

# Incremental Methods for Information Access in Context: The Role of Topic Descriptors and Discriminators

Carlos M. Lorenzetti – Rocío L. Cecchini  
Ana G. Maguitman



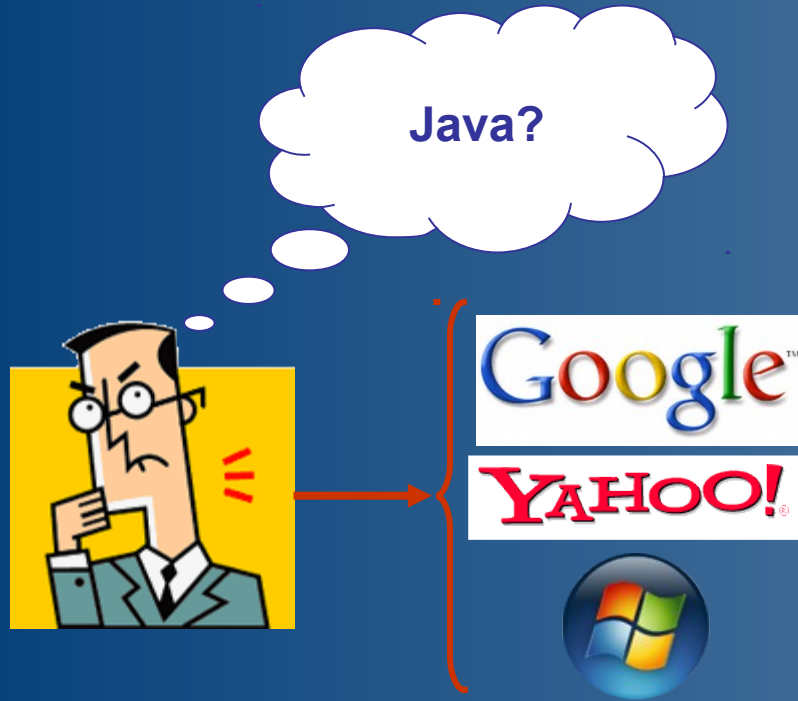
# Problemas: ambigüedad

Java?





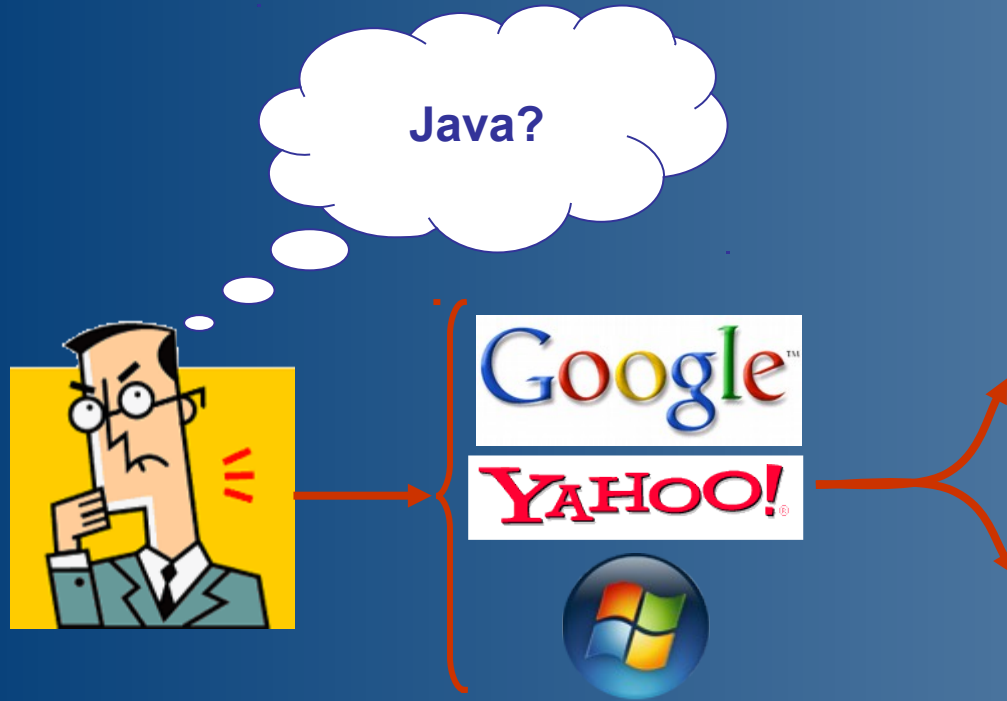
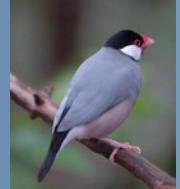
# Problemas: ambigüedad





# Problemas: ambigüedad

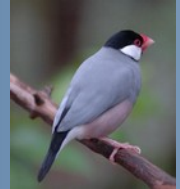
## Animales



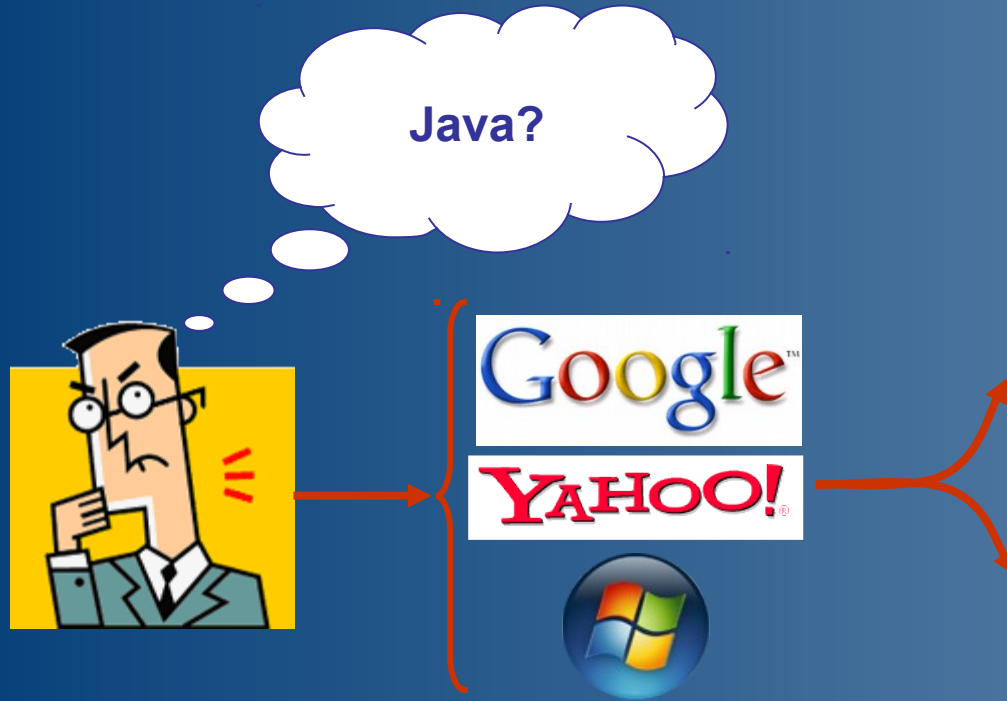


# Problemas: ambigüedad

Animales



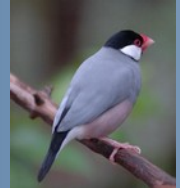
Computación





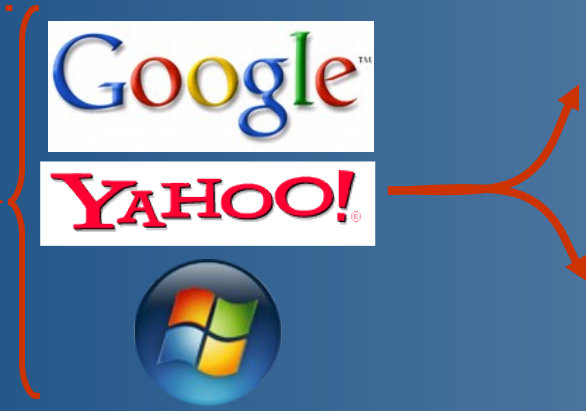
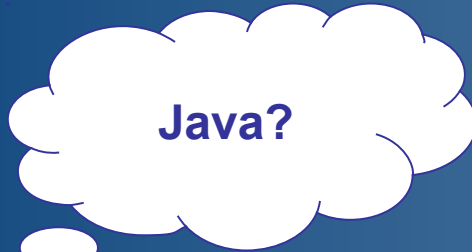
# Problemas: ambigüedad

Animales



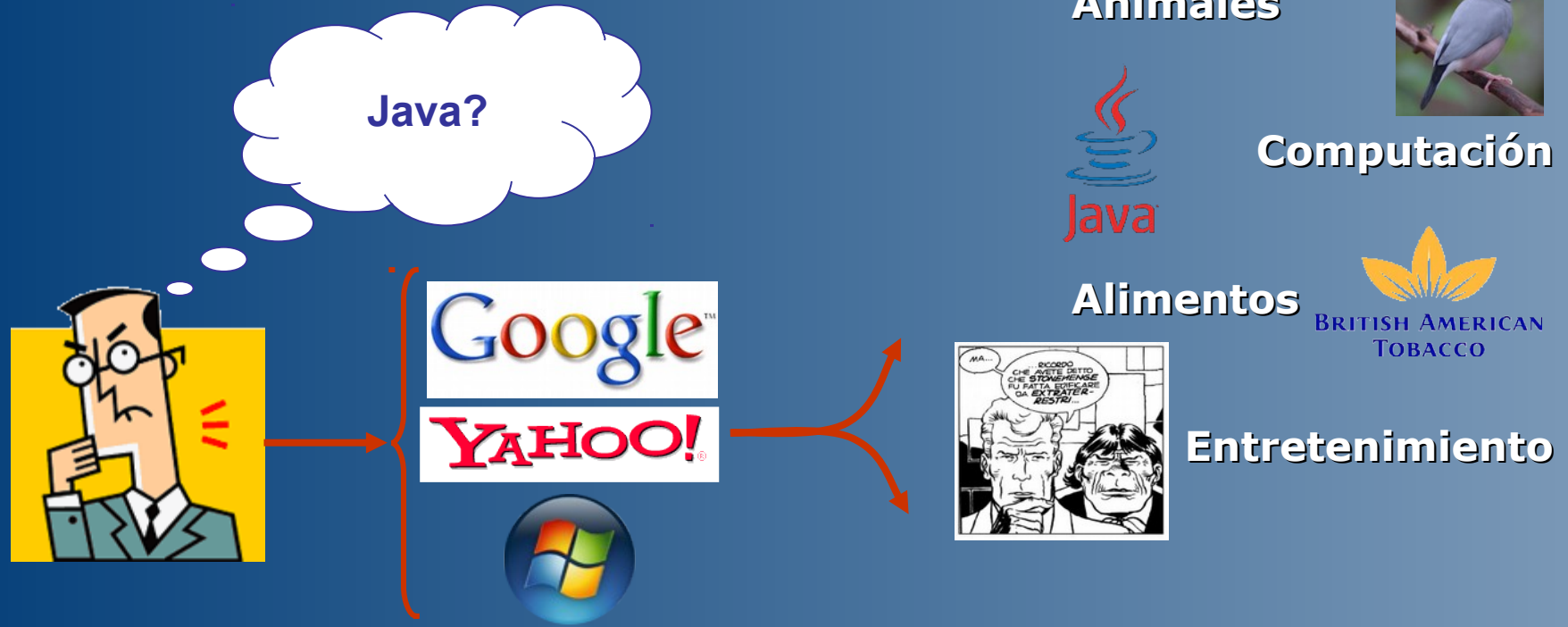
Computación

Alimentos

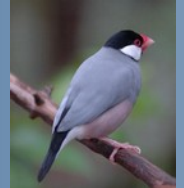




# Problemas: ambigüedad



Animales



Computación

Alimentos

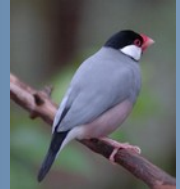


Entretenimiento



# Problemas: ambigüedad

Animales



Computación

Alimentos



BRITISH AMERICAN TOBACCO

Entretenimiento



Geografía



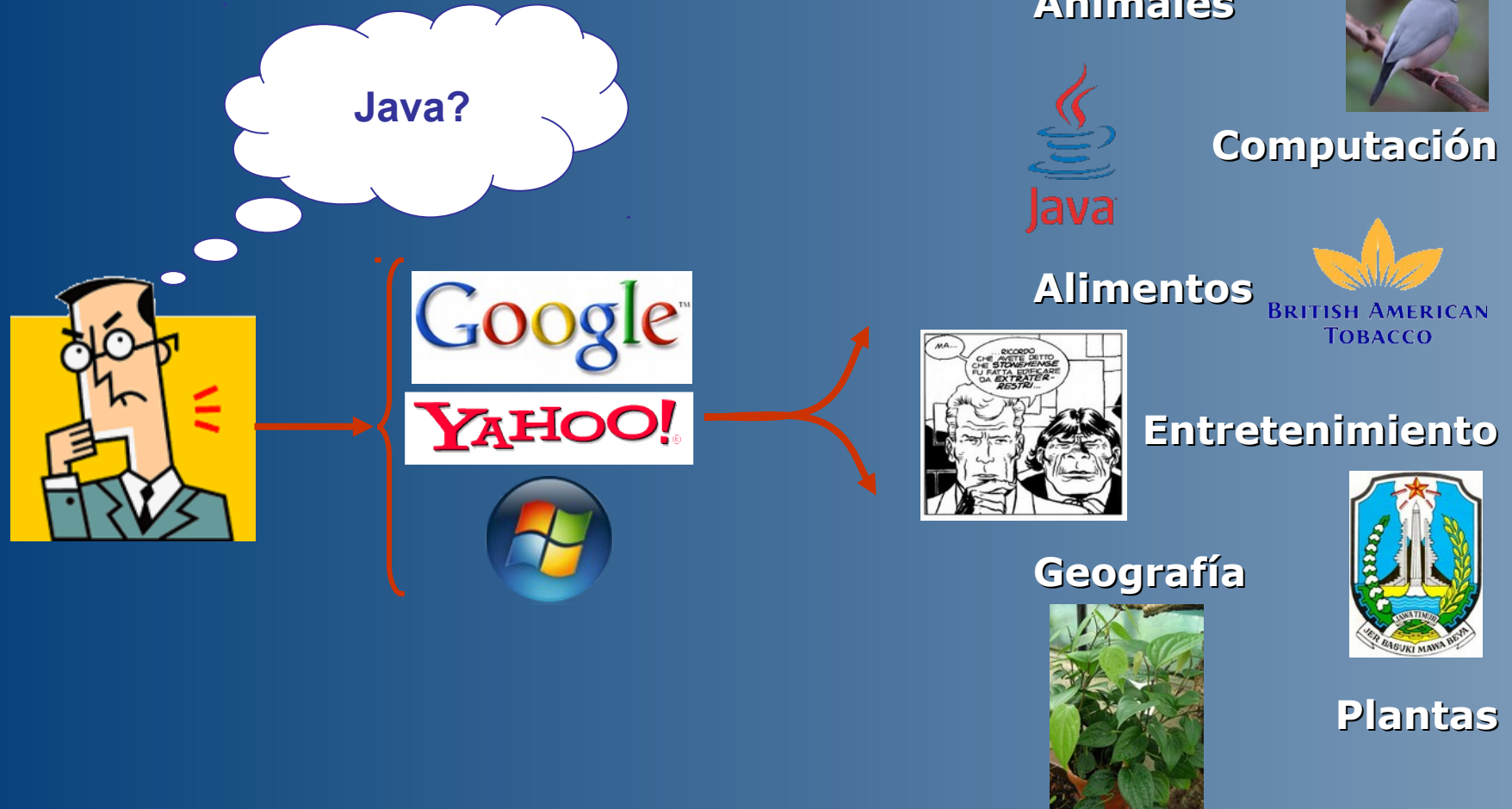
Java?







# Problemas: ambigüedad

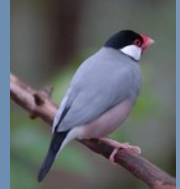




# Problemas: ambigüedad



**Animales**



**Computación**

**Alimentos**



BRITISH AMERICAN  
TOBACCO



**Entretenimiento**

**Geografía**



**Plantas**

**Barcos**





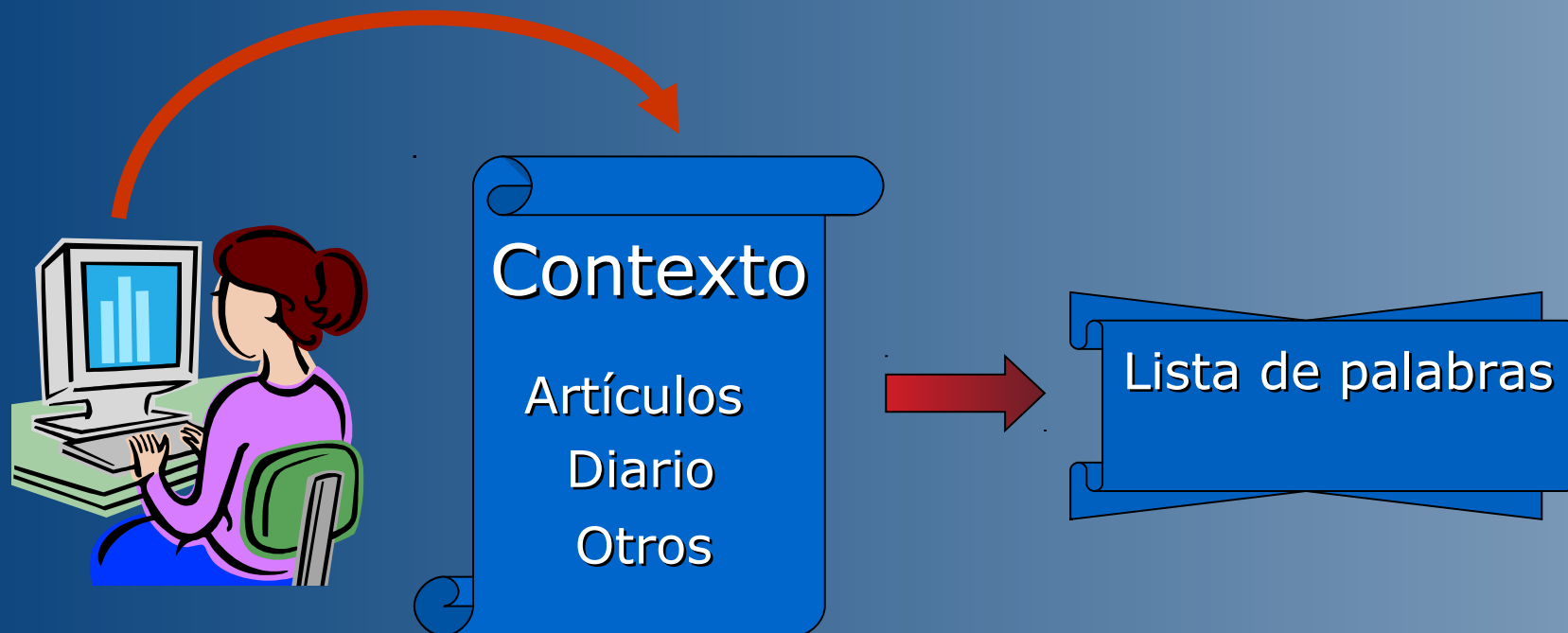
# Soluciones

## Proponemos:

- identificar términos específicos
- encontrar fuentes relevantes
- generar incrementalmente consultas

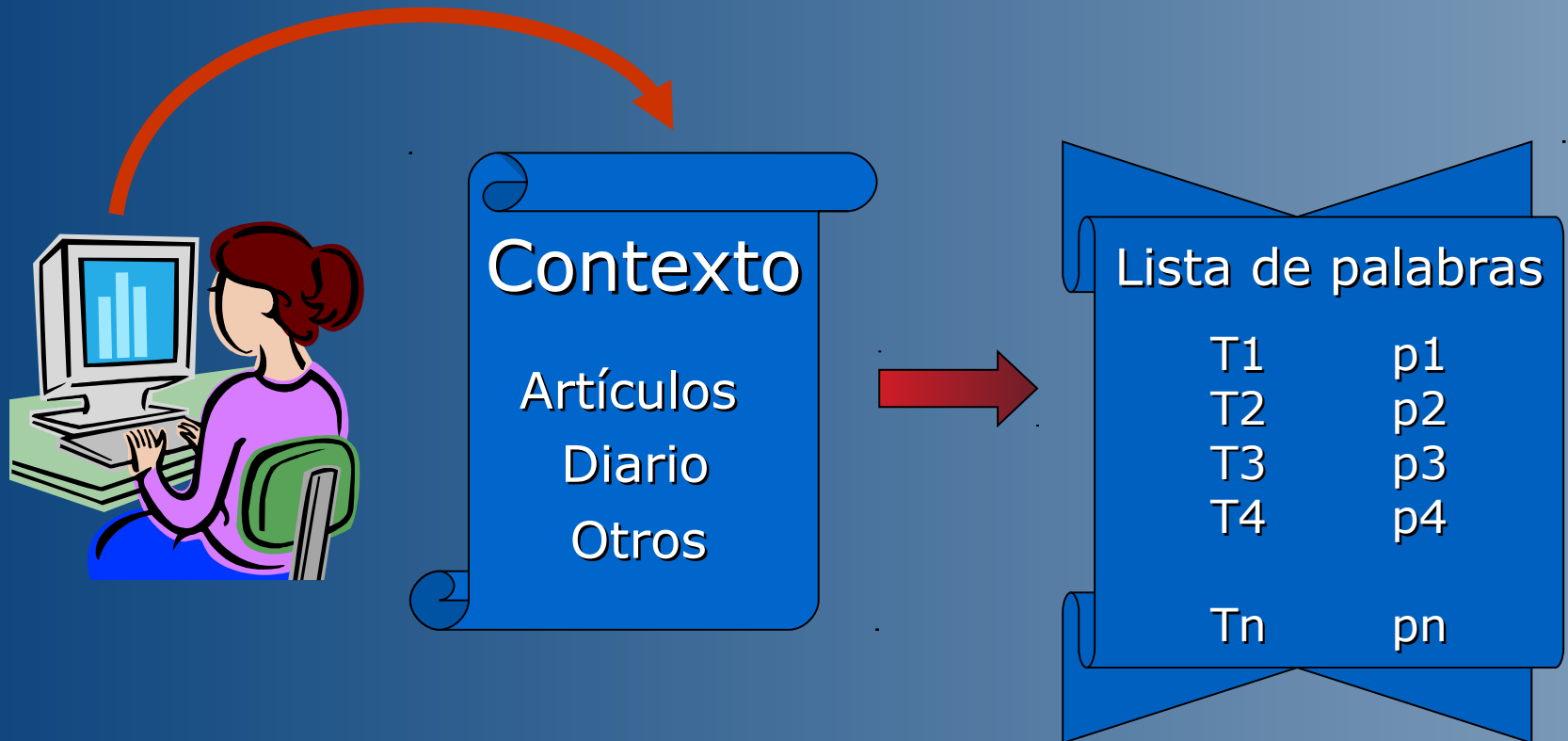


# Una solución: CONTEXTO





# Una solución: CONTEXTO





# Importancia de los términos

## Método tradicional: TF-IDF

emplea la forma más simple

$$TFIDF(d, t) = TF(d, t) \times IDF(t)$$



# Importancia de los términos

## Método tradicional: TF-IDF

emplea la forma más simple

$$TFIDF(d, t) = TF(d, t) \times IDF(t)$$

Cuenta las apariciones  
de un término en el  
documento

Penaliza a aquella  
palabras que son  
muy comunes



# Importancia de los términos

## Método Propuesto: Incremental

- *Descriptorios*

Términos que aparecen **muchas veces** en documentos de un mismo tópico:

*¿Sobre qué trata este tema?*

- *Discriminadores*

Términos que **sólo** aparecen en documentos de un mismo tópico:

*¿Qué palabras utilizo para encontrar información similar?*





# Descriptores y Discriminadores



Tópico: Máquina Virtual de Java



# Descriptores y Discriminadores



Tópico: Máquina Virtual de Java



# Descriptores y Discriminadores



Tópico: Máquina Virtual de Java



# Cálculo de Descriptores y Discriminadores



# Descriptores y Discriminadores en Documentos

Contexto Inicial		H			
		(1)	(2)	(3)	(4)
java	4	2	5	5	2
máquina	2	6	3	2	0
virtual	1	0	1	1	0
lenguaje	1	0	2	1	1
programación	3	0	2	2	0
café	0	3	0	0	3
isla	0	4	0	0	2
provincia	0	4	0	0	1
jvm	0	0	2	1	0
jdk	0	0	3	3	0

## Tópico: Máquina Virtual de Java

- (1) [espressotec.com](http://espressotec.com)
- (2) [netbeans.org](http://netbeans.org)
- (3) [sun.com](http://sun.com)
- (4) [wikitravel.org](http://wikitravel.org)

$$\mathbf{H}[d_i, t_j] = k$$

Cantidad de **ocurrencias** del término  $k$  en el documento  $i$



# Descriptores de Documentos

Contexto Inicial		$\lambda(d_0, t_j)$
java	4	0,718
máquina	2	0,359
virtual	1	0,180
lenguaje	1	0,180
programación	3	0,539
café	0	0,000
isla	0	0,000
provincia	0	0,000
jvm	0	0,000
jdk	0	0,000

Tópico: Máquina Virtual de Java

Poder **descriptivo** de un término de un **documento**

$$\lambda(d_i, t_j) = \frac{\mathbf{H}[i, j]}{\sqrt{\sum_{k=0}^{n-1} (\mathbf{H}[i, k])^2}}$$



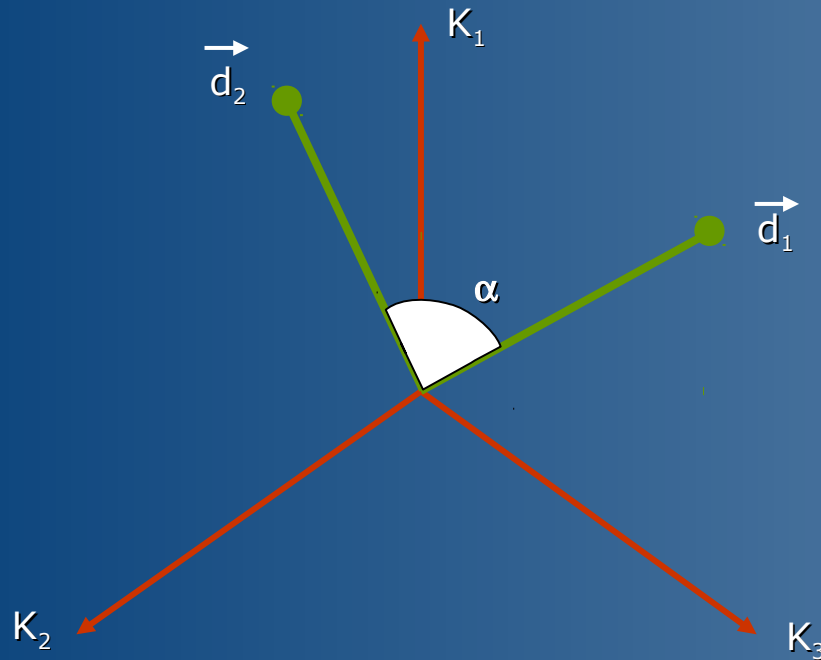
# Discriminadores de Documentos

Contexto Inicial		$\delta(t_i, d_0)$
java	4	0,447
máquina	2	0,500
virtual	1	0,577
lenguaje	1	0,500
programación	3	0,577
café	0	0,000
isla	0	0,000
provincia	0	0,000
jvm	0	0,000
jdk	0	0,000

Tópico: Máquina Virtual de Java

Poder **discriminante** de un término de un **documento**

$$\delta(t_i, d_j) = \frac{s(\mathbf{H}^T[i, j])}{\sqrt{\sum_{k=0}^{m-1} s(\mathbf{H}^T[i, k])}}$$



Criterio de comparación  
de documentos:  
*Similitud por coseno*

$$\sigma(d_i, d_j) = \sum_{k=0}^{n-1} (\lambda(d_i, t_k) \cdot \lambda(d_j, t_k))$$

Similitud entre  
documentos





# Descriptores de Tópicos

Contexto Inicial		$\Lambda(d_0, t_j)$
java	4	0,385
máquina	2	0,158
jdk	0	0,124
café	0	0,089
isla	0	0,064
programación	3	0,055
lenguaje	1	0,040
provincia	0	0,040
jvm	0	0,032
virtual	1	0,014

## Tópico: Máquina Virtual de Java

Poder **descriptivo** de un término en el tópico de un documento

$$\Lambda(d_i, t_j) = \frac{\sum_{k=0, k \neq i}^{m-1} (\sigma(d_i, d_k) \cdot \lambda(d_k, t_j)^2)}{\sum_{k=0, k \neq i}^{m-1} \sigma(d_i, d_k)}$$



# Discriminadores de Tópicos

Contexto Inicial		$\Delta(t_i, d_0)$
jvm	0	0,848
jdk	0	0,848
virtual	1	0,566
programación	3	0,566
máquina	2	0,524
lenguaje	1	0,517
java	4	0,493
café	0	0,385
isla	0	0,385
provincia	0	0,385

Tópico: Máquina Virtual de Java

Poder **discriminante** de un término en el tópico de un documento

$$\Delta(t_i, d_j) = \sum_{k=0, k \neq j}^{m-1} (\sigma(d_k, d_j) \cdot \delta(t_i, d_k)^2)$$

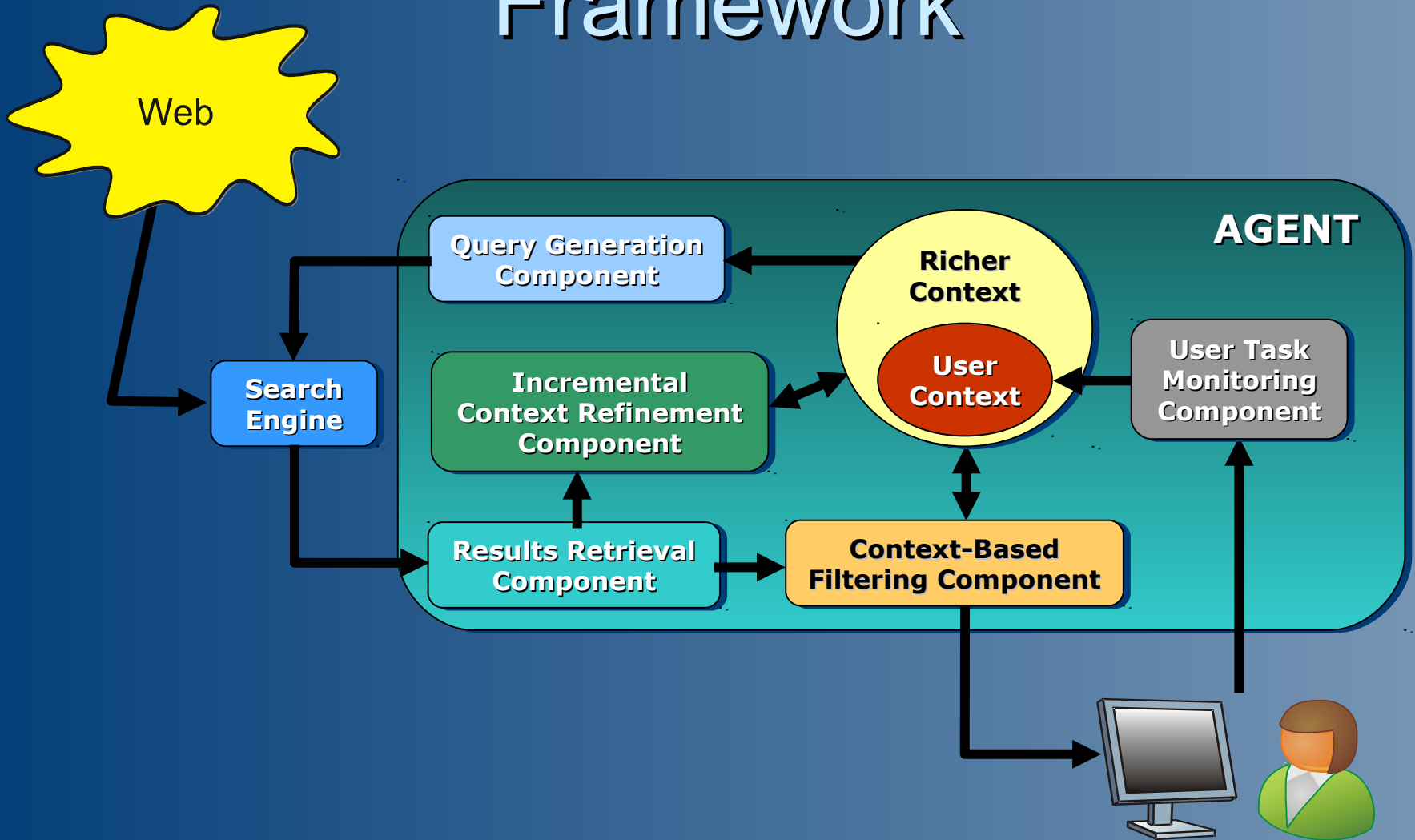


# Implementación



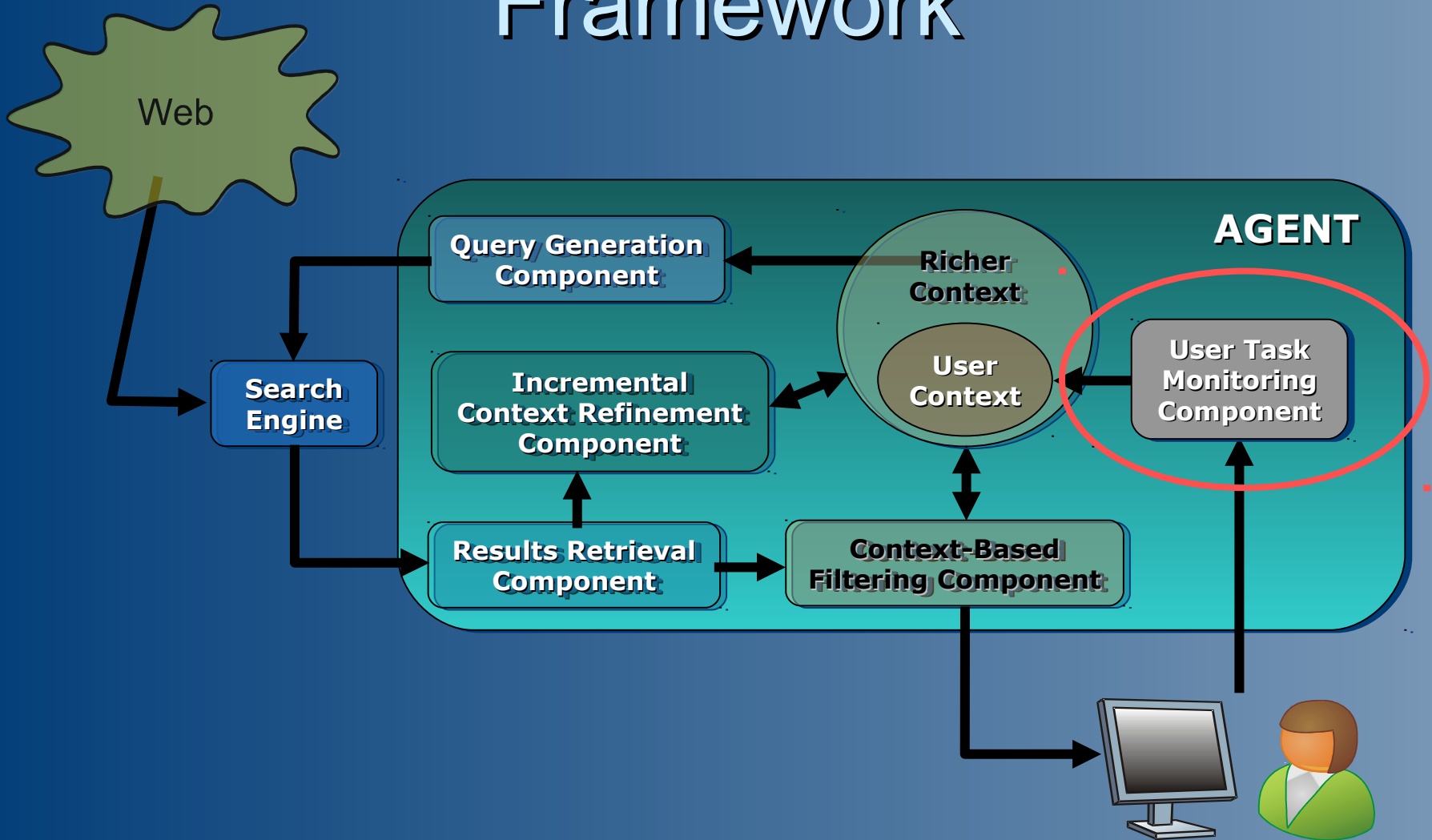


# Framework



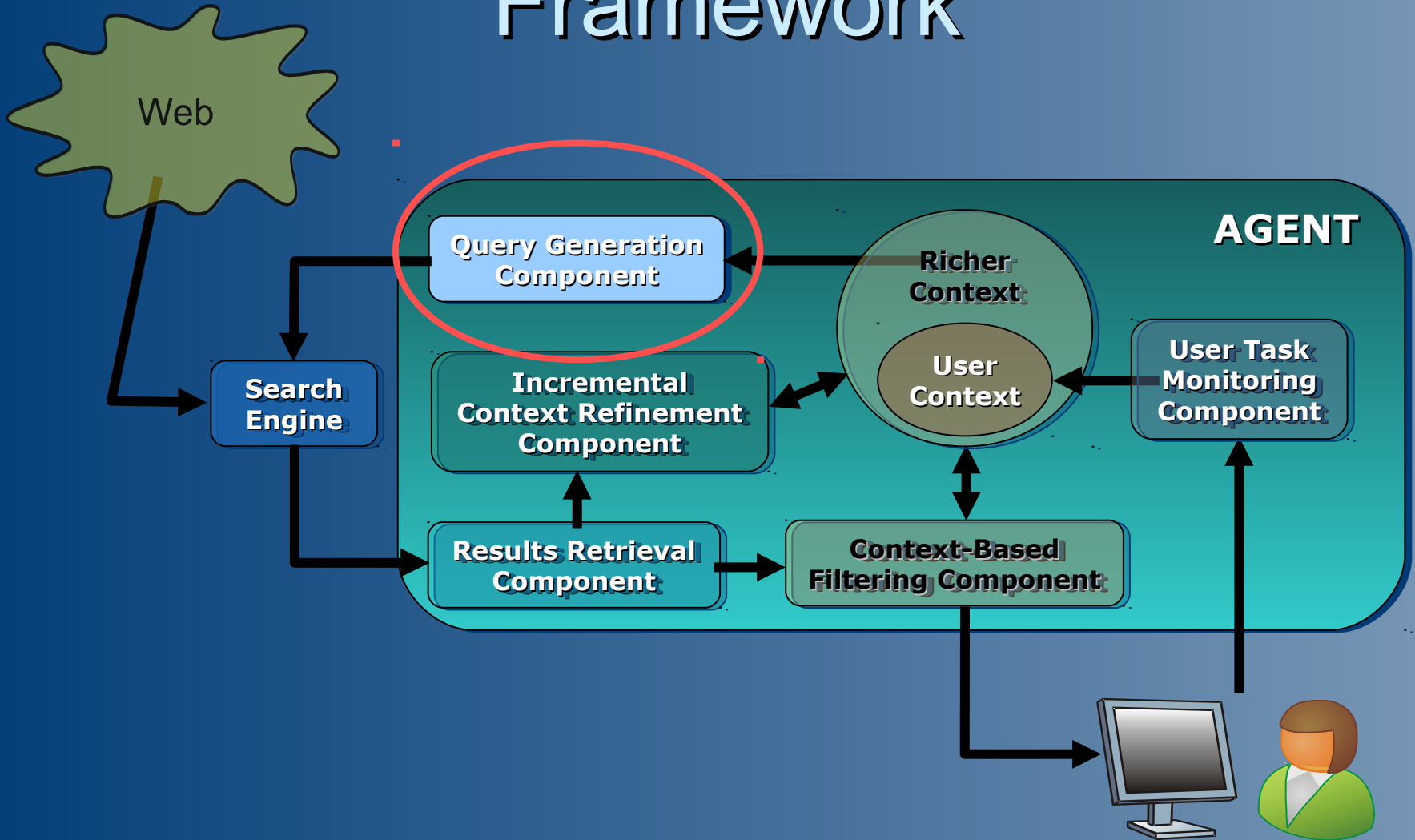


# Framework



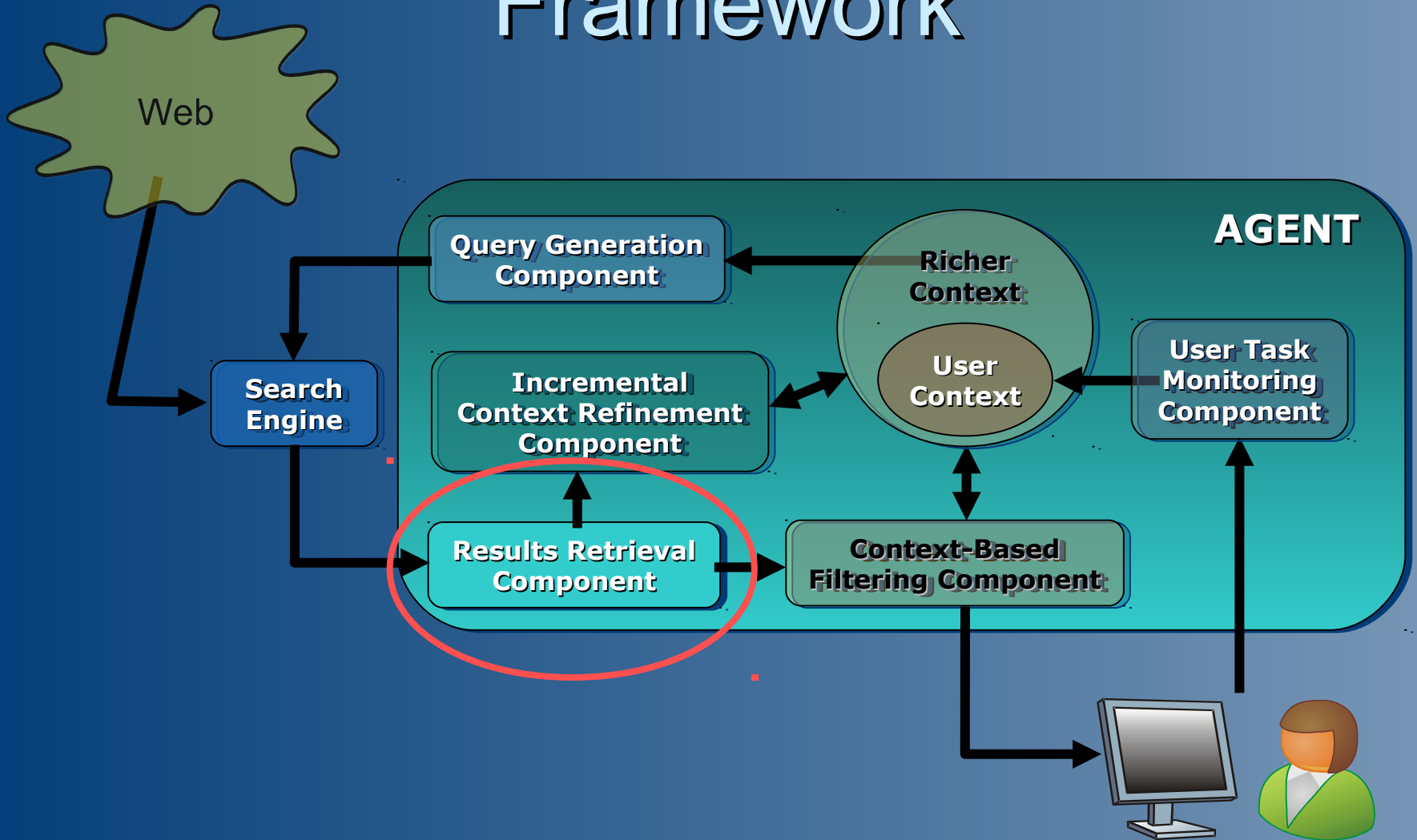


# Framework



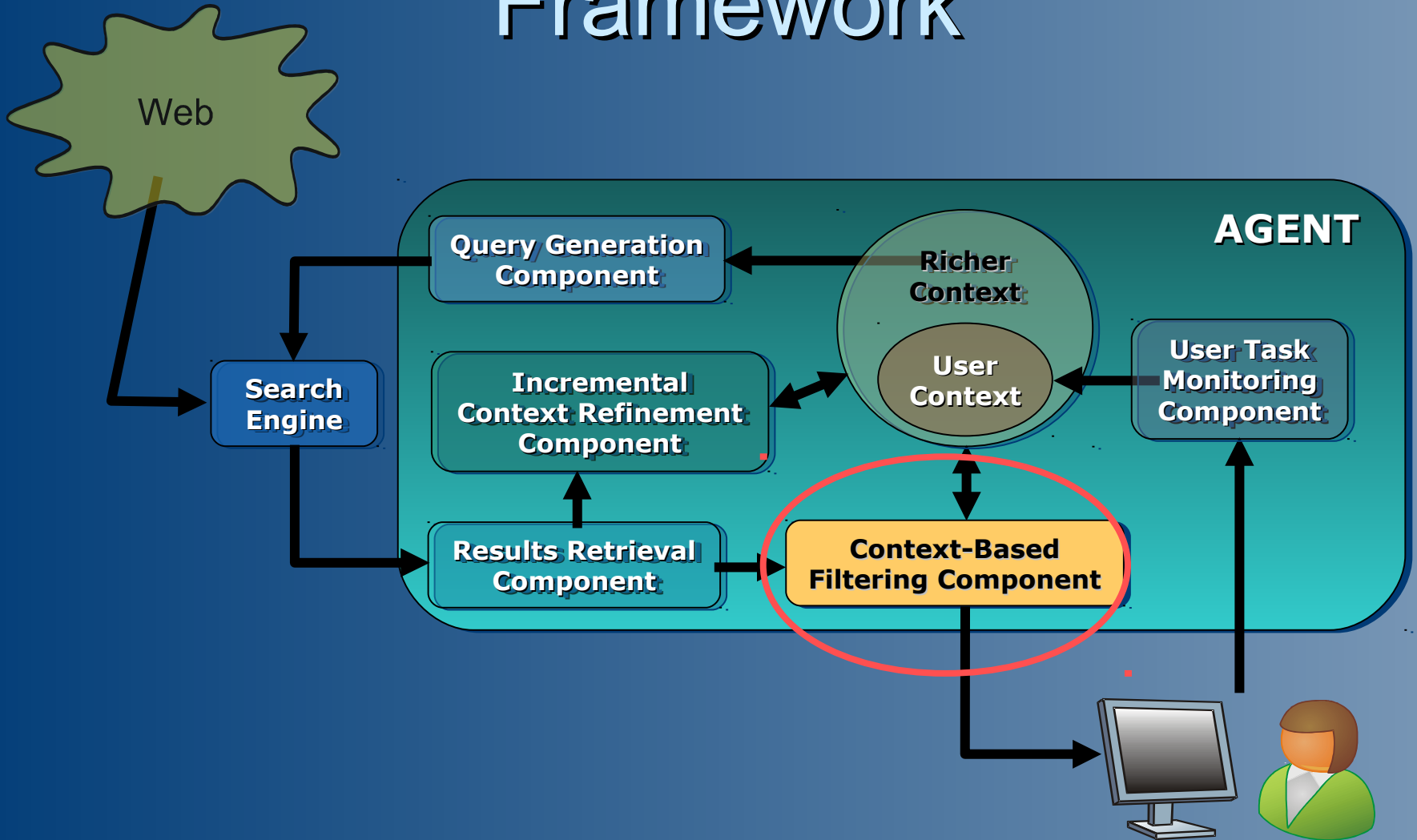


# Framework





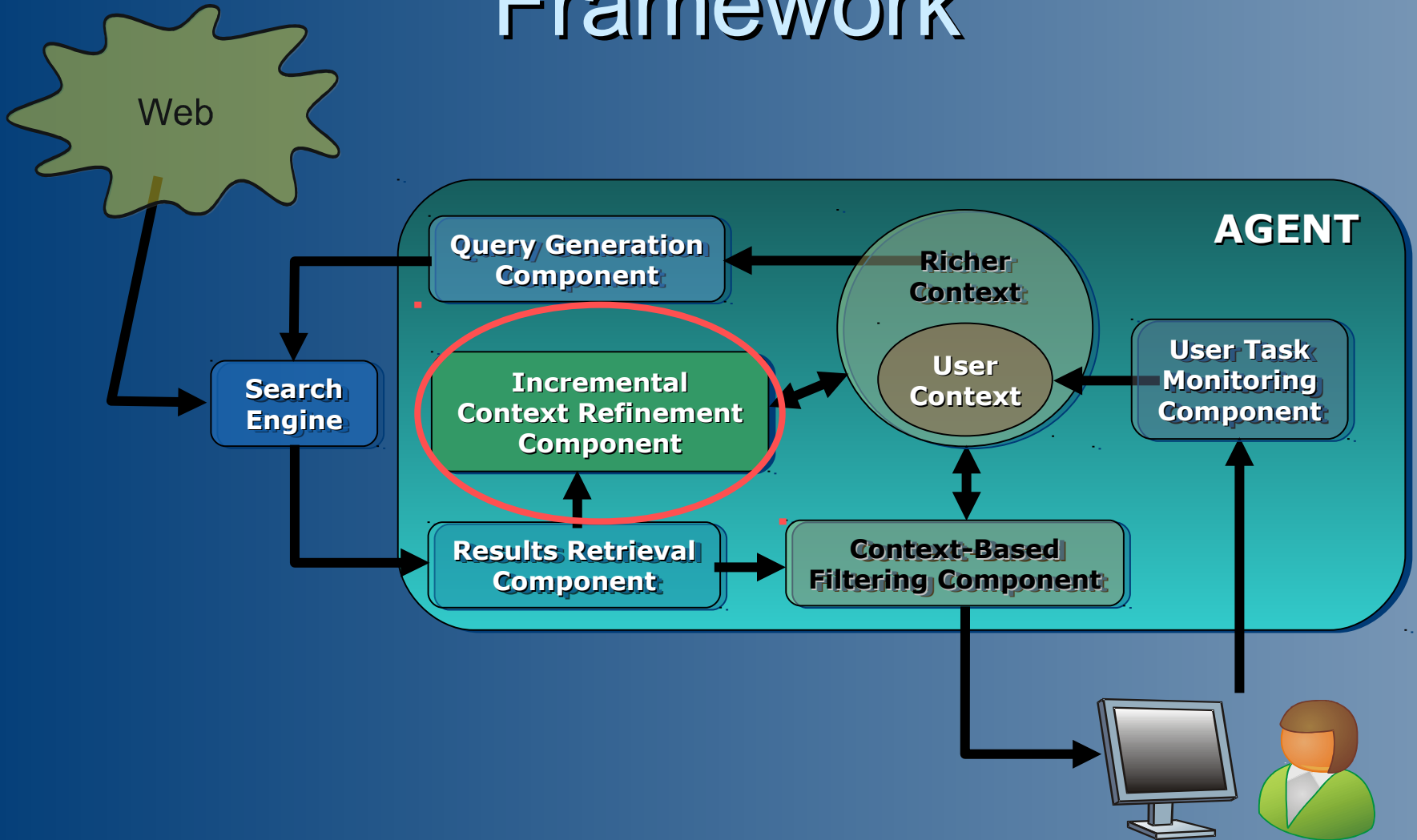
# Framework







# Framework





# Evaluación



# Evaluación

## Intelligent Incremental Method

1. Generar  $k$  consultas usando los términos del contexto
2. Enviar las  $Q(i)$  al motor de búsqueda
3. Obtener las respuestas y convertirlas a repres. vectorial
4. Generar una lista ordenada de **descriptores**,  $L'_\Delta$
5. Generar una lista ordenada de **discriminadores**,  $L'_\Delta$
6. Actualizar incrementalmente las listas  $L_\Delta$  y  $L'_\Delta$
7. Generar  $k$  consultas,  $Q(i) \leftarrow$  **una combinación de  $L_\Delta$  y  $L'_\Delta$**
8. Enviar las  $Q(i)$  al motor de búsqueda
9. Obtener las respuestas y convertirlas a repres. vectorial
10. Para cada respuesta, comprobar si es una **buena consulta**
11. Para cada **mala consulta**, tratar de reformularla
12. ir al paso 4



# Evaluación

## Naïve Method

1. Generar  $k$  consultas usando los términos del contexto
2. Enviar las  $Q(i)$  al motor de búsqueda
3. Obtener las respuestas y convertirlas a repres. Vectorial
4. Generar una lista ordenada términos  $L_{TF}$  por frecuencia
5. Generar  $k$  consultas,  $Q(i) \leftarrow$  una combinación de  $L_{TF}$
6. ir al paso 2



# Evaluación

## Contexto Inicial

- 15 páginas en inglés del DMOZ
- Tópicos: **Recreación, Negocios, Sociedad**

## Consulta

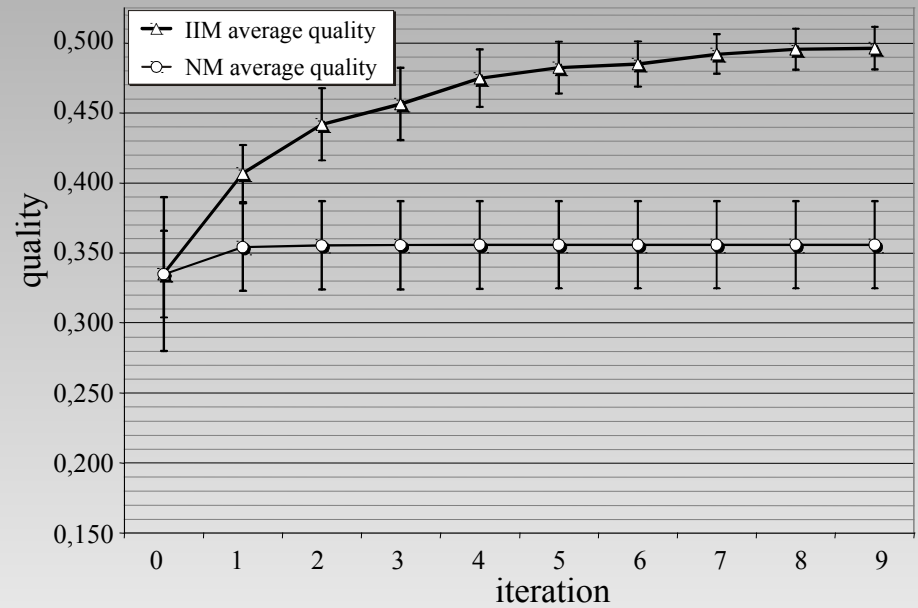
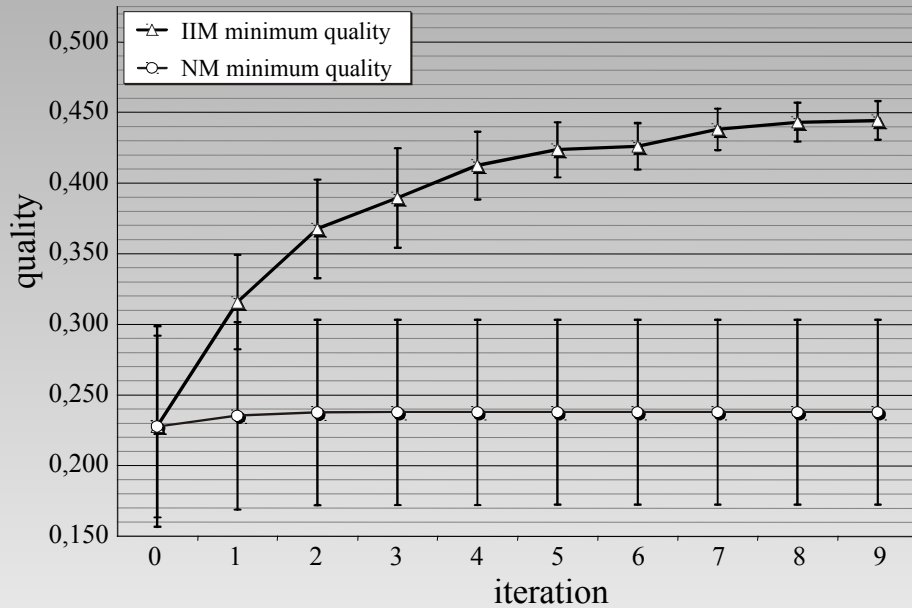
- **5** términos  $L_{\Delta}$  + **1** término  $L_{\Delta}$
- 20 consultas simultáneas
- $\forall \alpha = 0.4$
- Motor de búsqueda: Google

## Resultados analizados

- Similitudes promedio mínima y promedio por iteración



# Evaluación





# Trabajo a Futuro

- Evaluaciones intensivas con colecciones estándar (TREC, etc)
- Mejorar la pérdida del foco
- Evaluar métodos cualitativos para reordenar los resultados basándose en preferencias
- Pruebas con usuarios

# Incremental Methods for Information Access in Context: The Role of Topic Descriptors and Discriminators

Carlos M. Lorenzetti – Rocío L. Cecchini  
Ana G. Maguitman