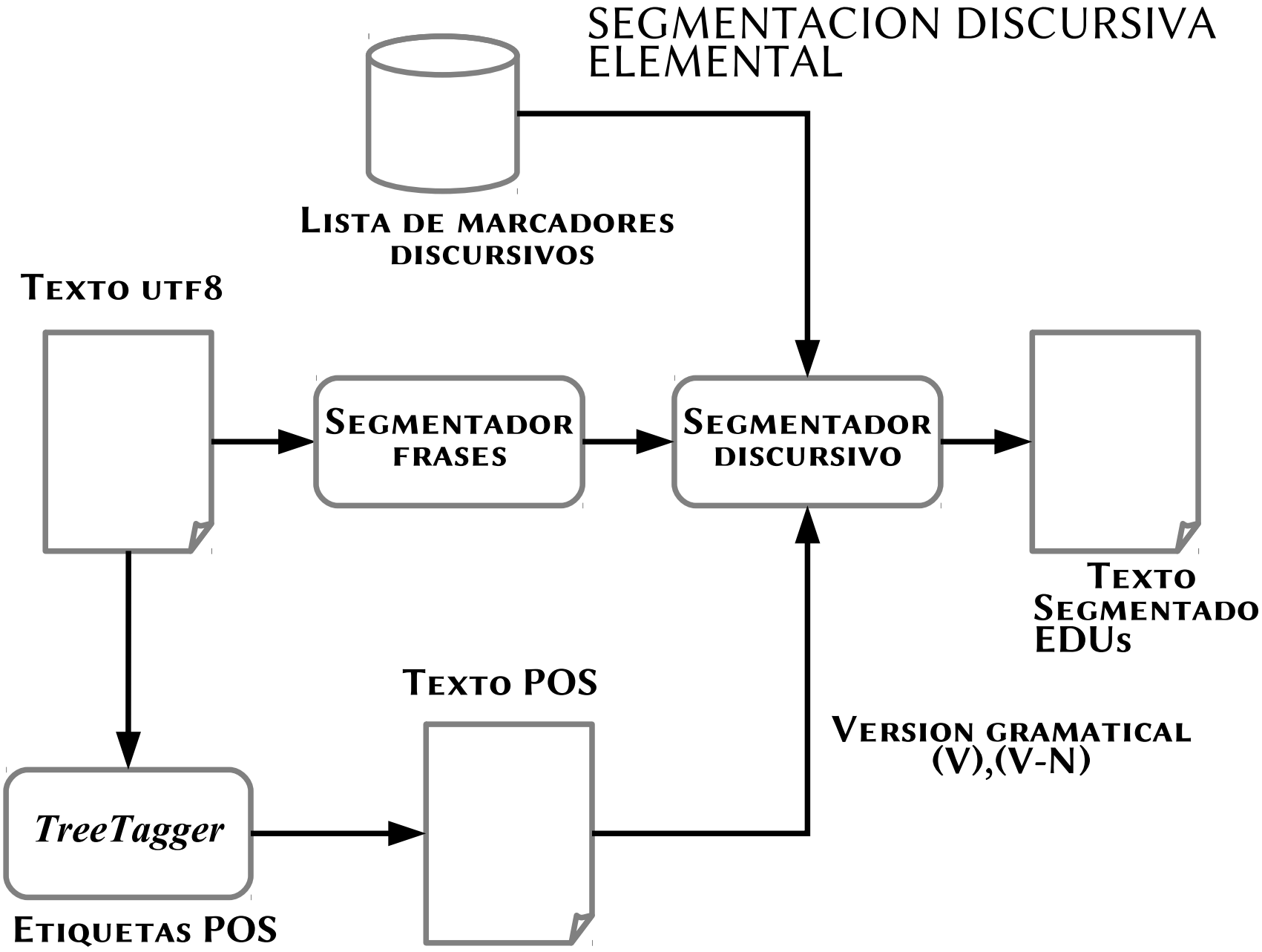
TEXTO POS

*TreeTagger*

ETIQUETAS POS



VERSION GRAMATICAL (V),(V-N)





# MODELOS DE LENGUAJE: PROBABILIDAD DE EXISTENCIA DE LA FRASE

$$\begin{aligned} P(\text{Juliette, prépare, un, gâteau, pour, le, manger}) &= P(\text{Juliette}) \times \\ &\quad P(\text{prépare} \mid \text{Juliette}) \times \\ &\quad P(\text{un} \mid \text{Juliette prépare}) \times \\ &\quad P(\text{gâteau} \mid \text{Juliette prépare un}) \times \\ &\quad P(\text{pour} \mid \text{Juliette prépare un gâteau}) \times \\ &\quad P(\text{le} \mid \text{Juliette prépare un gâteau pour}) \times \\ &\quad P(\text{manger} \mid \text{Juliette prépare un gâteau pour le}) \end{aligned}$$

$$P(w_1^n) = P(w_1) \times P(w_2|w_1) \times P(w_3|w_1^2) \times \dots \times P(w_n|w_1^{n-1}).$$

PROBABILIDADES ESTIMADAS SOBRE UN CORPUS  
REPRESENTATIVO : GOOGLE 5-GRAMMES EN FR ES

# COMPRESIÓN DE FRASES: MODELO LINEAL COMBINANDO PARÁMETROS

$$P_{\text{elim}}(s, \varphi) = \text{Ener}(s, \varphi) + \text{Gram}(s, \varphi) + \text{Seg}(s, \varphi) + \text{Lon}(s, \varphi)$$

Ener ~ informatividad

Gram ~ Gramaticalidad

Seg ~ Segmentador

Lon ~ Longitud

**GRAN INDEPENDENCIA DEL IDIOMA**

# EVALUACIÓN DE COMPRESIONES

## MÉTODOS CLÁSICOS

- DE TRADUCCIÓN: BLEU
- DE RESUMEN SEMI-AUTOMÁTICOS: ROUGE (LIN 2007)
- DE RESUMEN AUTOMÁTICOS: FRESA (TORRES ET AL 2010, SAGGION ET AL. 2011)

**TEST DE TURING REVISITADO (MOLINA, SANJUAN & TORRES, 2013)**

# Test de Turing (*The Imitation Game*)



Alan Turing et son test de éponyme

# **DESCUBRIMIENTO DE MAMUT EMOCIONA A CIENTÍFICOS**

## **(DOCUMENTO FUENTE)**

El inusual descubrimiento de los huesos de un mamut en una finca de Oskaloosa ha emocionado a los expertos que estudian la vida prehistórica por los descubrimientos científicos que podrían hacer se con la enorme bestia. El hallazgo es raro porque al parecer buena parte del esqueleto del animal se encuentra en excelente estado. Eso permite a los científicos recolectar polen y evidencias de otras plantas en el sitio de la excavación que podrían revelar detalles respecto del medio ambiente de Iowa hace más de 12000 años. Los científicos de la Universidad de Iowa planean escanear el lugar el viernes con un radar que penetra en la tierra para ver si pueden determinar qué tanto de los restos del mamut siguen enterrados. La excavación se extenderá por varios meses. Los huesos los encontró hace dos años el propietario del terreno, quien desea que tanto su nombre como la ubicación de la granja sean confidenciales para proteger el sitio.

# DESCUBRIMIENTO DE MAMUT EMOCIONA A CIENTÍFICOS

## (DOCUMENTO FUENTE)

El inusual descubrimiento de los huesos de un mamut en una finca de Oskaloosa ha emocionado a los expertos que estudian la vida prehistórica por los descubrimientos científicos que podrían hacer se con la enorme bestia. El hallazgo es raro porque al parecer buena parte del esqueleto del animal se encuentra en excelente estado. Eso permite a los científicos recolectar polen y evidencias de otras plantas en el sitio de la excavación que podrían revelar detalles respecto del medio ambiente de Iowa hace más de 12000 años. Los científicos de la Universidad de Iowa planean escanear el lugar el viernes con un radar que penetra en la tierra para ver si pueden determinar qué tanto de los restos del mamut siguen enterrados. La excavación se extenderá por varios meses. Los huesos los encontró hace dos años el propietario del terreno, quien desea que tanto su nombre como la ubicación de la granja sean confidenciales para proteger el sitio.

# DESCUBRIMIENTO DE MAMUT EMOCIONA A CIENTÍFICOS

## (RESUMEN POR COMPRESION)

---

Los científicos de la Universidad de Iowa planean escanear el lugar el viernes con un radar que penetra en la tierra. La excavación se extenderá por varios meses. Los huesos los encontró hace dos años el propietario del terreno.

---

*20 % DEL TAMANO ORIGINAL, SEGMENTADOR DISEG  
RESUMIDOR ENERTEX*

# TEST DE TURING

Réponse du juge	Origine du résumé	
	Humain	Machine
Humain	a	b
Machine	c	d

Juge id	Contingence		$p$	$H_0$
Juge 1	4	0	0.030	faux
	2	6		
Juge 4	1	5	0.998	vrai
	5	1		
Juge 7	5	5	0.772	vrai
	1	1		

$H_0$  : INDEPENDENCIA : NO HAY ASOCIACION ENTRE EL ORIGEN DEL RESUMEN Y LAS RESPUESTAS

$H_1$  : EL JUEZ IDENTIFICA EL ORIGEN DEL RESUMEN

RESULTADOS : SOBRE 54 HUMANOS, 53 FUERON INCAPACES DE DESCUBRIR EL ORIGEN ARTIFICIAL DE LOS RESUMENES

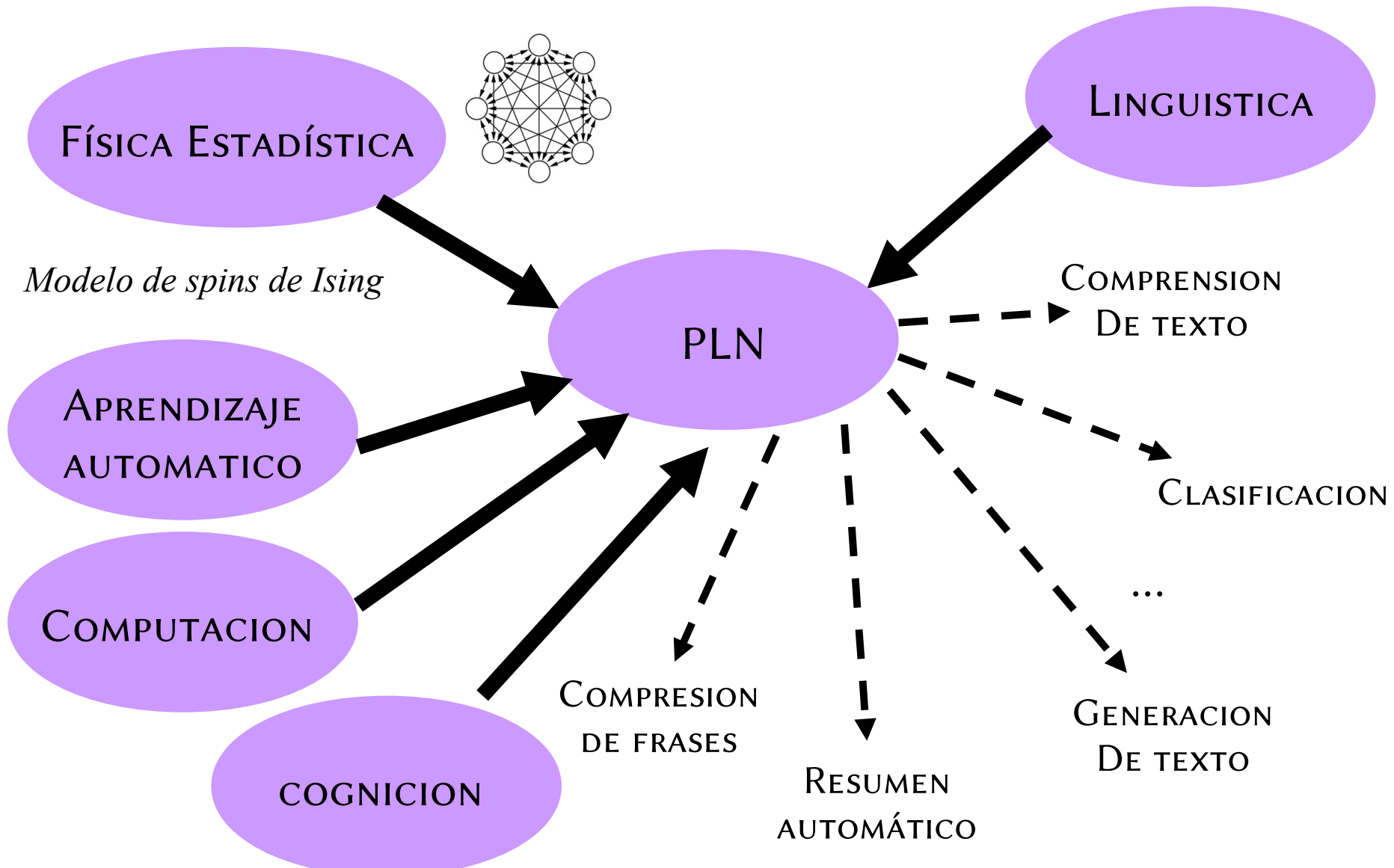
**(P-VALUE = 0.496 > 0.05 SE ACEPTA  $H_0$ )**



LOS DOCUMENTOS SON INFORMATIVOS,  
PERO...  
POSEEN PROPIEDADES FISICAS...?  
VOLUMEN? MASA? LONGITUD?  
ENERGIA...?

PROBABLEMENTE PUEDAN USARSE UNA  
TRANSPOSICION DE IDEAS  
DE OTRAS DISCIPLINAS...

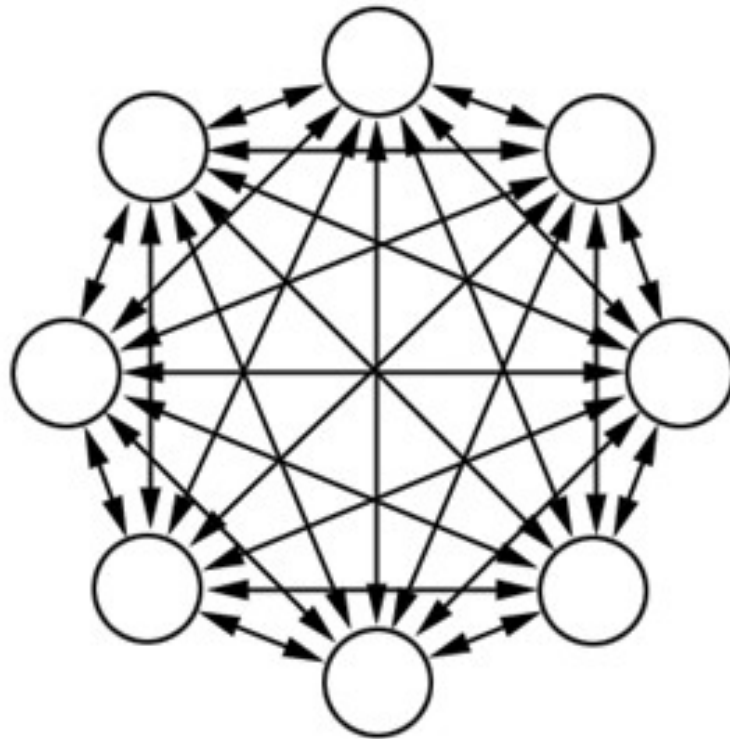
# Como estudiar el lenguaje humano?



# **INGREDIENTES BASICOS...**

## **MECANICA ESTADISTICA :**

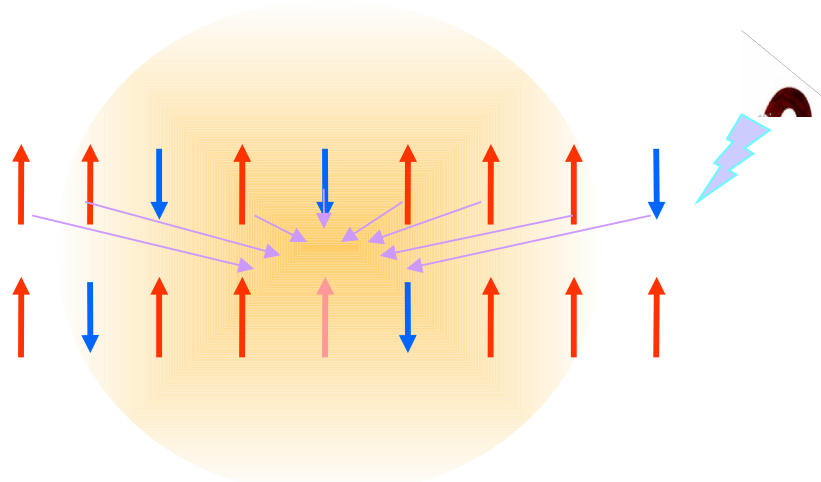
### **ENERGIA TEXTUAL**



# ENERGÍA DEL SISTEMA

$$E = E(\text{INTERACCIONES}) + E(\text{CAMPO})$$

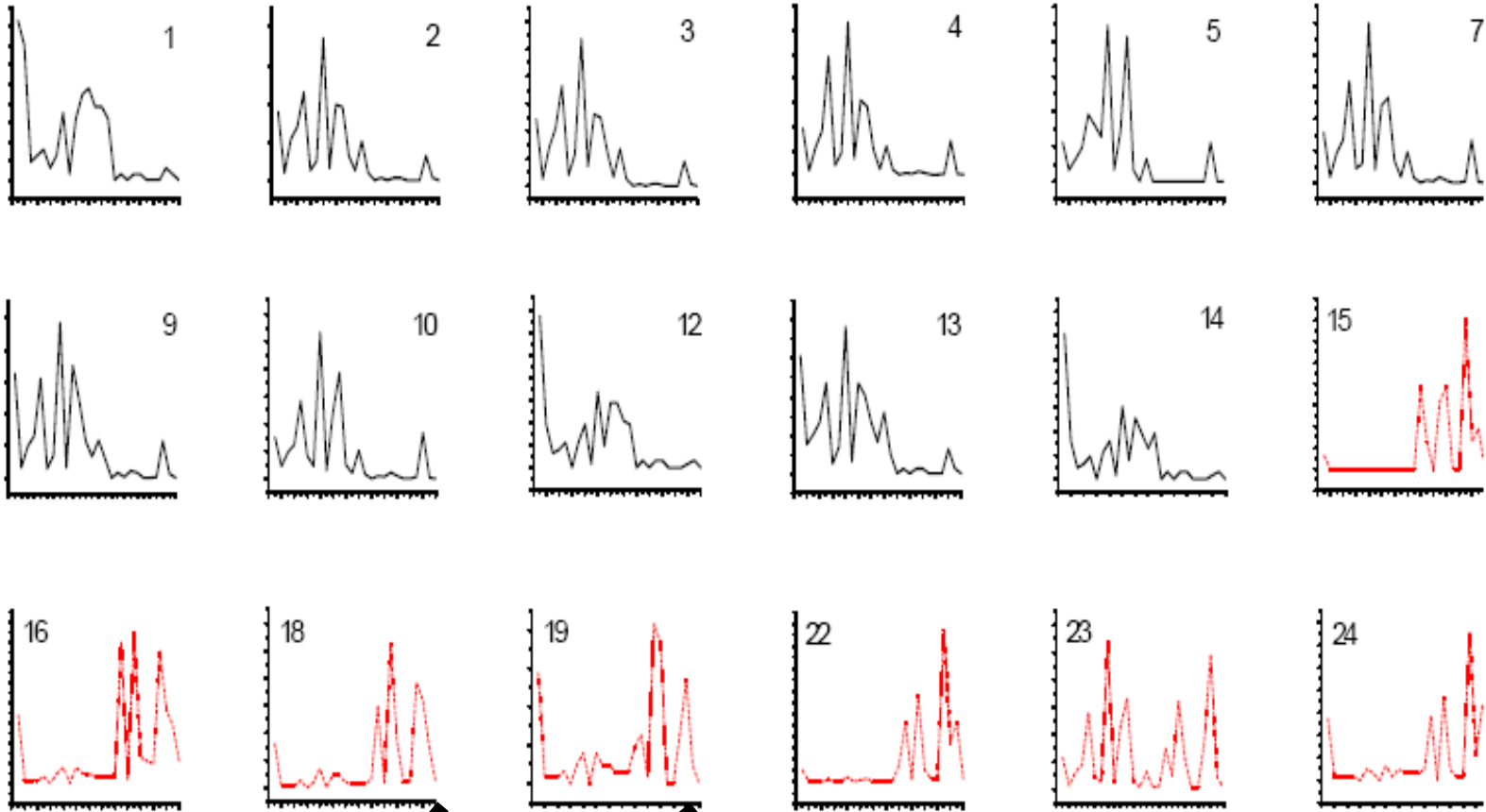
$$E_{ij} = J_{ij} s_i s_j \quad + \quad E_i = H s_i$$
$$J_{ij} = J_{ji}$$



CONFIGURACIÓN DE SPIN FINAL : MINIMIZACIÓN DE  $E$

$p(\text{estado del sistema}) = f(E, T, Z)$  ;  $Z$ =función de partición ;  
T = temperatura

# ENERGÍA TEXTUAL



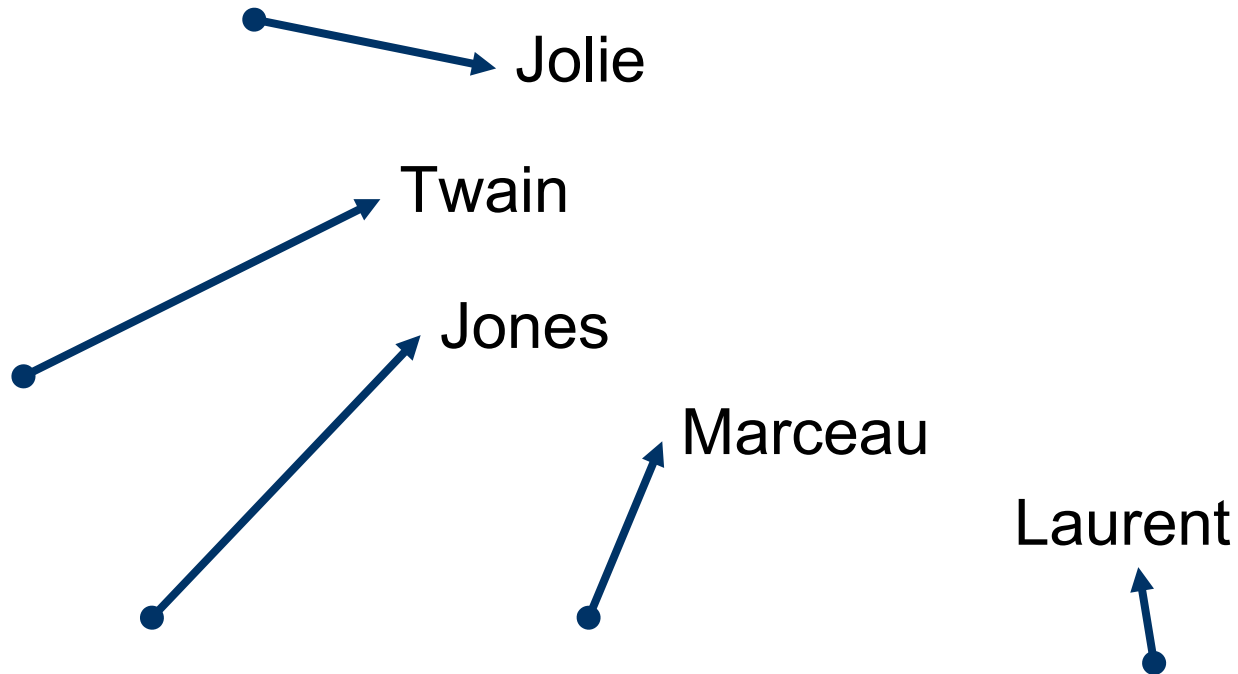
$|E^\mu|$  de frases :  
Resumen automatico

Concordancia de curvas :  
Segmentacion tematica

# INGREDIENTES BASICOS:

IDEAS DE COGNICION...

MEMORIAS ASOCIATIVAS



# MEMORIAS ASOCIATIVAS



Angelina Jolie

Shania Twain



Catherine Z Jones

Sophie Marceau



Mélanie Laurent



# **INGREDIENTES BASICOS**

**LINGUISTICA**

**ANALISIS DISCURSIVO**

**CORPUS**

**COMPUTACION**

**MATEMATICAS**

**INGENIERIA...**



# APLICACIONES: DETECCIÓN DE FRONTERAS TEMÁTICAS

- SEPARAR TEMÁTICAMENTE DOCUMENTOS
- CORPUS
  - POLITICA | CIENCIA | ARTE | DEPORTES | CULTURA
- TAREA CLÁSICA DE PLN
- INDEPENDIENTE DEL IDIOMA (TRILINGÜE EN/FR/ES)

# **APLICACIONES: DETECCIÓN DE SIMILITUD TEXTUAL**

**LA ENERGIA TEXTUAL PUEDE SERVIR PARA  
DETECTAR SIMILITUD TEXTUAL...**

**EN PARTICULAR PARAFRASIS, INDEPENDIENTE  
DEL IDIOMA, TEMATICA Y CONTENIDO**

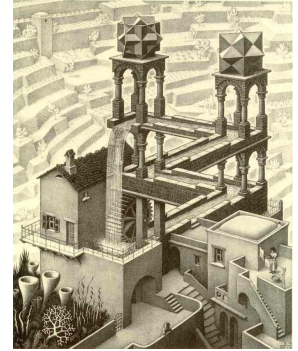
POR MI RAZA HABLARA EL ESPIRITU  
EL ESPIRITU VA A HABLAR POR MI RAZA  
QUIEN HABLARA POR MI RAZA SERA EL ESPIRITU  
EL ESPIRITU SERA QUIEN HABLE POR NOSOTROS

...

# **MAS APLICACIONES**

**RESUMEN AUTOMATICO,  
GENERACION DE TEXTO,  
CLASIFICACION (TWEETS, BLOGS,  
DOCUMENTOS,...),  
IDENTIFICACION DE ESTILOS,  
ANALISIS AUTOMATICO DE CV,  
...**

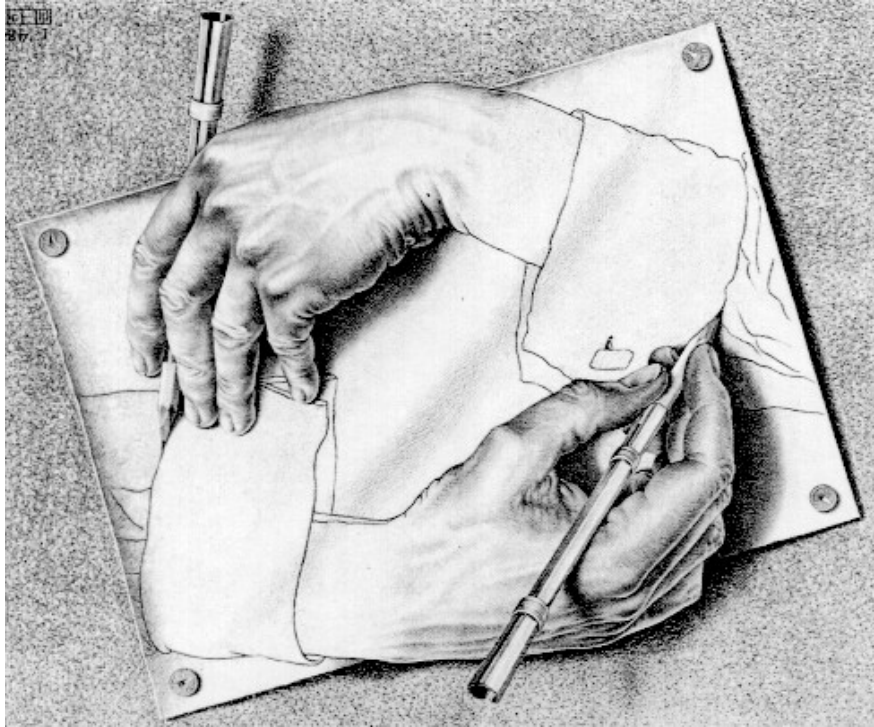
# CONCLUSIONES...



NO SABEMOS ESCRIBIR PROGRAMAS QUE COMPRENDAN EL TEXTO COMO LO HACE UN HUMANO...

PROBABLEMENTE NO NECESITAMOS (O NO PODEMOS) ESCRIBIR PROGRAMAS QUE VERDADERAMENTE COMPRENDAN EL TEXTO

NECESITAMOS ÚNICAMENTE ESCRIBIR PROGRAMAS QUE RAZONABLEMENTE PROCESEN MASAS DE DOCUMENTOS EN LUGAR DE LAS PERSONAS... Y QUE LO HAGAN BIEN Y RÁPIDAMENTE



**Merci beaucoup!**

***Avez-vous des questions?***

**`juan-manuel.torres@univ-avignon.fr`**  
**`http://juanmanuel.torres.free.fr/`**