

Computación distribuida usando redes P2P: Caso práctico en el área de recuperación de Información

Jorge Camargo
`jor-cama@uniandes.edu.co`

2a Jornada de Computación Paralela

April 24, 2008

Agenda

- 1 Introducción
- 2 Conceptos
- 3 IR en redes P2P
- 4 Retos



Introducción

- Motores de búsqueda
- Volumen de información
- Contenidos altamente dinámicos
- Transparencia, eficiencia, flexibilidad, escalabilidad y fiabilidad?
- Sistemas distribuidos: Redes P2P



Agenda

- 1 Introducción
- 2 Conceptos
 - Recuperación de información
 - Redes P2P
- 3 IR en redes P2P
- 4 Retos



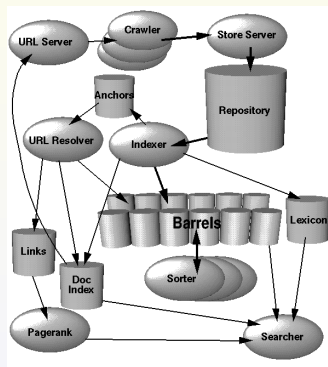
Agenda

- 1 Introducción
- 2 Conceptos
 - Recuperación de información
 - Redes P2P
- 3 IR en redes P2P
- 4 Retos



Motores de búsqueda

Arquitectura Google



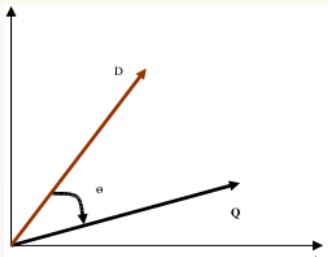
[Sergey1998]



Vector Space Model

En VSM los documentos se modelan como un vector n-dimensional

$$d_i \longrightarrow \vec{d_i} = (w(t_1, d_i), w(t_2, d_i), w(t_3, d_i) \cdots w(t_k, d_i))$$

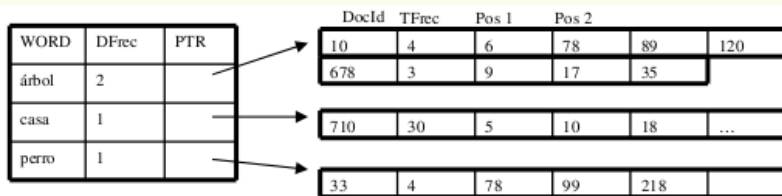


$$sim(Q, D_j) = \frac{Q \bullet D_j}{|Q| \cdot |D_j|}$$



Vector Space Model

Anatomía de un índice invertido



Otros conceptos de IR

- Stemming
- Stopwords
- Sinonimia
- Polisemia
- Recall
- Precision



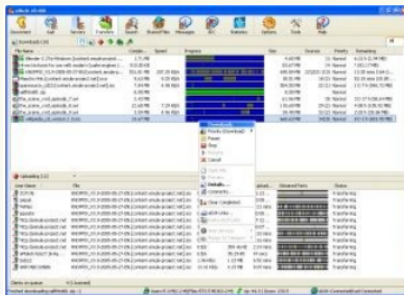
Agenda

- 1 Introducción
- 2 Conceptos
 - Recuperación de información
 - Redes P2P
- 3 IR en redes P2P
- 4 Retos



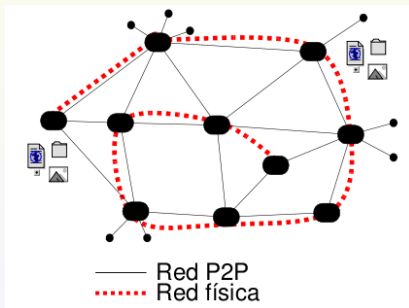
Qué se conoce como una red P2P

Alguien conoce a esta mula?

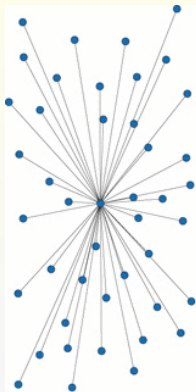


Qué es una red P2P

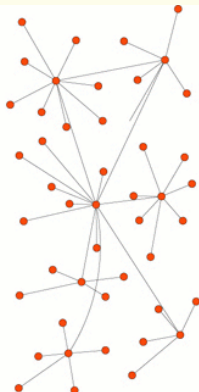
- Cada nodo hace la función de servidor y cliente (servant)
- Los nodos ingresan y salen de la red en cualquier momento
- Sistema distribuido



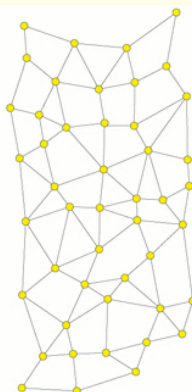
Topologías de red P2P



RED CENTRALIZADA



RED DESCENTRALIZADA



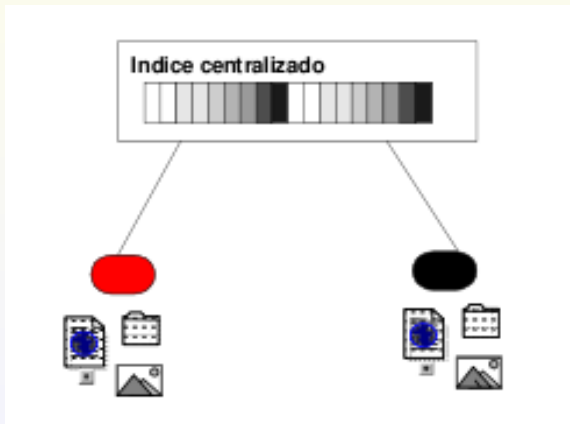
RED DISTRIBUIDA

[Wikipedia]



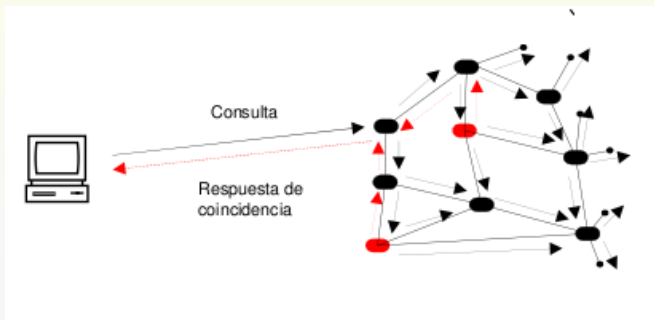
Redes P2P centralizadas

Basadas en un índice central (Napster y Audiogalaxy)



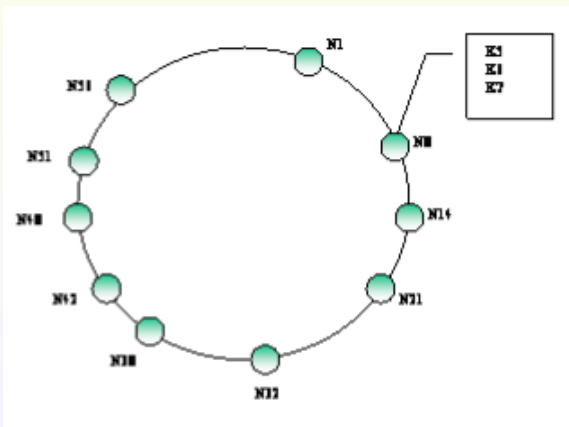
Redes P2P no estructuradas

Descentralizadas utilizando inundación de mensajes (Gnutella)



Redes P2P estructuradas

Descentralizadas utilizando tablas hash distribuidas (Chord, Pastry, CAN)



Agenda

- 1 Introducción
- 2 Conceptos
- 3 IR en redes P2P
 - IR en redes no estructuradas
 - IR en redes estructuradas
 - Nuestra propuesta
- 4 Retos



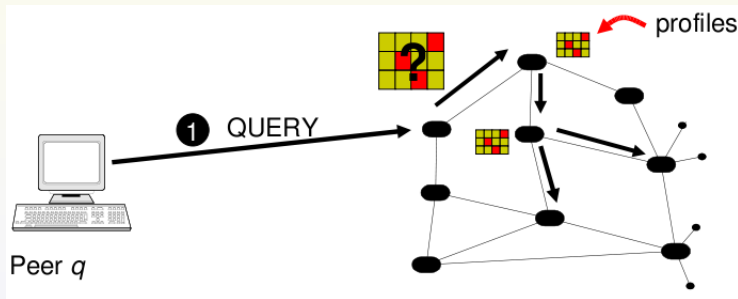
Agenda

- 1 Introducción
- 2 Conceptos
- 3 IR en redes P2P
 - IR en redes no estructuradas
 - IR en redes estructuradas
 - Nuestra propuesta
- 4 Retos



Mecanismo Inteligente de Búsqueda

- Inundación de q
- Cada nodo almacena un perfil de consultas repondidas
- Algunos nodos nunca serán contactados (TTL)



[Zeinalipour2002]



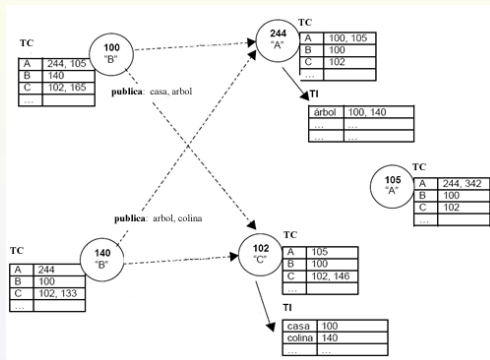
Agenda

- 1 Introducción
- 2 Conceptos
- 3 IR en redes P2P
 - IR en redes no estructuradas
 - IR en redes estructuradas
 - Nuestra propuesta
- 4 Retos



aDICS

- Cada nodo administra una letra del alfabeto
- Cada nodo reporta términos a nodos responsables
- Intersección de nodos que responden



[Tolosa2004]



KSS

- A, B, C y D son las palabras de un documento docID
- Se crean entradas al índice para las 6 combinaciones
- Se calcula el hash (SHA1) de la concatenación de cada par de términos para publicarlos en CHORD
- En una búsqueda de ABCD se localiza al nodo responsable del índice de AB, dicho nodo busca localmente para luego reenviar su resultado al nodo responsable por CD. Este último realiza la intersección de los nodos encontrados en AB y CD.

[Gnawali2002]



Agenda

- 1 Introducción
- 2 Conceptos
- 3 IR en redes P2P
 - IR en redes no estructuradas
 - IR en redes estructuradas
 - Nuestra propuesta
- 4 Retos



APSE: A P2P Search Engine

Motor de búsqueda distribuido basado en una red P2P estructurada, utilizando técnicas de VSM para recuperación de información

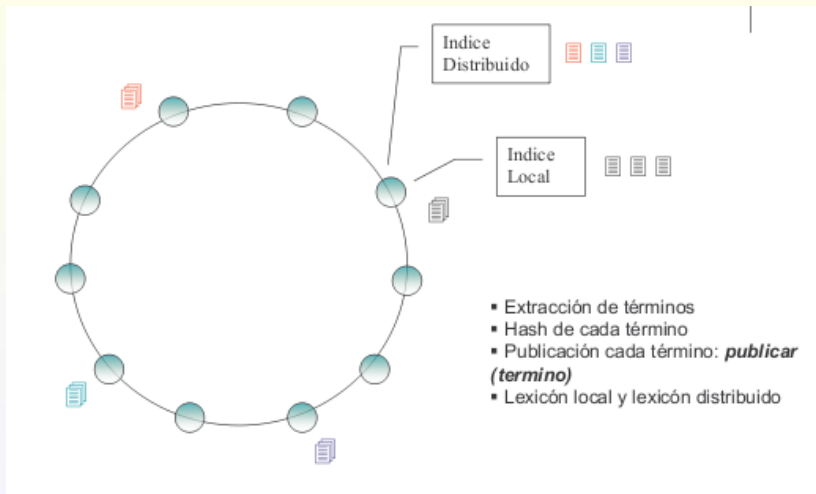
[Camargo2006]



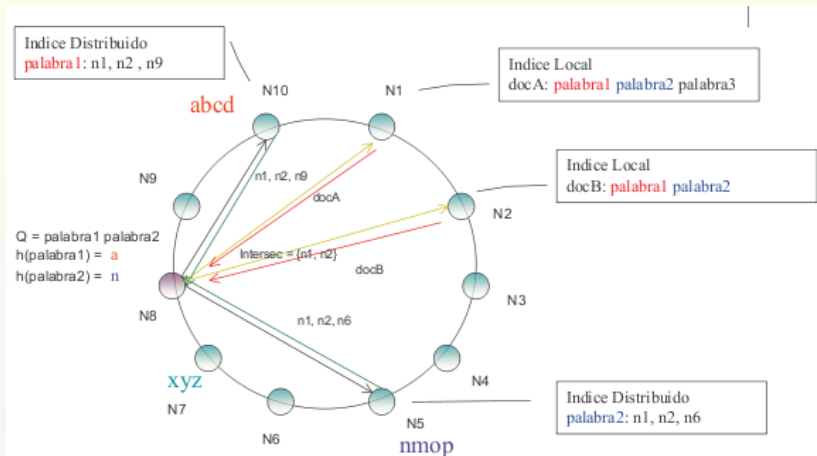
Componentes



Indexación



Consultas



Mensajes

Cantidad de mensajes para resolver una consulta:

$$\#msgs = 2(n + insec(n, k))$$

en donde n es el número de nodos e $insec(n, k)$ es el número de nodos intersección para una consulta de k términos. Por ejemplo, para una red con $n=104$ una consulta con $k=4$ términos, el número máximo de mensajes es $\#msgs=2*(104+4)=216$.



Experimentación

- Lenguaje de programación Java
- Basado el API FreePastry ampliamente utilizado por la comunidad
- Los experimentos fueron realizados en una red LAN 10/100 de 20 PCs (2G RAM, PIV 3.4 GHz)
- Corpus de Reuters-21578 generada con dataGen [Zeinalipour2002]



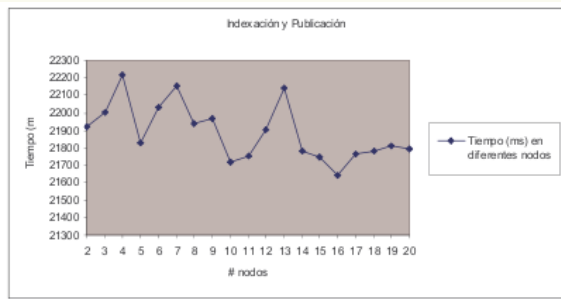
Metodología

- Red controlada
- Dataset de consultas elegidas aleatoriamente, de un término hasta 10 términos
- Promedio de mediciones
- 1 nodo de bootstrap
- Recolección automática de logs (scripts de shell)

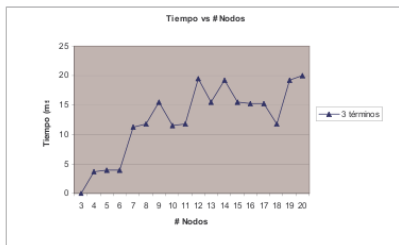


Tiempo de indexación

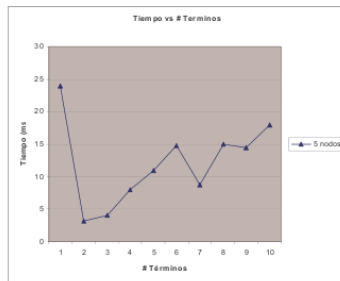
Tiempo utilizado por APSE para indexar el dataset (un nodo) en configuraciones desde 3 hasta 20 nodos



Tiempo de consulta



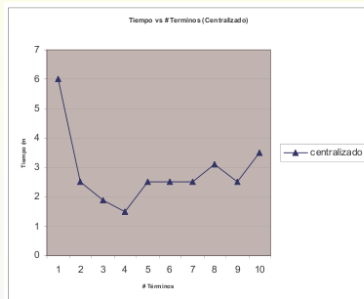
Tiempo utilizado para resolver una consulta de 3 términos incrementando la cantidad de nodos en el sistema



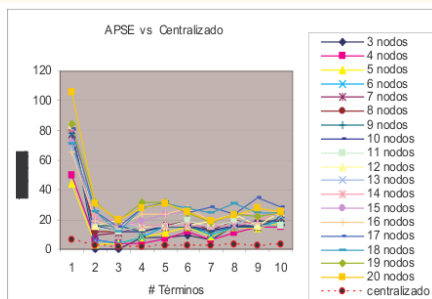
Tiempo utilizado para resolver una consulta en una configuración del sistema con 5 nodos incrementando la cantidad de términos de una consulta



Tiempo de consulta



Tiempo utilizado por el sistema centralizado para consultas incrementando la cantidad de términos



Tiempo de respuesta del sistema centralizado versus APSE



Cantidad de mensajes

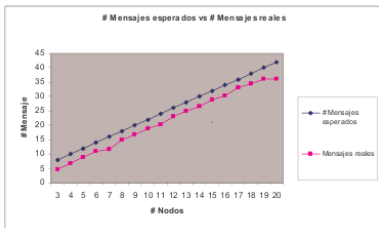
Cantidad de mensajes para resolver una consulta:

$$\#msgs = 2(n + insec(n, k))$$

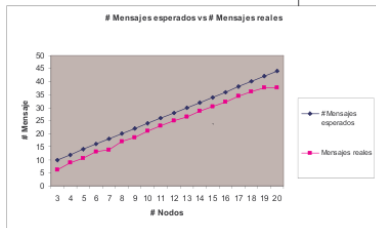
en donde n es el número de nodos e $insec(n, k)$ es el número de nodos intersección para una consulta de k términos. Por ejemplo, para una red con $n=104$ una consulta con $k=4$ términos, el número máximo de mensajes es $\#msgs=2*(104+4)=216$.



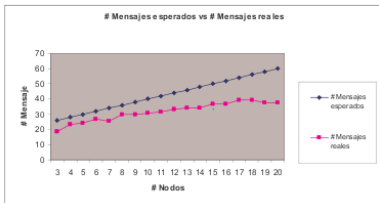
Cantidad de mensajes



Cantidad de mensajes esperados versus cantidad mensajes obtenidos en la experimentación para consultas de 1 término



Cantidad de mensajes esperados versus cantidad mensajes obtenidos en la experimentación para consultas de 2 términos



Cantidad de mensajes esperados versus cantidad mensajes obtenidos en la experimentación para consultas de 20 términos

$$\# \text{ msgs} = 2 (n + \ln \sec(n, k))$$



Agenda

- 1 Introducción
- 2 Conceptos
- 3 IR en redes P2P
- 4 Retos



Retos para sistemas distribuidos

Sistemas distribuidos

- TCP/IP es el más adecuado?
- Qué pasó finalmente con el esperado Internet 2?
- Otras maneras más eficientes para distribuir en índice (Grid Computing)

Recuperación de información

- Cómo indexar objetos diferentes a texto (imágenes, videos, música, secuencias ADN)
- Cómo visualizar colecciones de objetos con estructura compleja



Preguntas

Preguntas



Referencias I



[Sergey, 1998] Sergey Brin y Lawrence Page.

The Anatomy of a Large-Scale Hypertextual Web Search Engine, 1998



[Zeinalipour, 2002] Zeinalipour.

Zeinalipour, University of California Riverside. CIKM 2002



[Tolosa, 2004] Gabriel Tolosa, Fernando Bordignon y Jorge A Peri.

aDICS: Modelo de índice Distribuido sobre una Red P2P para Búsquedas por Contenido. X Congreso Argentino de Ciencias de la Computación, Octubre 2004



[Wikipedia, 2008] Wikipedia

Peer-to-peer, consultado en abril de 2008



Referencias II



[Gnawali, 2002] O. Gnawali.

A Keyword Set Search System for PeertoPeer Networks.
Master's thesis, Massachusetts Institute of Technology, 2002



[Camargo, 2006] Jorge Camargo y Francisco Rueda

APSE: A P2P Search Engine. Tesis de maestría. Universidad
de Los Andes, 2006

