

# Webmining



una breve introducción y la  
perspectiva Java

# Próximamente charlas

- Análisis y auditoría de performance  (A. Popovski) 22/09/05 (19:00hs)  
en internet
- Filtro semántico e IR  (N. Di Tada) 06/10/05 (19:00hs)
- NN con Java  (D. López De Luise) 20/10/05 (19:00hs)
- Java y el XML  (D. López De Luise) 27/10/05 (19:00hs)
- La IA en Java  (D. López De Luise) 17/11/05 (19:00hs)
- Cognitive Memory  (S. Piedrahita) 21/11/05 (19:00hs)
- Webbrowsing con Java  (D. López De Luise) 24/11/05 (19:00hs)

informes:

[upgrade@palermo.edu](mailto:upgrade@palermo.edu)

[sec.argentina@ieee.org](mailto:sec.argentina@ieee.org)



# Objetivo

Presentar los conceptos fundamentales de Webmining y su aplicación desde la perspectiva del lenguaje Java.

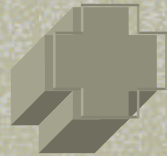
# Temario

- Introducción al webmining.
- Problemas a resolver del WWW.
- WM como una solución posible.
- Incumbencias y tipos de WM.
- Diferencias con IR, IE y ML.
- Relación con Agentes.
- Por qué WM y Java.
- Java en WM: VVV y C-BIRD.
- Java y metadatos.

# El Webmining

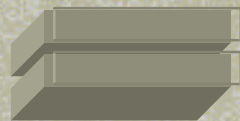
**DM**

minería de datos sobre BD



**WWW**

interconexión de BD



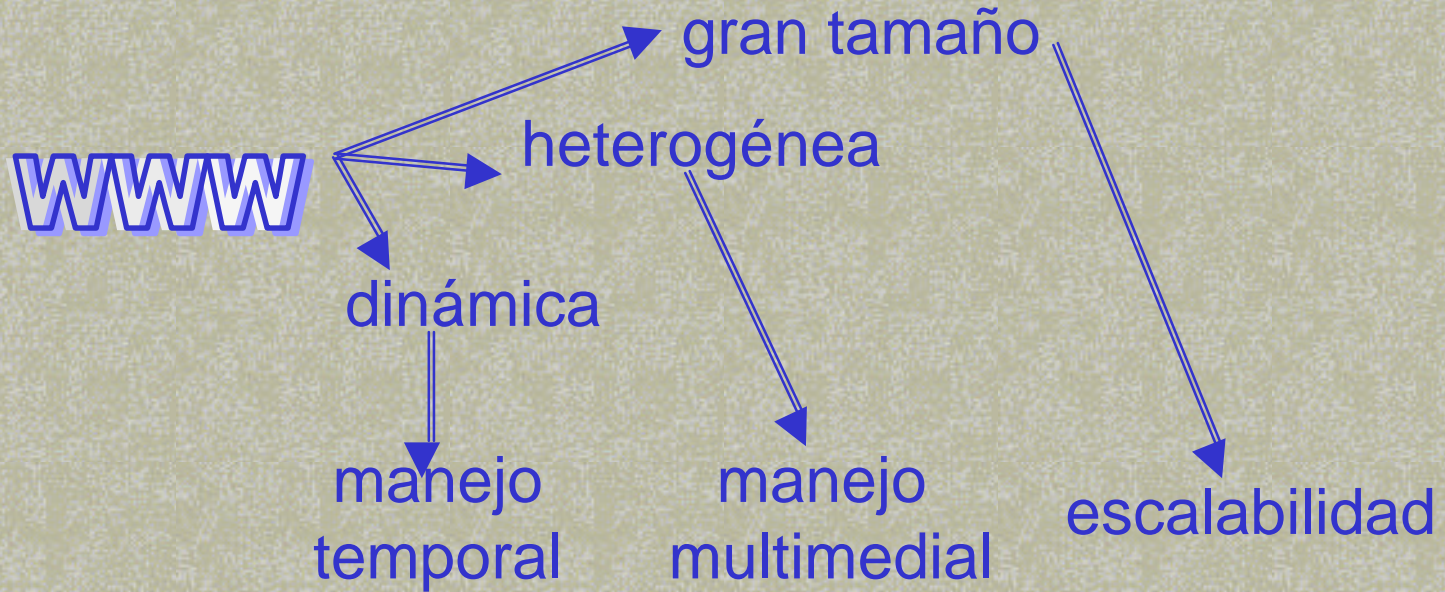
**WWM**

minería de datos sobre www



# WWW

## sus problemas



### problemas

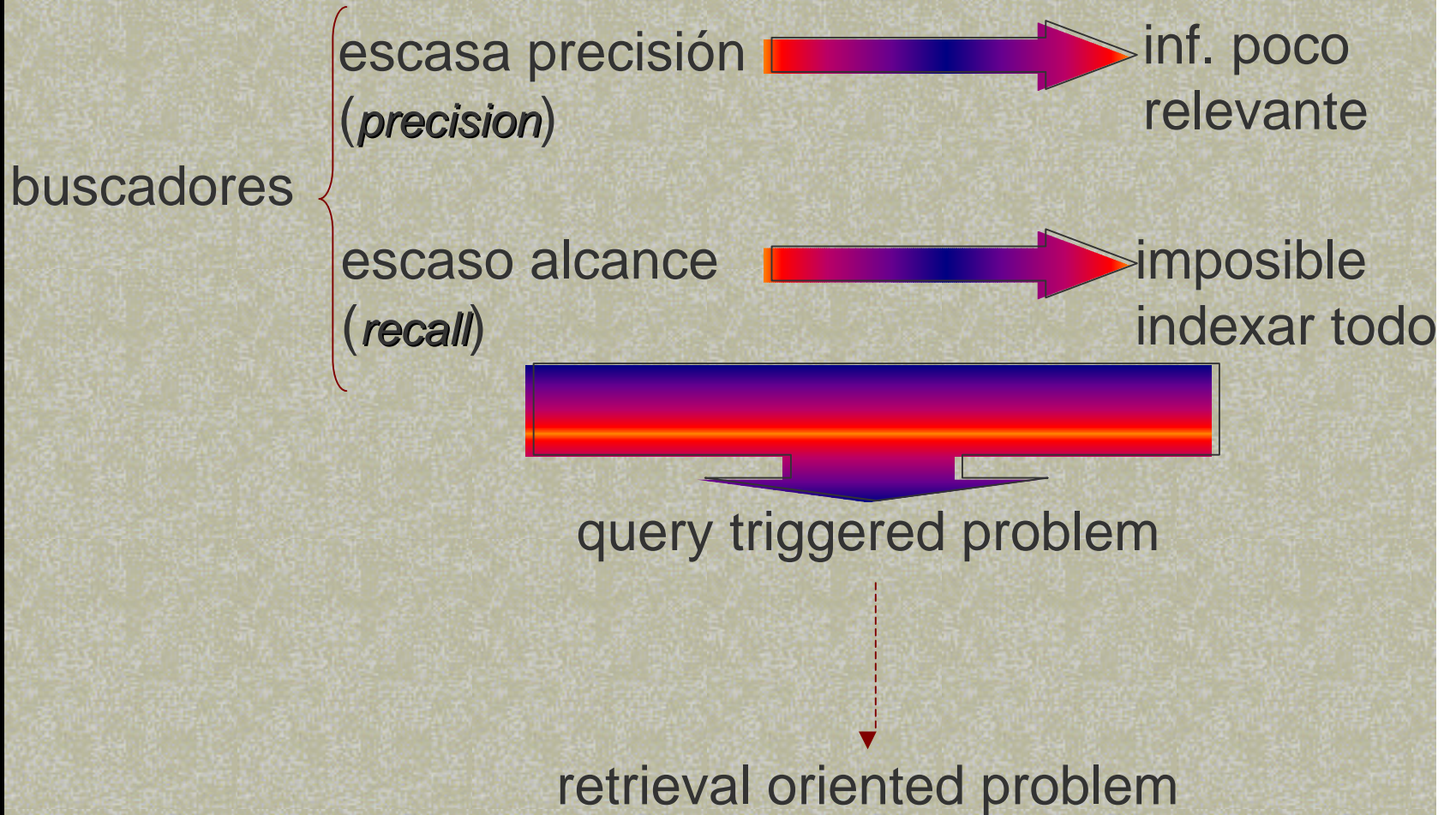
1 Hallar información relevante

2 Aprender acerca de consumidores o usuarios individuales

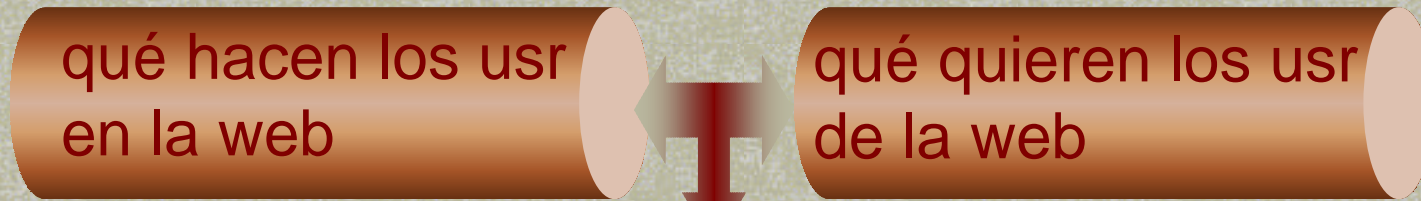
3 Personalización de la información

4 Crear información a partir de la almacenada

1 Hallar información relevante



2 Aprender acerca de consumidores o usuarios individuales



determinar qué inf. debe adaptarse a un grupo

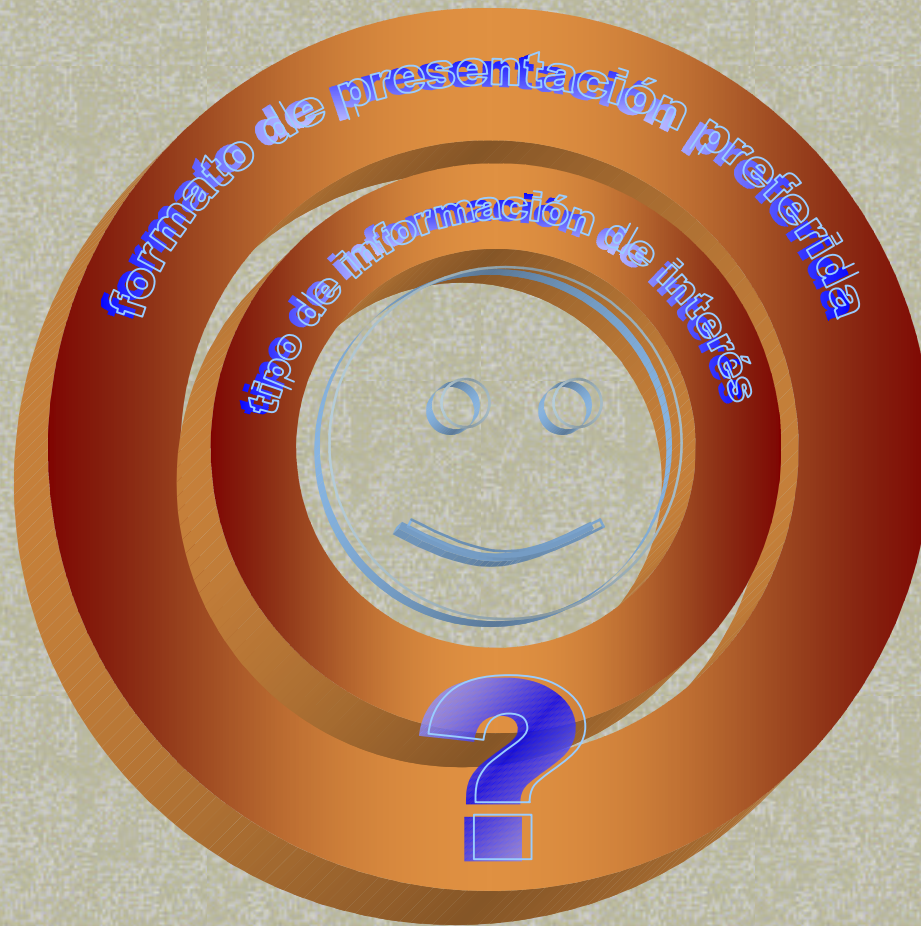
determinar qué inf. debe adaptarse a un usuario

diseño eficiente de los websites

administración eficiente de los websites

uso de la inf. del usuario para marketing

### 3 Personalización de la información



4 Crear información a partir de la almacenada

dada una colección de datos en la web



el problema es sacar información

relevante no obvia

a partir de ella

**data triggered problem**

*toma de decisiones*

data mining oriented problem



# El Webmining:

- **como una de las soluciones**
- **incumbencias**
- **tipos**
- **diferencias con IR, IE, ML**
- **relación con Agentes**



**El Webmining:**

***como una solución***



# El Webmining: como una de las soluciones





El Webmining:

*incumbencias*

# El Webmining: incumbencias

técnicas que permiten descubrir inf. desconocida no evidente relevante

## definición

aplicación de técnicas de DM

para

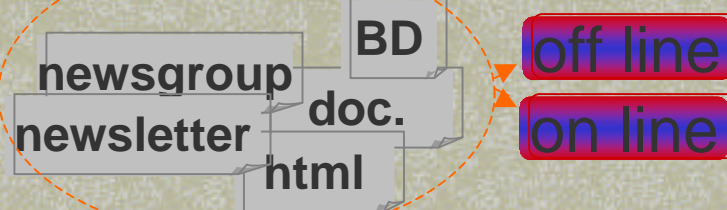
descubrir inf.

extraer inf.

desde

docs. Web

servs. Web



## tareas (Etzioni)

detección de los recursos potencialmente útiles

selección automática y preprocesado de esos recursos

detección automática de patrones (en el mismo site o entre sites)

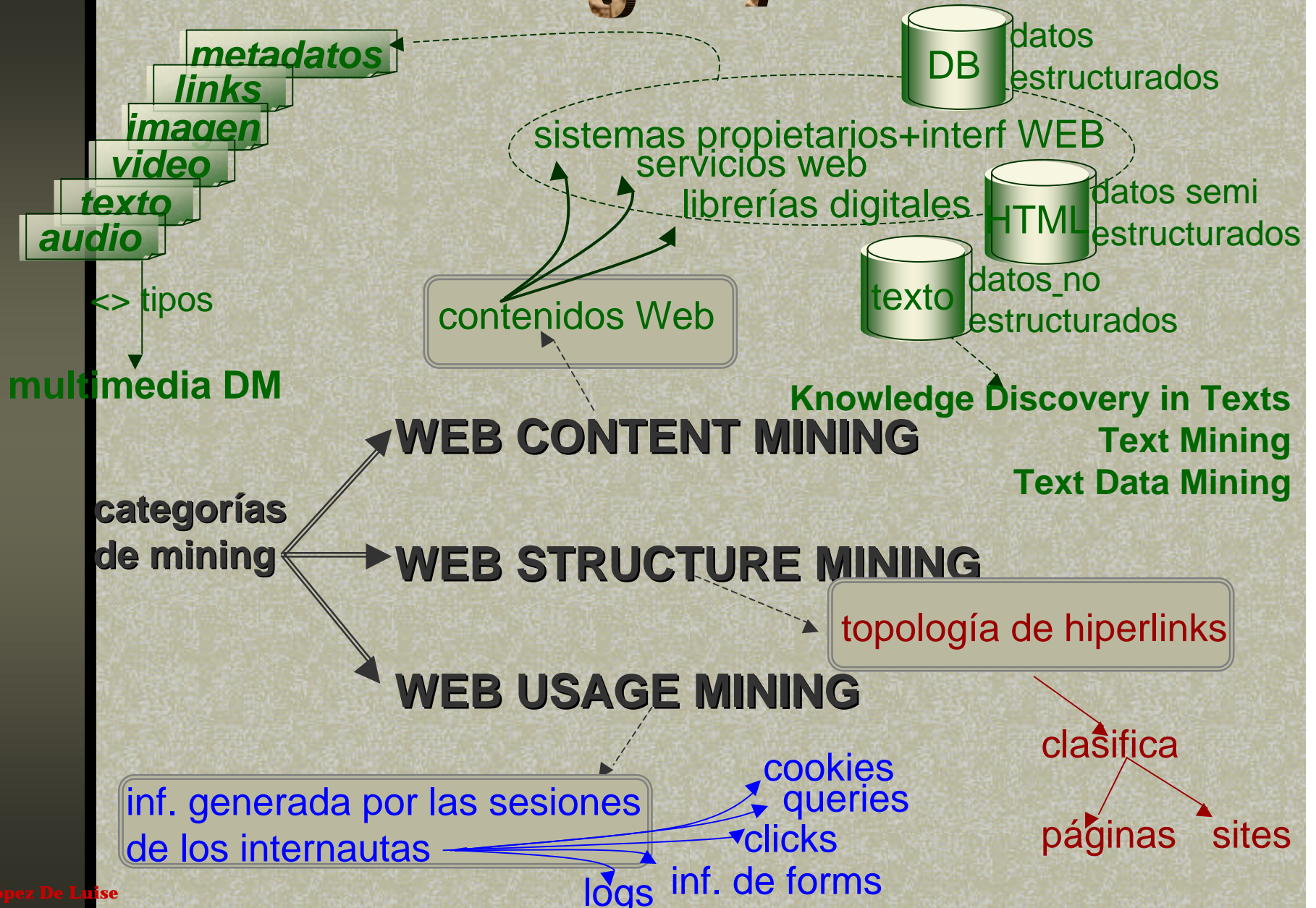
validación e interpretación de los patrones



El Webmining:

*tipos*

# El Webmining: tipos





# El Webmining:

diferencias con IR IE ML

# El Webmining: relación con IR

*WM*



*Information Retrieval*

extracción de inf.  
no evidente + relevante  
a partir de datos

recuperación automática de inf.  
procurando  
máx. cantidad docs relevantes  
+  
min. cantidad docs no relevantes

*tareas*

*tareas*

clasificación  
categorización  
filtrado de información  
formalización de patrones  
...

clasificación  
categorización  
filtrado de información  
indexación  
interfaces usr

# El Webmining: relación con IE

*WM*



*Information Extraction*

extracción de inf.  
no evidente + relevante  
a partir de datos

transforma un conj de datos  
(a veces con IR)  
simplificando su interpretación

*tareas*

detecta patrones relevantes  
procesos automáticos  
escalable

detecta hechos relevantes  
procesos automáticos o no  
a veces no escalable

# El Webmining: relación con ML.

*WM*



*Machine Learning*

extracción de inf.  
no evidente + relevante  
a partir de datos

usa aprendizaje automático  
con distintos fines  
(ej. optimización, filtrados, etc.)

*tareas*

optimización del proceso  
de clasificación en mining  
de textos



# El Webmining:

*relación con Agentes*

# El Webmining: relación con Agentes

WM



Agentes SW

extracción de inf.  
no evidente + relevante  
a partir de datos

*puede hacerla con*

pueden usar DM para  
alcanzar su objetivo

tareas

tipos

móviles

distribuidos

interfaces

content mining

structure mining

usage mining

combinación

content-based filters

reputation-based filters

collaborative/social based f.

event – based filters

hybrid filters



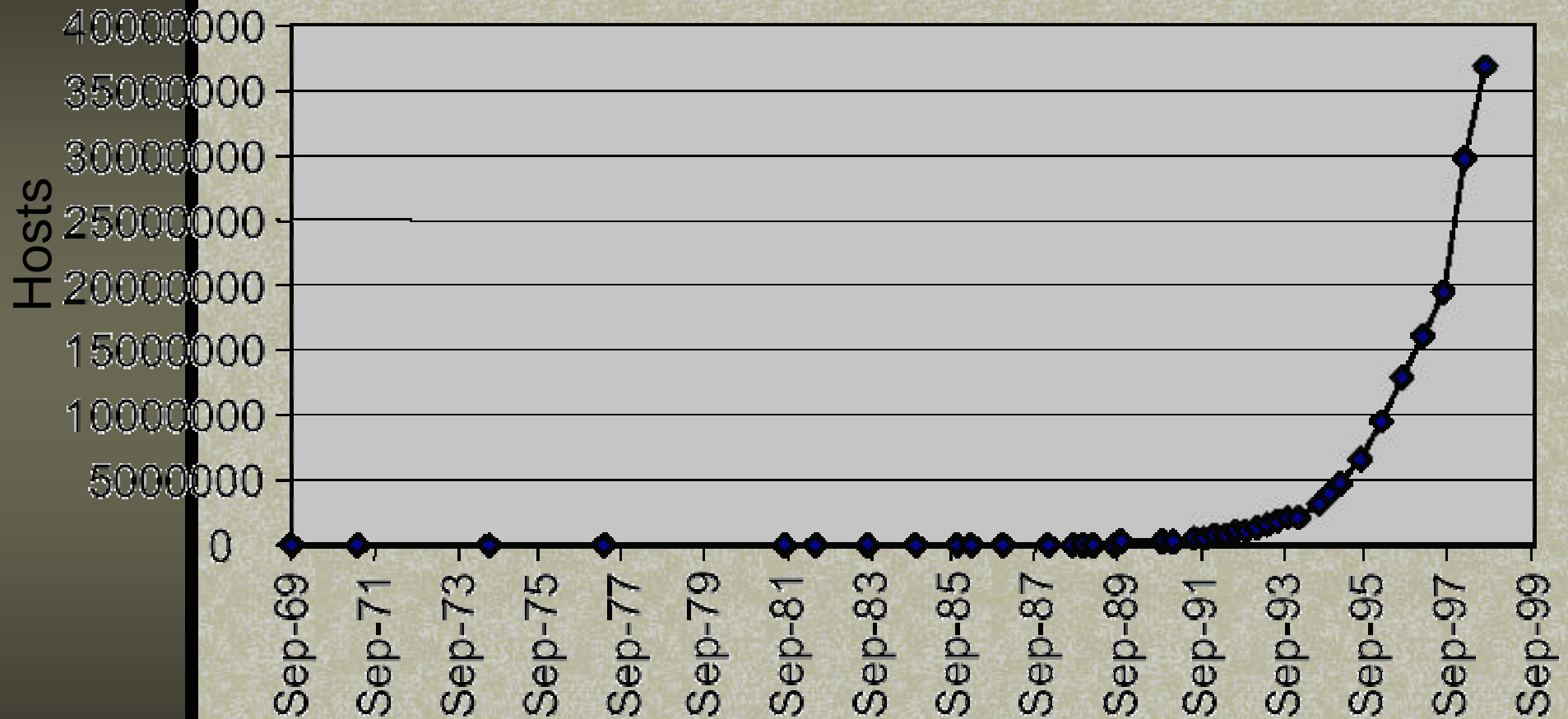
- ***Por qué hacemos WMI?***

- ***Por qué usar Java?***



# Por qué WMI?

## a. Crecimiento físico de Internet

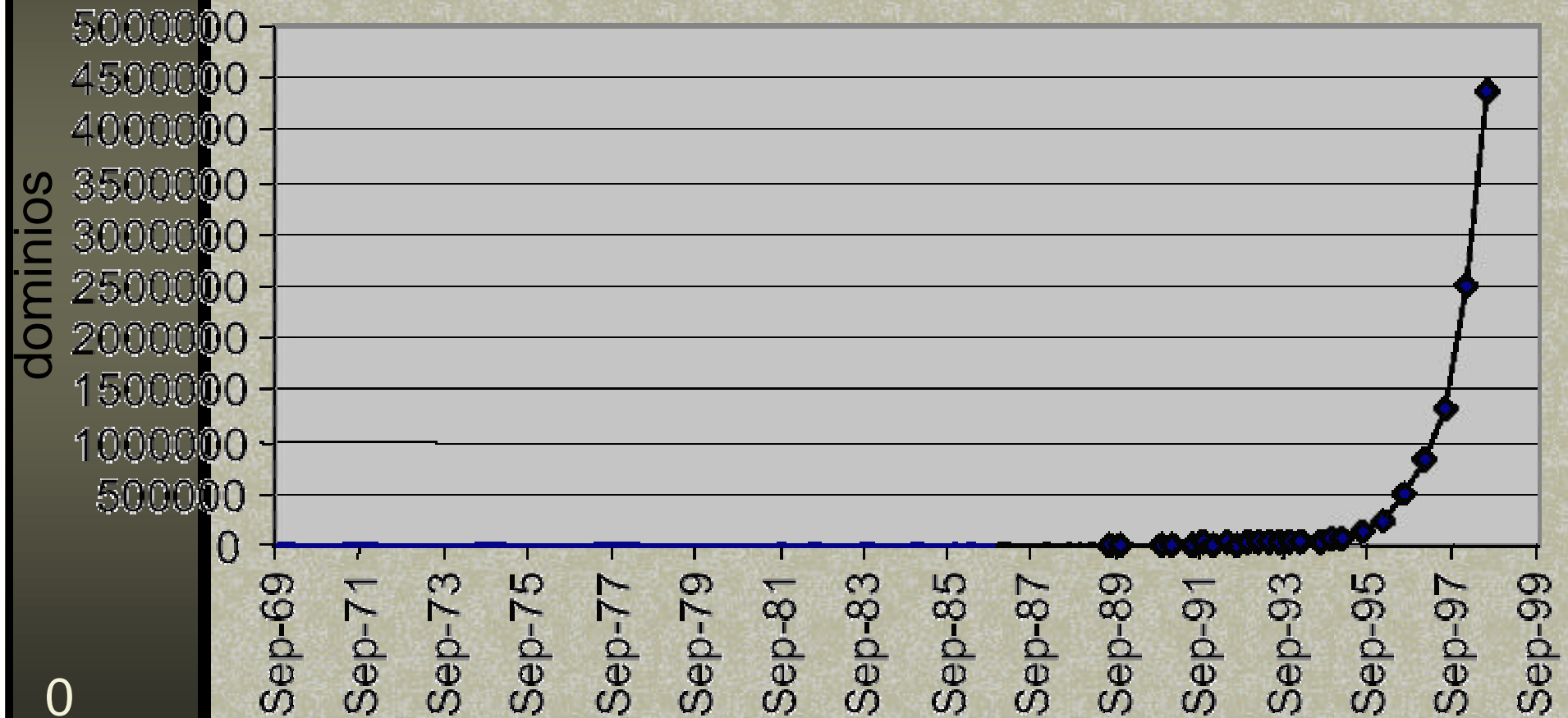


0



# Por qué WMI?

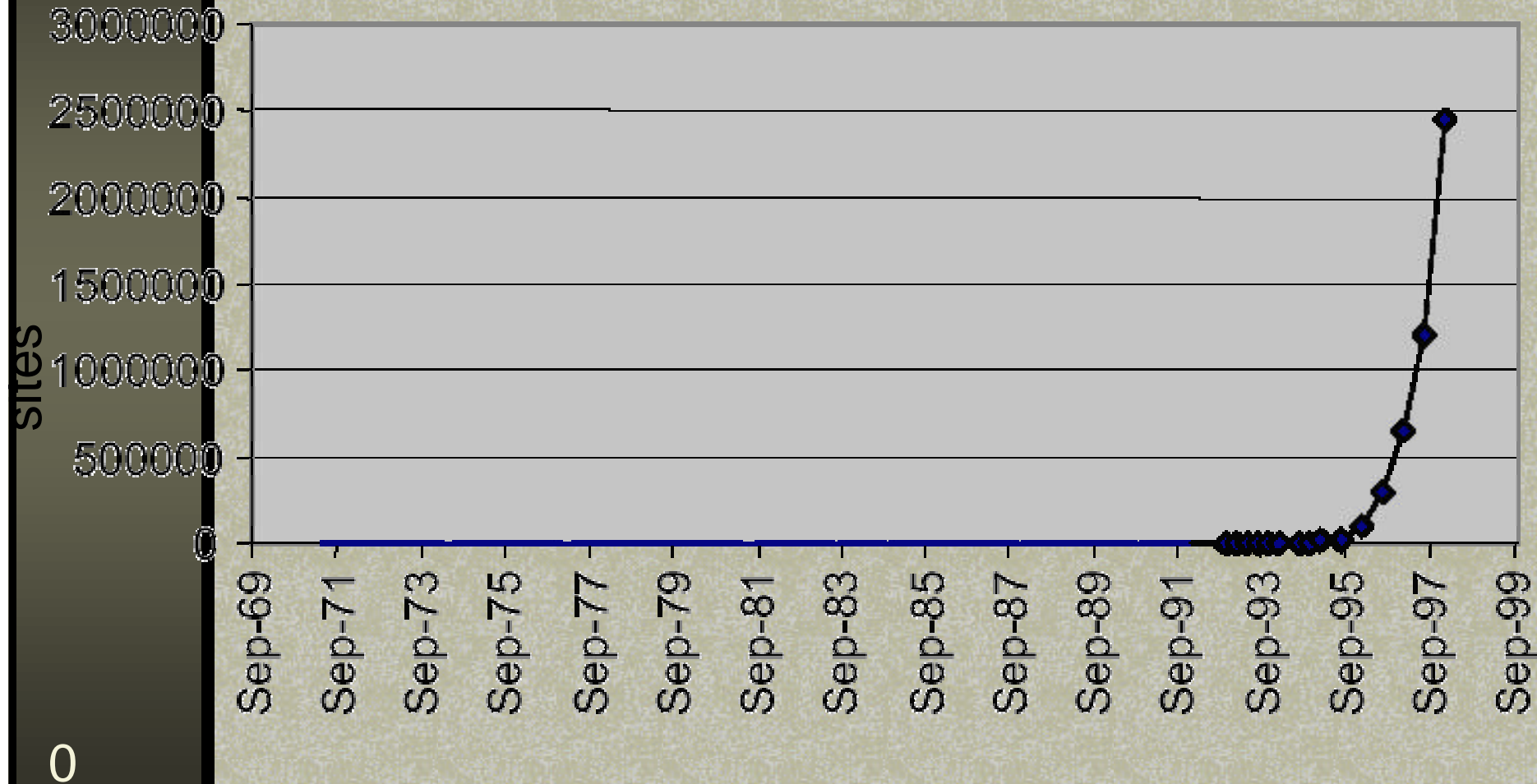
## b. Crecimiento de cantidad de dominios





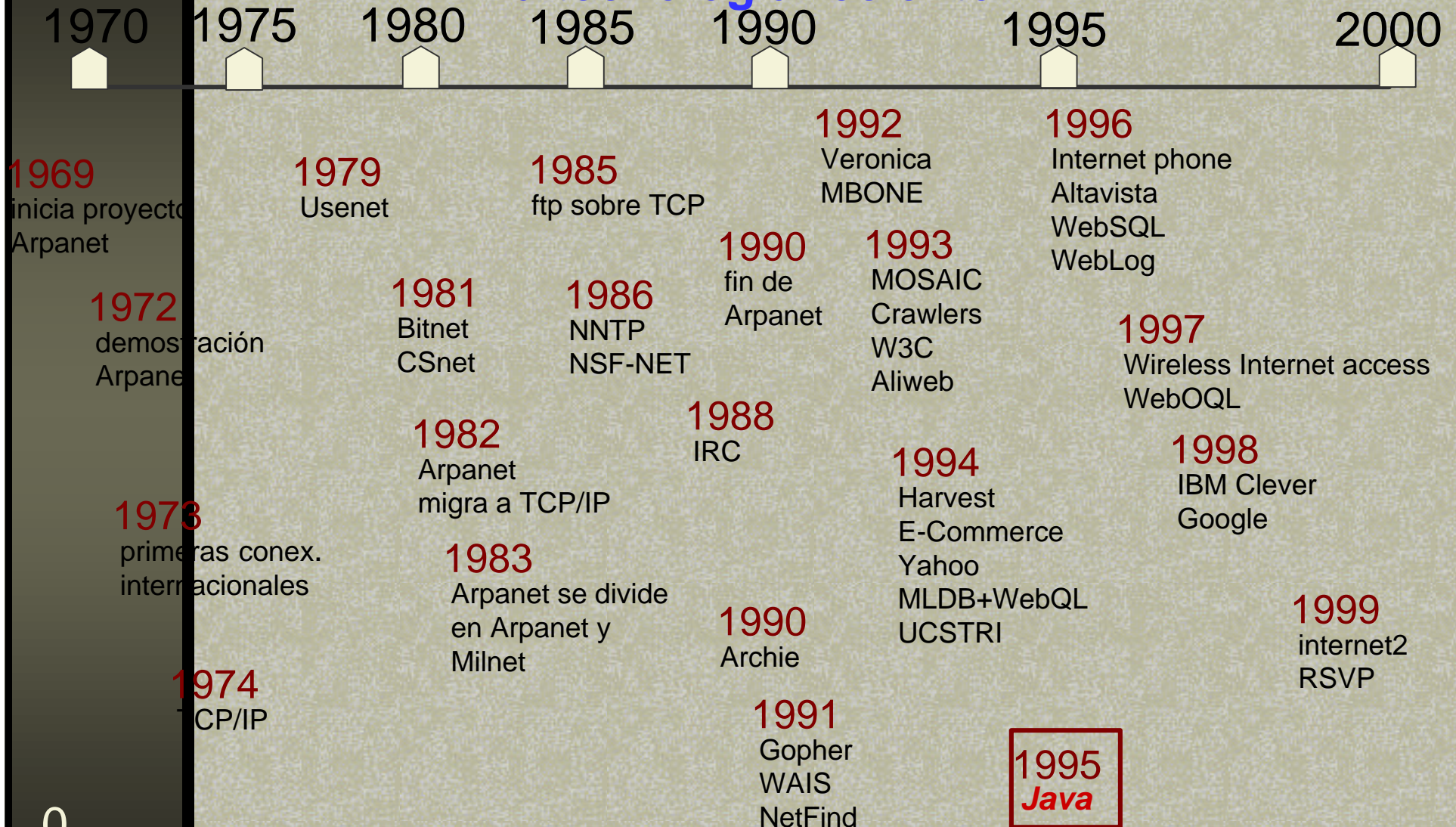
# Por qué WMI?

## c.Crecimiento en cantidad de sites



# Por qué Java?

a. tecnología reciente





# Por qué Java?

## b. tecnología flexible

- Diseñado para resolver problemas comunes de los programadores profesionales y mejorar su productividad
- Solución dentro del propio espacio de solución  
=> una línea de código puede hacer muchas cosas
- Soluciones más sencillas y leíbles por un entendido  
=> mejor mantenibilidad
- Uso sencillo y flexible de librerías  
(mejor programa ? mejor uso de librerías)  
=> el pgmador se desentiende de tareas menores  
(ej. inicialización y cleanup)



# Por qué Java?

c. tecnología transportable

- manejo localizado de errores  
=> fácil detección del código involucrado en su manejo
- Java fue diseñado con la complejidad necesaria para ser escalable.
- Java fue diseñado para vivir dentro y fuera de internet sin grandes cambios



# WM + JAVA

*Algunos casos*



# JAVA + WME

- **el caso C-BIRD**
- **el caso VWV**
- **el uso de metadatos**



JAVA + WME

C-BIRD

# JAVA + WML el caso C-BIRD

<http://jupiter.cs.sfu.ca/cbird/java/> (IE 4.0 or Netscape 4.0).

Web content mining  
s/imágenes

C-BIRD

IR

Index of /VML/cbird

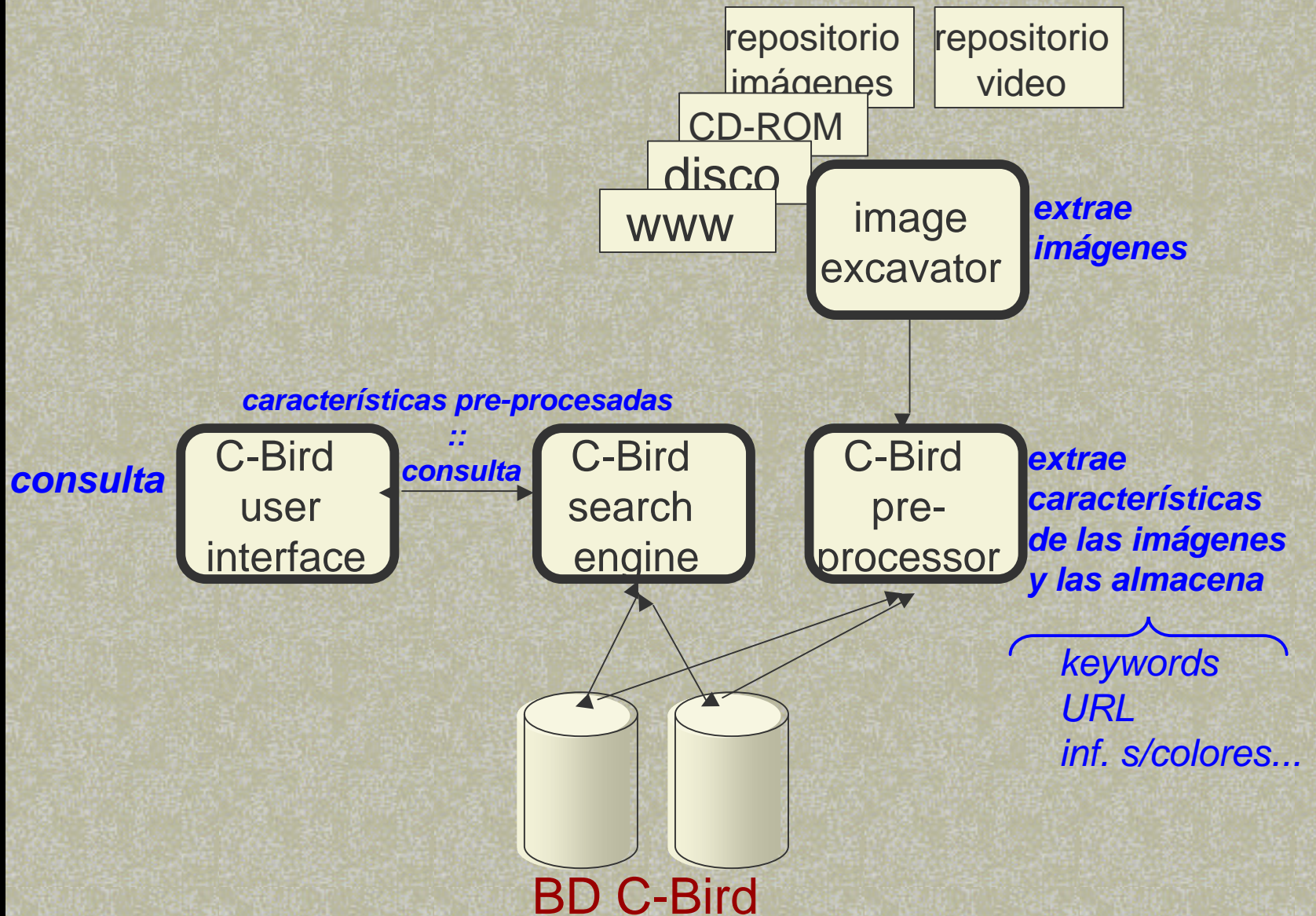
Name	Last modified	Size	Description
<a href="#">Parent Directory</a>		-	
<a href="#">CBIRD.cgi</a>	30-Oct-1997 20:23	133K	
<a href="#">CBIRD_GetImage.cgi</a>	26-Sep-1997 10:59	13K	
<a href="#">CBS.cgi</a>	30-Oct-1997 20:23	133K	
<a href="#">CBS.cgi.old</a>	06-Aug-1997 10:33	119K	
<a href="#">CBS_GetImage.cgi</a>	26-Sep-1997 10:59	13K	
<a href="#">CBS_SMC96.ps.Z</a>	04-Dec-1996 09:47	912K	
<a href="#">CbirdAbout.html</a>	18-Aug-1997 14:48	1.0K	
<a href="#">CbirdHelp.html</a>	13-Aug-1997 11:46	4.1K	
<a href="#">bin/</a>	19-Oct-1997 17:41	-	
<a href="#">cbird.cgi</a>	30-Oct-1997 20:23	133K	
<a href="#">colors/</a>	25-Sep-1997 15:00	-	

**Estructura**

- html
- Java user-interface
- C++ preprocessor
- Images DB

0

# JAVA + WML el caso C-BIRD



# JAVA + WWT: el caso C-BIRD

interface usuario  
*tipos de busqueda*

- búsqueda por histogramas de colores
- búsqueda por color layout  
(1x1, 2x2, 3x3, 4x4, 8x8)
- búsqueda por porcentaje (?5 colores)
- búsqueda por ?/? de metadatos
- búsqueda por illumination invariance  
(cromaticidad) (img normalizada)
- combinaciones y otras



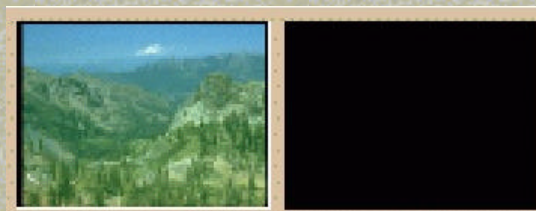
búsqueda por color layout



# JAVA + WINE el caso C-BIRD



búsqueda con características de la imagen



búsqueda con características de la imagen y **keywords** sacadas del contexto



JAVA + WME

WWW

# JAVA + WML: el caso VWV

Information browsing

Virtual Web View

Resource discovering

(Osmar Zaiane et al.)

## Estructura



query language

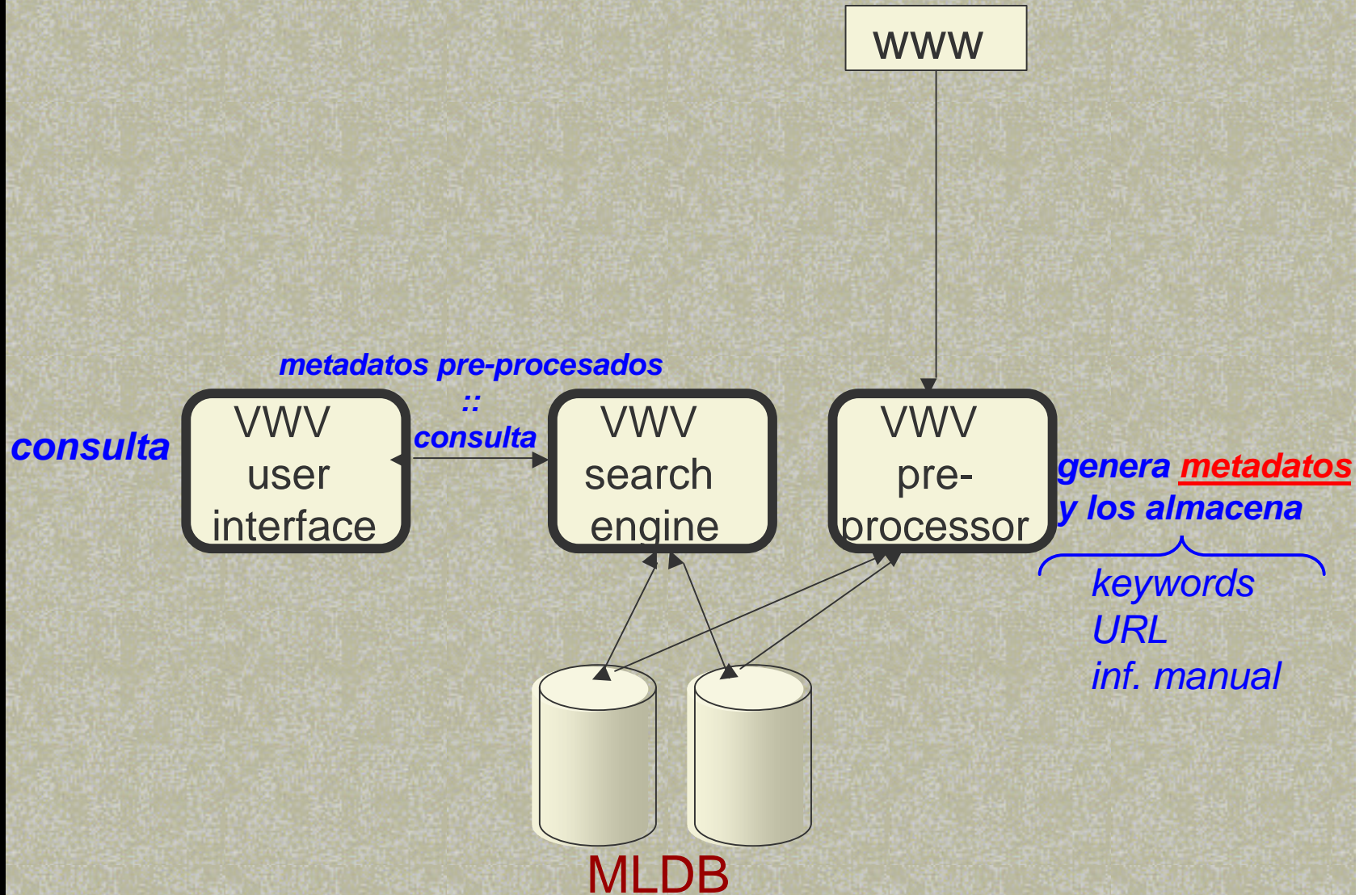
Web Modeling Language  
aplicación del XML

notación para especificar  
conceptualmente Web sites complejos

def

capas de abstracción sucesivas  
datos relativamente estructurados,  
manejables y administrable por DB<sup>s</sup>

# JAVA + WME: el caso VWV





JAVA + WME

METADATOS

# JAVA + XML: metadatos en WWW

keywords

metadatos

def: son datos acerca de los datos.

charla Java y el XML

XML  
+  
Java

def: lenguaje flexible del tipo "markup" (basado en tags)

API XML

JAXP

def: API p/procesar XML

Interpreta y presenta XML

Indep de cualquier procesador XML

parser XML

SAX DOM

Document Object Model

representación XML estándar de contenidos y modelos de documentos (Bindings para JavaScript, Java, C++, etc.)

def: verifica si el XML está bien formado

JXerces

basado en

Std Application

XML

administrador de docs XML basado en eventos

permite manejo secuencial sin necesidad de cargarlo en MC

# JAVA + WNE

## metadatos en WWW (WWW como web semántica)

*def* (Tim Berners-Lee): extensión de la web que permite a humanos y computadoras trabajar en cooperación

### Web semántica

- Knowledge Acquisitions
- Knowledge Representations
- Agent Systems
- Multi Agent Systems
- Ontology

charla IA en Java  
charla NN en Java

incorpora el paradigma de **Soft-Computing**

*Machine Intelligence JSRs* (Java Specification Request)  
en la JCP (Java Community Process)



FIN