



HY463 - Συστήματα Ανάκτησης Πληροφοριών
Information Retrieval (IR) Systems

**Ανάκτηση Πληροφοριών
& Συστήματα Ομοτίμων**
(Peer-to-Peer Systems)



Γιάννης Τζίτζικας

Διάλεξη : 18-19

Ημερομηνία : 1-6-2007

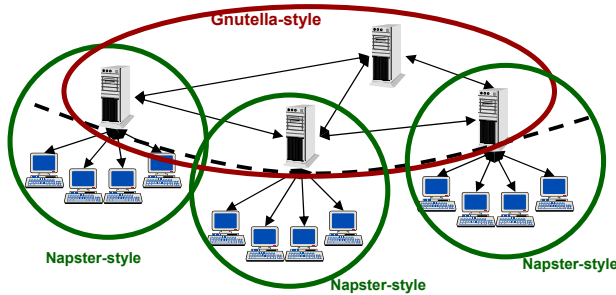
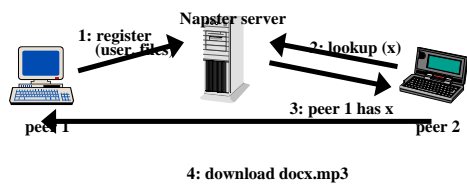


**Τι διαφέρει η Ανάκτηση σε
P2P συστήματα από την Κατακεκομημένη Ανάκτηση;**

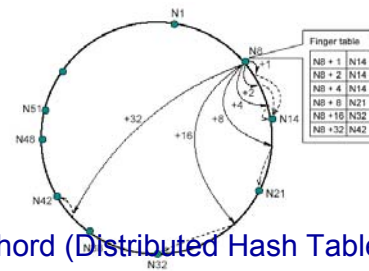
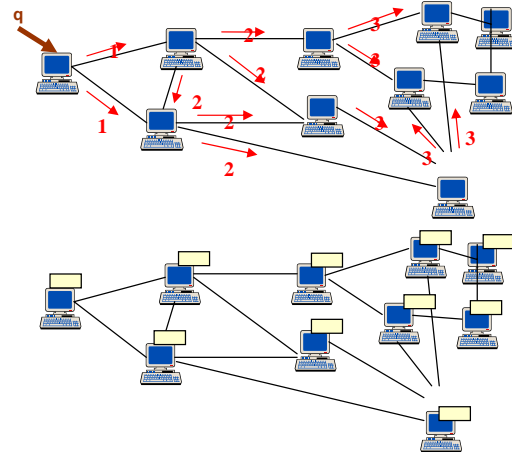
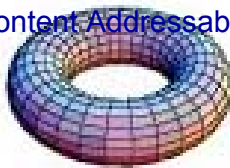
Η ανάκτηση πληροφοριών σε συστήματα ομοτίμων είναι μια περίπτωση κατακεκομημένης ανάκτησης

Ιδιαιτερότητες των ομοτίμων συστημάτων:

- Υπερβολικά μεγάλος αριθμός πηγών (peers)
- Μεγαλύτερη αυτονομία πηγών
- Έλλειψη Σταθερότητας, Ελέγχου, Προβλεψιμότητας
 - (not stable, controllable, unpredictable)
- Επιτακτική ανάγκη για μείωση του κόστους επικοινωνίας



CAN (Content Addressable Network)



Chord (Distributed Hash Table -DHT)

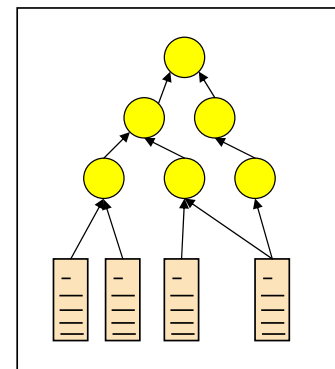
P2P and IR: Περίπτωση: Κατηγοριοποιημένα Έγγραφα

Έστω ότι κάθε έγγραφο είναι ταξινομημένο σε μια κατηγορία ενός ελεγχόμενου ευρετηρίου (ODP, Yahoo!). Ο χρήστης κάνει αναζήτηση δίνοντας μια κατηγορία

έγγραφο \approx mp3 αρχείο
κατηγορία εγγράφου \approx τίτλος του mp3 αρχείου

Άρα μπορούμε να φτιάξουμε ένα ομότιμο σύστημα

- τύπου Napster (Hybrid P2P)
- τύπου Gnutella (Pure P2P)
- τύπου Kazaa (Hierarchical P2P)
- τύπου Freenet (Structured P2P)
- τύπου Chord (Structured P2P)
- τύπου CAN (Structured P2P)



Βέβαια έτσι θα μπορούμε να κάνουμε ανάκτηση μόνο βάσει του τίτλου (και όχι βάσει του περιεχομένου).

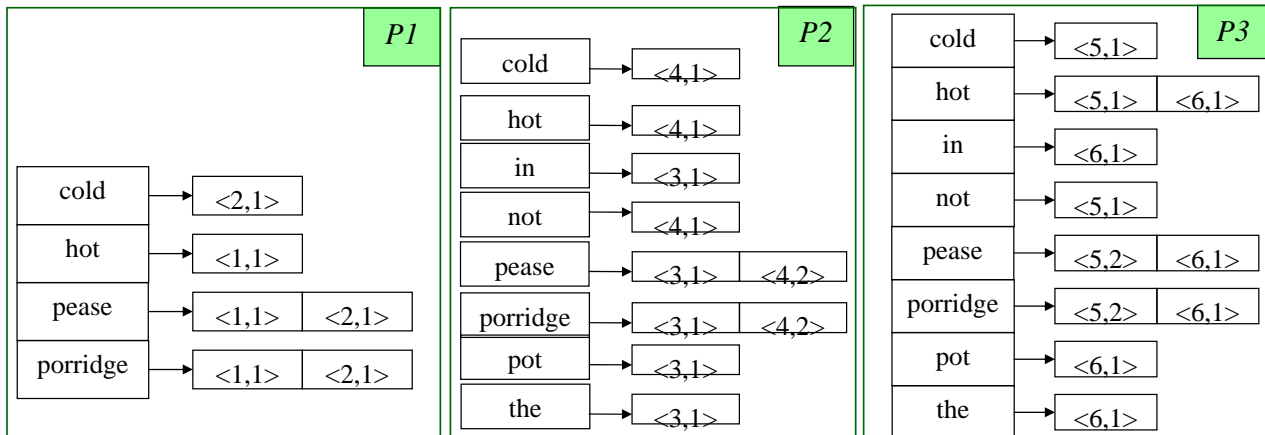


P2P and Statistical IR

Τυπικό Ευρετήριο P2P

p	k	d
121.111.86	"Singing in the Rain"	SR.mp3
121.111.86	«Υπάρχω»	stelios.mp3
222.18.78	"Singing in the Rain"	SingRain.mp3

Τυπικό Ευρετήριο IRS
(document partitioning)

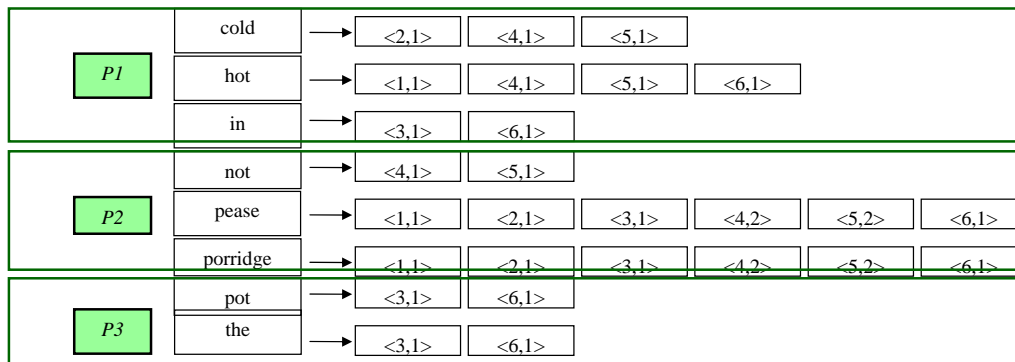


P2P and Statistical IR

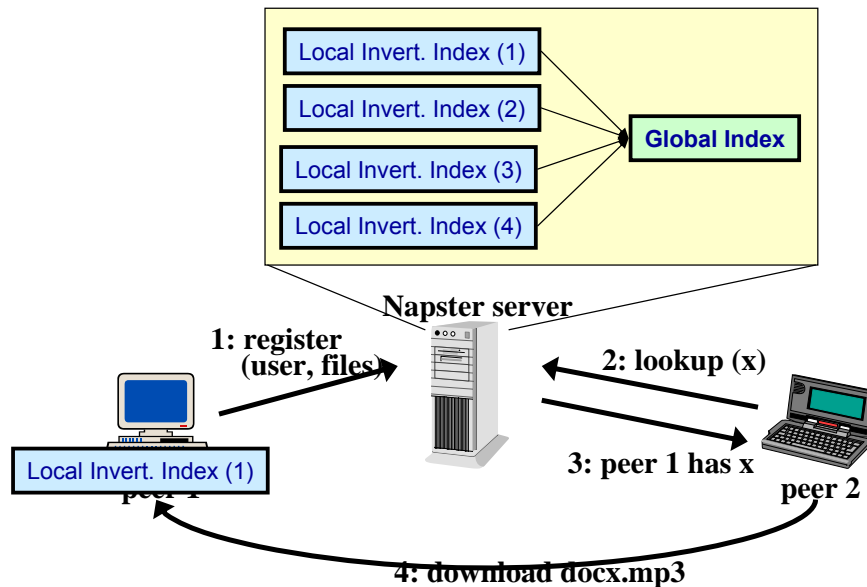
Τυπικό Ευρετήριο P2P

p	k	d
121.111.86	"Singing in the Rain"	SR.mp3
121.111.86	«Υπάρχω»	stelios.mp3
222.18.78	"Singing in the Rain"	SingRain.mp3

Τυπικό Ευρετήριο IRS
(term partitioning)



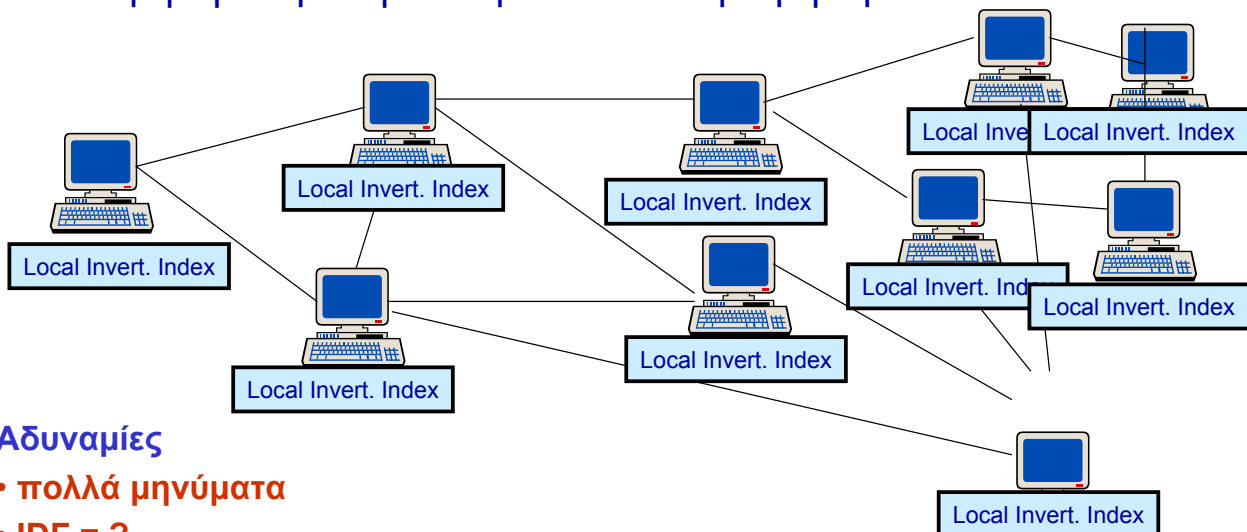
- Ένας κεντρικός εξυπηρετητής αποθηκεύει όλα τα ανεστραμμένα ευρετήρια των κόμβων



αδυναμίες:

- ο εξυπηρετητής χρειάζεται πολύ χώρο
- χρονοβόρο upload των ευρετηρίων στον εξυπηρετητή,
- το κόστος αποτίμησης επερωτήσεων πάει εξ' ολοκλήρου στον εξυπηρετητή
- ≈ Google, χωρίς το crawling (συλλογή σελίδων) και έχοντας έτοιμα κομμάτια του ευρετηρίου

- Κάθε κόμβος συντηρεί το ανεστραμμένο ευρετήριο των εγγράφων του.
- Αποτίμηση επερωτήσεων με κατακλυσμό μηνυμάτων



Αδυναμίες

- πολλά μηνύματα
- IDF = ?



P2P and IR: (Gnutella-style)

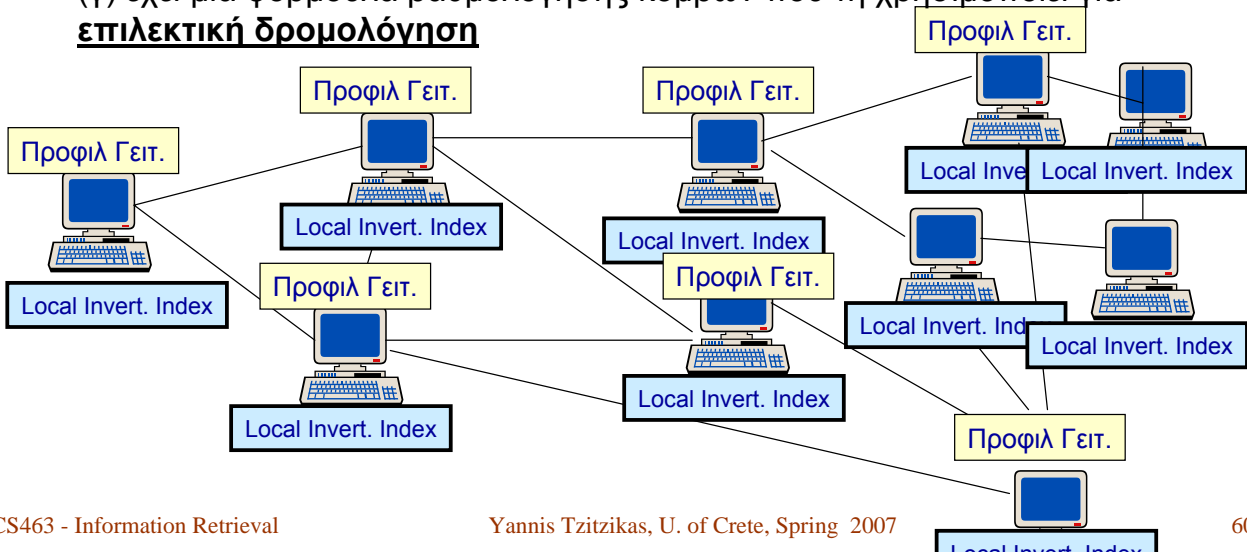
Παραλλαγές του Κατακλυσμού μηνυμάτων

- BFS: Breadth First Search (=Gnutella)
- RBFS: κάθε κόμβος προωθεί ένα μήνυμα σε ένα τυχαίο ποσοστό (π.χ. 20%) των γνωστών του κόμβων
 - + πιθανοκρατικός αλγόριθμος
 - - μπορεί το μήνυμα να μην πάει σε κόμβους που έχουν συναφή αντικείμενα
- 1-Random Walker:
 - κάθε κόμβος προωθεί ένα μήνυμα σε έναν τυχαία επιλεγμένο κόμβο από τους γνωστούς του
- k-Random Walkers:
 - κάθε κόμβος προωθεί ένα μήνυμα σε κ τυχαία επιλεγμένους κόμβους από τους γνωστούς του
 - + λιγότερα μηνύματα από το RDFS



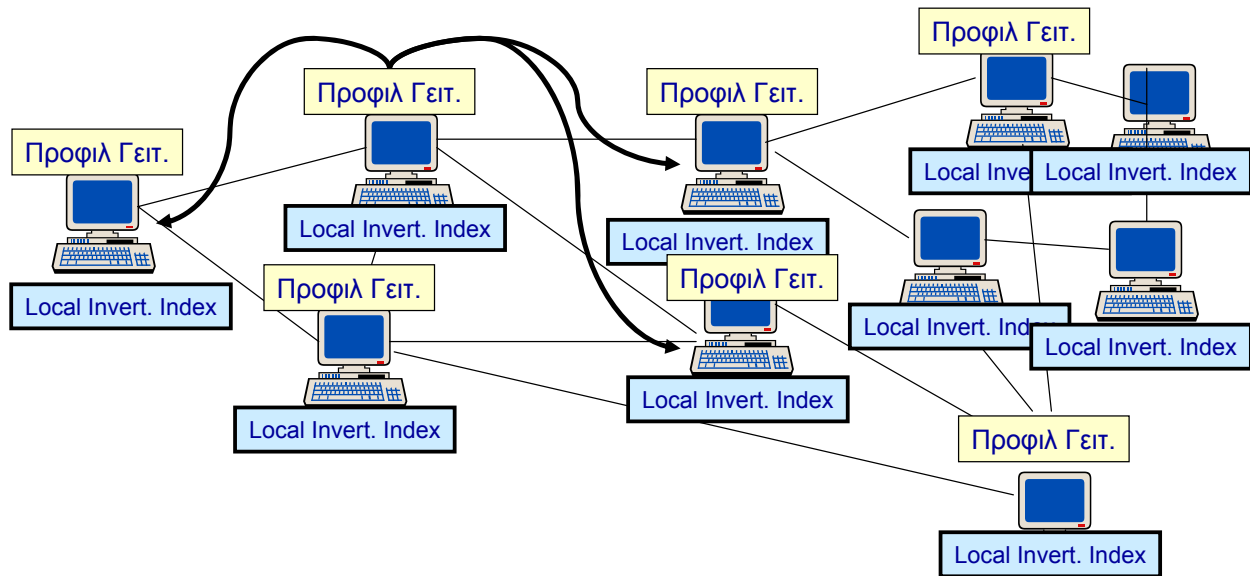
P2P and IR (Freenet-style)

- Κάθε κόμβος:
 - (α) συντηρεί το ανεστραμμένο ευρετήριο των εγγράφων του.
 - (β) φτιάχνει ένα προφίλ των γειτόνων του βασισμένο στις επερωτήσεις του παρελθόντος
 - (γ) έχει μια φόρμουλα βαθμολόγησης κόμβων που τη χρησιμοποιεί για επιλεκτική δρομολόγηση





P2P and IR (Freenet-style)



Προφίλ Γειτόνων βάσει των προηγούμενων απαντήσεων

p	q	$ ans(q) $	LRU (Least Recently Used) deletion policy
121.111.86	"Singing in the Rain"	4	
121.111.86	«Υπάρχω»	12	
222.18.78	"One" U2One.mp3	6	
131.161.86	"Dog song"	7	
131.111.86	«Υπήρξα»	4	
122.118.73	"Pop2" 12	8	
144.44.46	"Singing in the Rain"	5	
155.11.86	«Υπάρχω»	7	

Local Invert. Index

- Το **προφίλ** είναι τριάδες της μορφής $(p_j, q, |ans(p_j, q)|)$
 - όπου p_j ένας γείτονας, q μια επερώτηση που απήντησε αυτός ο γείτονας, και $|ans(p_j, q)|$ το μέγεθος της απάντησης
 - LRU update policy



Προφίλ Γειτόνων και Δρομολόγηση: >RES (περισσότερα αποτελέσματα)

p	q	ans(q)
121.111.86	"Singing in the Rain"	4
121.111.86	«Υπάρχω»	12
222.18.78	"One" U2One.mp3	6
131.161.86	"Dog song"	7
131.111.86	«Υπήρξα»	4
122.118.73	"Pop2" 12	8
144.44.46	"Singing in the Rain"	5
155.11.86	«Υπάρχω»	7

Σκορ(121.111.86)=16
Σκορ(122.118.73)=8

Local Invert. Index

- Για την δρομολόγηση μιας επερώτησης επιλέγονται εκείνοι οι γείτονες που έχουν δώσει τα περισσότερα αποτελέσματα στο παρελθόν (εξ ου και το όνομα >RES). Συγκεκριμένα στις προηγούμενες n επερωτήσεις.
- Το **σκορ ενός γείτονα p_j είναι**
 - $Score(p_j) = \sum \{ |ans(p_j, q_j)| \mid q_j \text{ answered by } p_j \text{ in the past} \}$



Προφίλ Γειτόνων και Δρομολόγηση: >RES και ομοιότητα επερωτήσεων

p	q	ans(q)
121.111.86	"Singing in the Rain"	4
121.111.86	«Υπάρχω»	12
222.18.78	"One" U2One.mp3	6
131.161.86	"Dog song"	7
131.111.86	«Υπήρξα»	4
122.118.73	"Pop2" 12	8
144.44.46	"Singing in the Rain"	5
155.11.86	«Υπάρχω»	7

Local Invert. Index

- Για την δρομολόγηση μιας επερώτησης q επιλέγονται εκείνοι οι γείτονες που έχουν δώσει τα περισσότερα αποτελέσματα στο παρελθόν (>RES) σε **επερωτήσεις που είναι κοντινές με το q**



Προφίλ Γειτόνων και Δρομολόγηση:
>RES και ομοιότητα επερωτήσεων

Το σκορ ενός γείτονα p_j δοθείσας επερώτησης q , είναι:

>RES

$$\text{Score}(p_j) = \sum \{ |\text{ans}(p_j, q_j)| \mid q_j \text{ answered by } p_j \text{ in the past} \}$$

>RES και ομοιότητα επερωτήσεων

$$\text{Score}(p_j, q) = \sum \{ |\text{ans}(p_j, q_j)| * \text{sim}(q_j, q)^\alpha \mid q_j \text{ answered by } p_j \text{ in the past} \}$$

- $\text{sim}(q_j, q)$: Π.χ. ομοιότητα συνημίτονου
- α : παράμετρος για το καθορισμό της σπουδαιότητας μεταξύ συνάφειας και μεγέθους απάντησης



Προφίλ Γειτόνων και Δρομολόγηση:
>RES και ομοιότητα επερωτήσεων

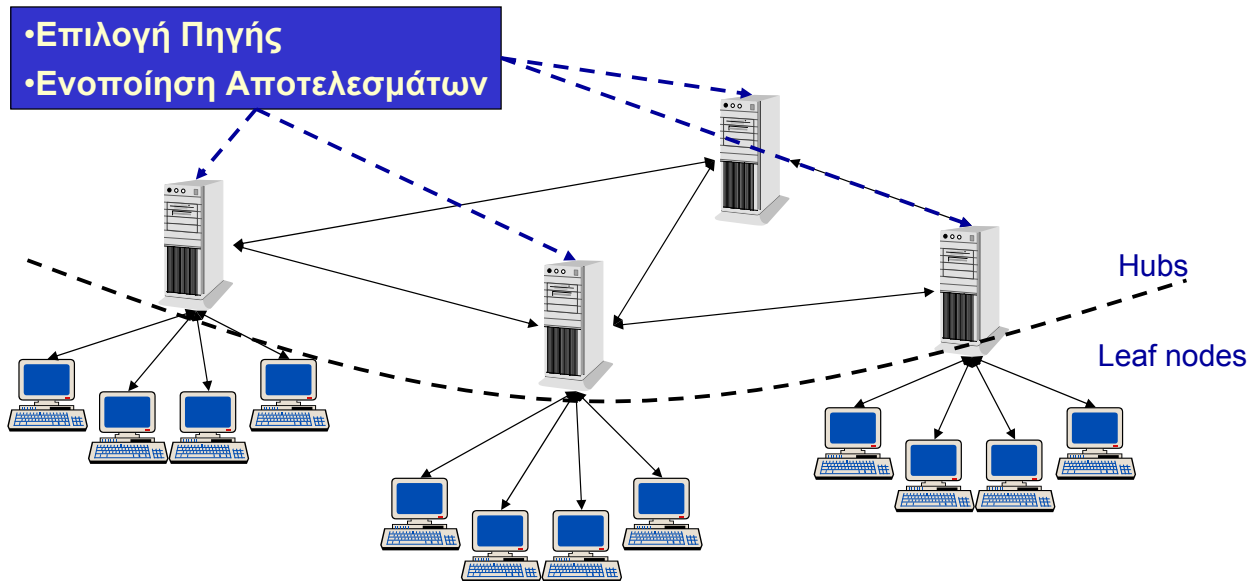


- Πότε αυτή η προσέγγιση είναι καλή;
- Απ: Όταν τα έγγραφα του κάθε κόμβου είναι σημασιολογικά κοντινά
 - Ποια η διαφορά με το Freenet ?
- Επειδή αυτό όμως δεν συμβαίνει πάντα η επερώτηση προωθείται και σε έναν τυχαία επιλεγμένο γείτονα. // αυτό επίσης συμβάλει στην καλή εκκίνηση του συστήματος

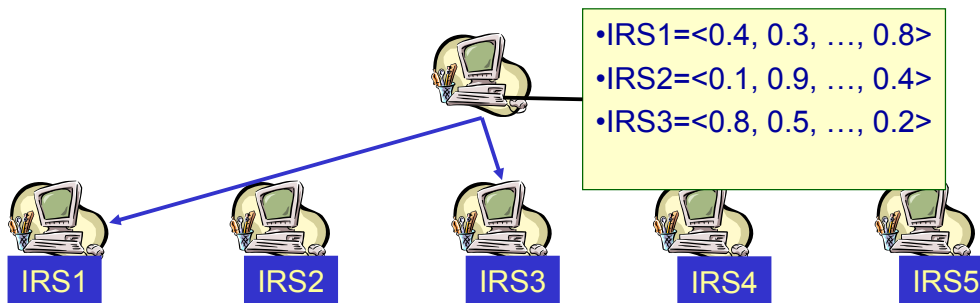


Ανάκτηση Κειμένων σε **Ιεραρχικά Ομότιμα Συστήματα** (Kazaa-style IR)

Γενική Ιδέα: Κάνουμε ό,τι και στην κατακευμαμένη, απλά εδώ έχουμε πολλούς μεσίτες
Κάθε μεσίτης (εδώ super-peer) έχει μια περιγραφή των περιεχομένων των υποκείμενων κόμβων



Επανάληψη: Επιλογή Πηγής με **Διανύσματα Πηγών**



- Βλέπουμε **κάθε συλλογή** ως ένα **μεγάλο έγγραφο**
- Φτιάχνουμε ένα **διάνυσμα για κάθε συλλογή** (τύπου TF-IDF)
 - tf_{ij} : συνολικές εμφανίσεις του όρου i στη συλλογή j
 - idf_i : $\log(N/n_i)$, όπου N το πλήθος των συλλογών, και n_i το πλήθος των συλλογών που έχουν τον όρο i
- Υπολογίζουμε το **βαθμό ομοιότητας** κάθε νέας επερώτησης με το **διάνυσμα κάθε συλλογής** (π.χ. ομοιότητα συνημίτονου)
- **Διατάσσουμε** τις συλλογές και **επιλέγουμε τις κορυφαίες**

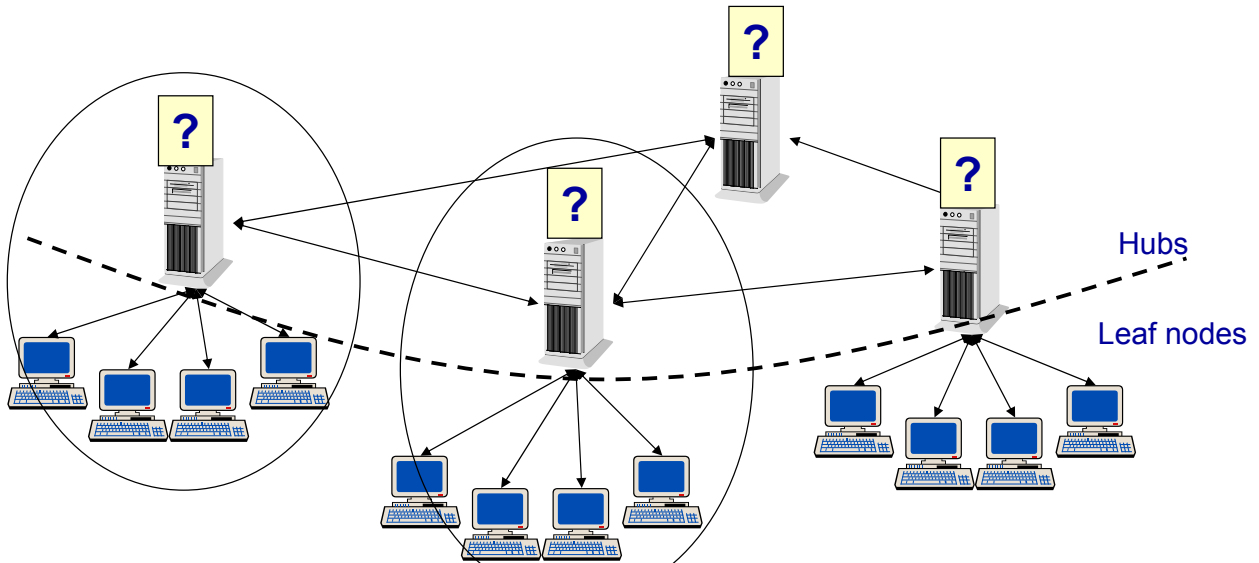
Εναλλακτικά: Αντί για ένα, μπορούμε να περιγράψουμε κάθε πηγή με K διανύσματα



Ανάκτηση Κειμένων σε *Ιεραρχικά Ομότιμα Συστήματα* (Kazaa-style IR)

Περιγραφή των περιεχομένων των φύλλων

Ανάγκη για μείωση του αποθηκευτικού χώρου στα Hubs



Ανάκτηση Κειμένων σε *Ιεραρχικά Ομότιμα Συστήματα*

Επιλογές

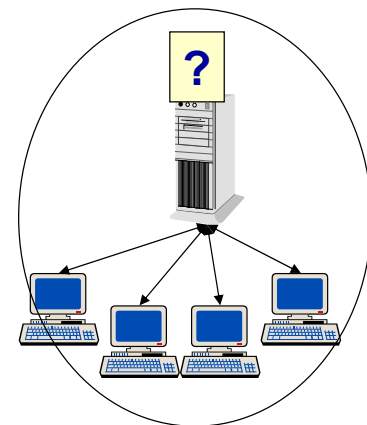
1/ Λεξιλόγια των υποκείμενων κόμβων +
συχνότητες εμφάνισής τους

- (δεν ξέρουμε το καθολικό λεξιλόγιο για να φτιάξουμε το διάνυσμα πηγής)

2/ Λεξιλόγια των υποκείμενων κόμβων

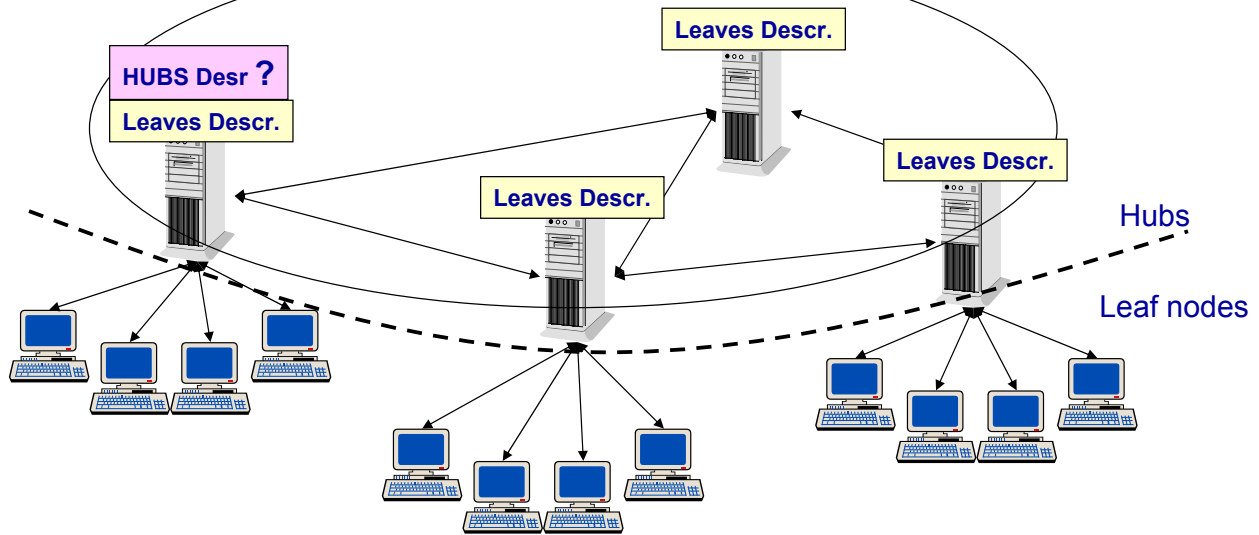
3/ Λέξεις που εμφανίζονται πάνω από 1
φορά + συχνότητές τους

- λόγω του νόμου του Zipf, ο απαιτούμενος αποθηκευτικός χώρος μειώνεται στο μισό





Περιγραφή των περιεχομένων των άλλων Hub ?



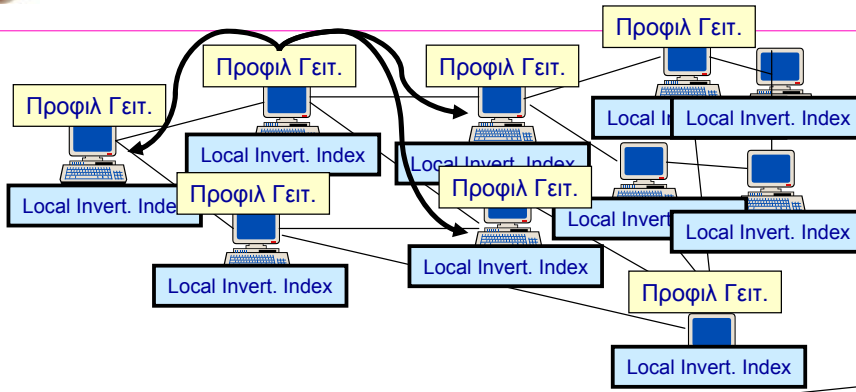
- Η περιγραφή ενός HUB είναι η ένωση των περιγραφών των υποκείμενων του κόμβων (Πρόβλημα: χώρος)
- Καταγραφή προηγούμενων επερωτήσεων που έχουν απαντηθεί
 - π.χ. >RES και ομοιότητα επερώτησης



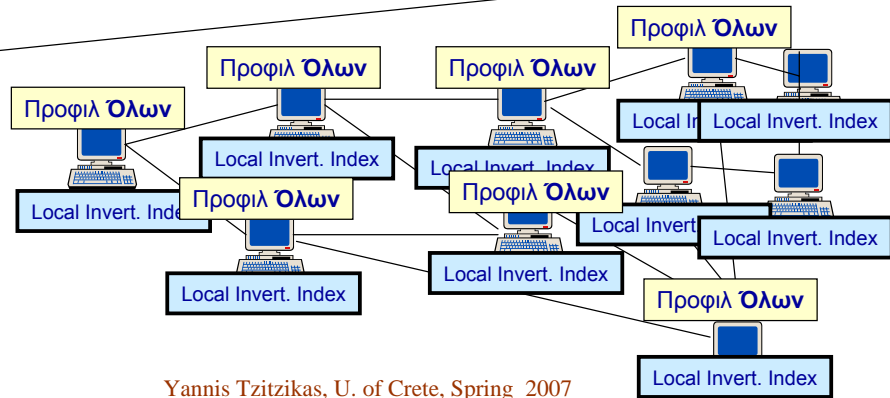
- A Client node sends its query to each of its connecting hubs.
- A hub that receives the query uses its **resource selection algorithm** to rank and select one or more neighboring leaf nodes as well as hubs, and routes the query to them if the message's TTL hasn't reached 0.
- A leaf node that receives the query message uses its **document retrieval algorithm** to generate a relevance ranking of its documents and responds with a queryhit message to include a list of top-ranked documents.
- Each top-level hub (the hub that connects directly to the client node that issues the request) collects the queryhit messages and uses its **result merging algorithm** to merge the documents retrieved from multiple leaf nodes into a single, integrated ranked list and returns it to the client node.
- If the client node issues the request to more than one hub, then it also needs to merge results returned by multiple toplevel hubs.



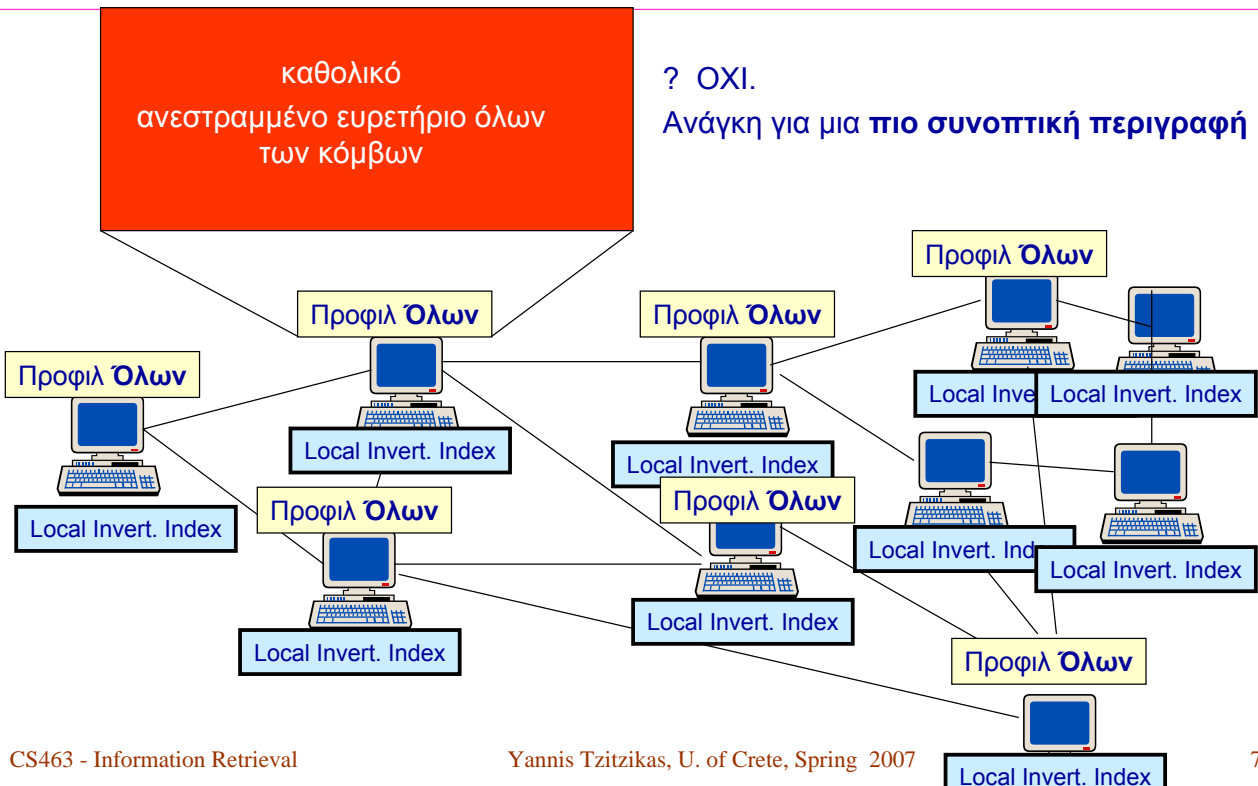
P2P and IR: Το σύστημα PlanetP



PlanetP



P2P and IR: Το σύστημα PlanetP





P2P and IR: Το σύστημα PlanetP Επιλογή Πηγής ?

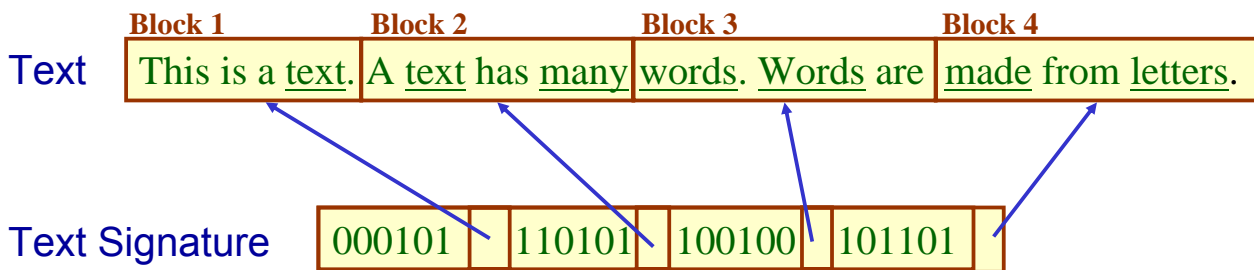
Θυμηθείτε:
Κατανομή Συναφών Εγγράφων
(Relevant document distribution (RDD))
Διανύσματα Πηγών

- Στο σύστημα PlanetP το **λεξιλόγιο του κάθε κόμβου** (όχι οι λίστες των εμφανίσεων των όρων) περιγράφεται με ένα **Bloom φίλτρο**
- Τα φίλτρα Bloom ομοιάζουν με την τεχνική των αρχείων υπογραφών (signature files).



Επανάληψη: Αρχεία Υπογραφών (Signature files)

b=3 (3 words per block) **B=6** (bit masks of 6 bits)



Signature Function

$h(\text{text}) = 000101$
 $h(\text{many}) = 110000$
 $h(\text{words}) = 100100$
 $h(\text{made}) = 001100$
 $h(\text{letters}) = 100001$



Bloom filters [Burton Bloom 1970] Συμπαγής Κωδικοποίηση Συνόλων

Ένα σύνολο κωδικοποιείται σε ένα δυαδικό διάνυσμα των m-bits

κ συναρτήσεις κατακερματισμού h_1, h_2, \dots, h_k , με πεδίο τιμών το $\{1, \dots, m\}$

Κωδικοποίηση στοιχείου:

$BF(\{a\}) =$ διάνυσμα με άσσους στις θέσεις $h_1(a), h_2(a), \dots, h_k(a)$

Κωδικοποίηση συνόλου:

$BF(\{a_1, a_2\}) = BF(\{a_1\})$ **BITwiseOR** $BF(\{a_2\})$

Πως βρίσκω αν ένα στοιχείο b ανήκει στο σύνολο A ?

1/ Υπολογίζω το BloomFilter του b

2/ Κοιτάζω αν οι άσσοι του $BF(b)$ υπάρχουν στο $BF(A)$

Αν όχι, τότε σίγουρα το b δεν ανήκει στο A

Αν ναι, τότε ανήκει αλλά μπορεί και να μην ανήκει (false positive)

Όσο μεγαλύτερο είναι το m , τόσο μικρότερη η πιθανότητα για false positives



Bloom filters: Παράδειγμα

$m=14, k=3$



$hash_1("apples") = 3$

$hash_2("apples") = 12$

$hash_3("apples") = 11$

$\{apples\} =$

$hash_1("plums") = 11$

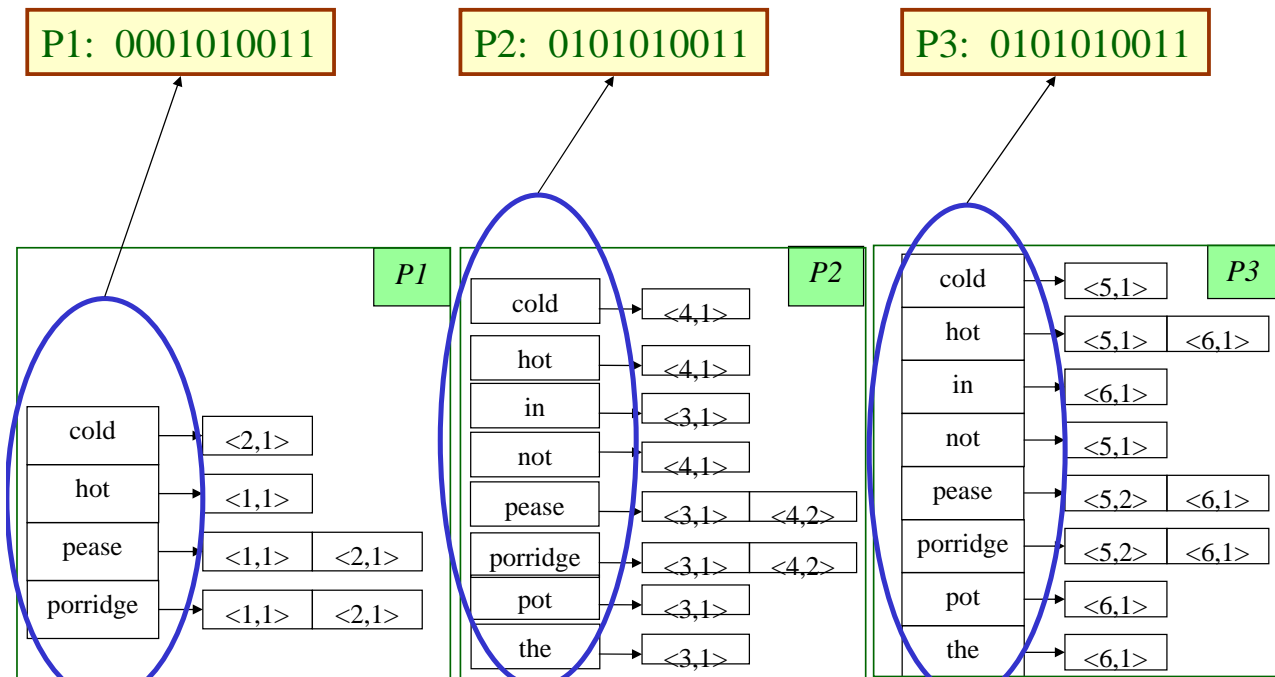
$hash_2("plums") = 1$

$hash_3("plums") = 8$

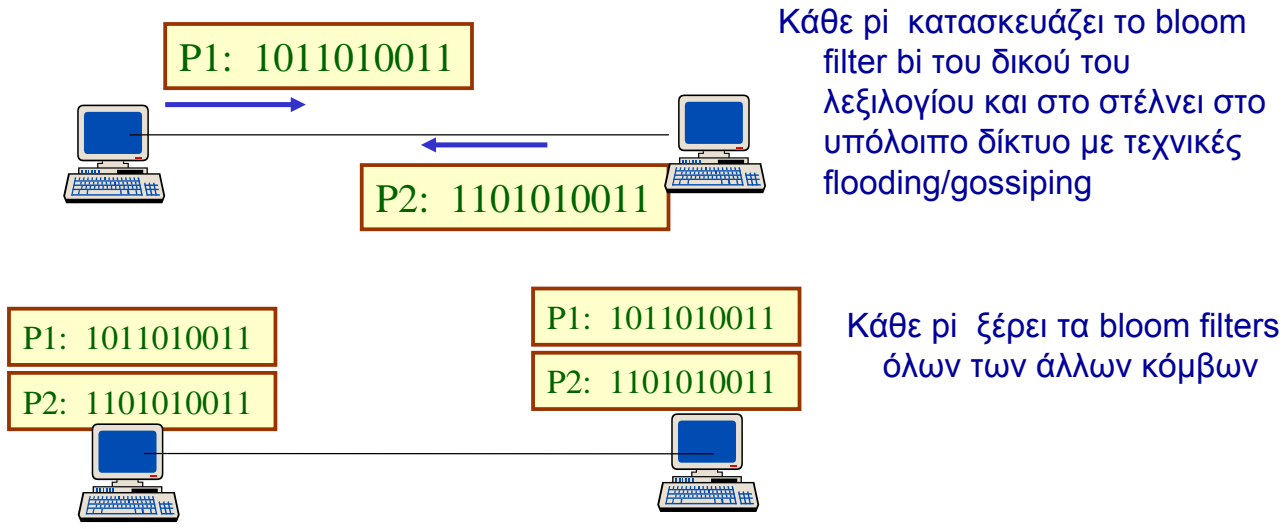
$\{apples, plums\} =$



Περιγραφή των λεξιλογίων με Bloom filters



P2P and IR: Το σύστημα PlanetP



Έτσι κάθε ρι μπορεί να βρεί τους κόμβους που έχουν έναν συγκεκριμένο όρο (άρα μπορεί να προσεγγίσει το καθολικό ευρετήριο)



Bloom filters in PlanetP: Πόσο μεγάλα είναι;

AP89 Collection (Associated Press articles of 1989 from TREC):

84,678 documents, 129,603 words, collection size 266 MB

Num. Peers	Memory used (MB)	% of collection size
10	0.45	0.18%
100	1.79	0.70%
1000	4.48	1.76%

1000 Nodes: => about 4500 terms per peer

Bloom filters with less than 5% false positives =>

Bloom filter size for the vocabulary of one peer: 4.6 KB

Total size of bloom filters of peers : 4.6 MBytes

Γιατί το μέγεθος
αυξάνει με το πλήθος
των κόμβων;



PlanetP: Τρόπος ενημέρωσης των κόμβων (Gossiping algorithms)

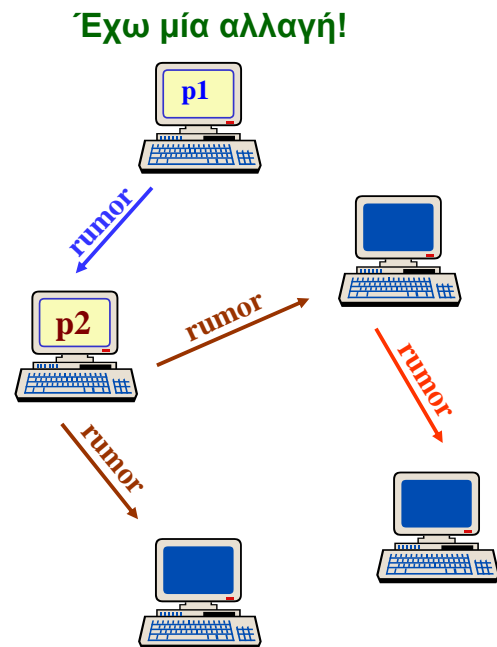
- Η μετάδοση των Bloom filters σε όλο το δίκτυο καθώς και η ενημέρωση των κόμβων (για νέα δεδομένα, είσοδο/έξοδο κόμβων) μπορεί να γίνει με ποικίλους **αλγορίθμους gossiping**:
 - rumoring algorithm
 - anti-entropy algorithm
 - partial anti-entropy algorithm
 -



(Gossiping algorithms) Rumoring (φημολογία)

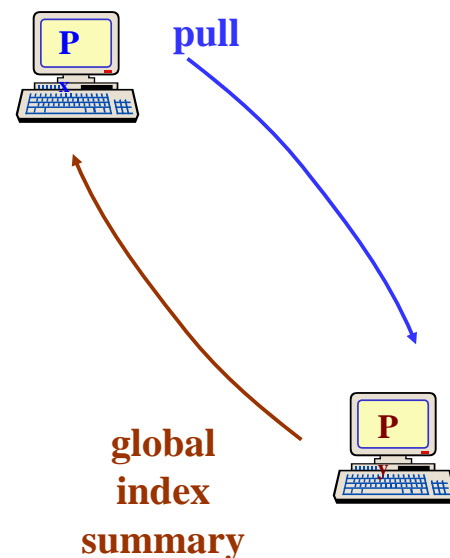
Ο p1 έχει μια αλλαγή:

- κάθε X δευτερόλεπτα, ο p1 στέλνει ένα μήνυμα με την αλλαγή σε έναν τυχαία επιλεγμένο κόμβο p2
- Αν ο p2 δεν ήξερε αυτήν την πληροφορία, τότε αρχίζει να κάνει ό,τι και ο p1
- Ο p1 σταματάει να στέλνει μηνύματα μόνο αν η συνεχόμενοι κόμβοι του πουν ότι ήταν ήδη ενημερωμένοι της αλλαγής.



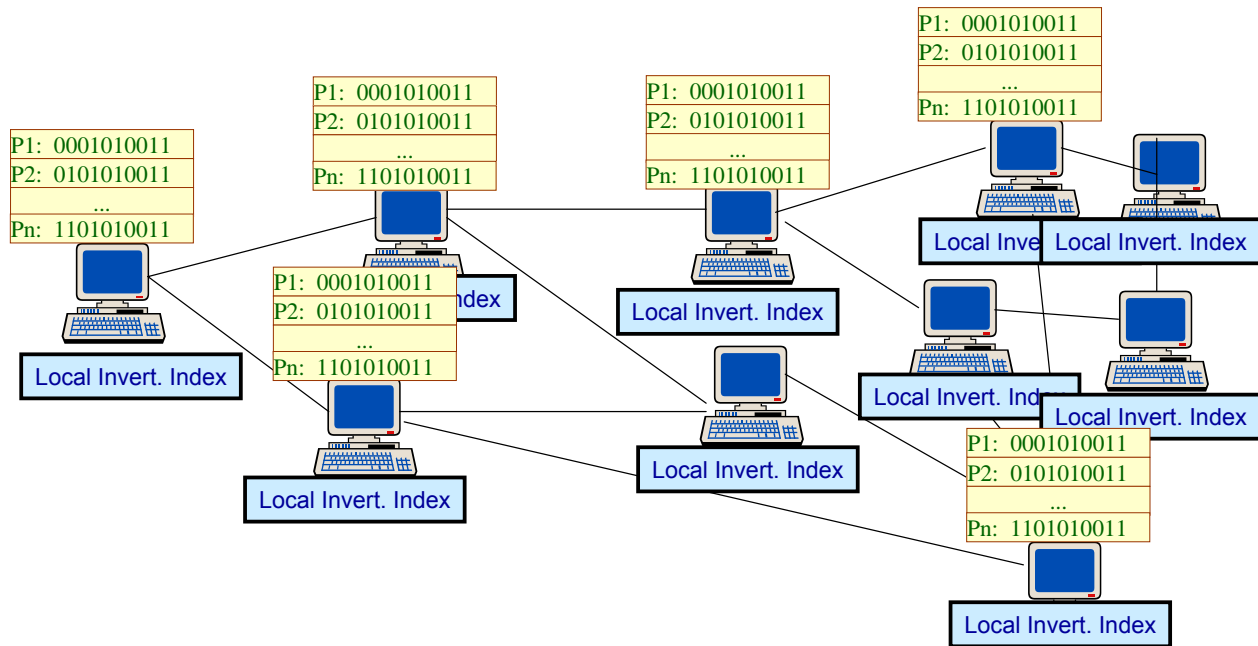
(Gossiping algorithms) anti-entropy

- Κάθε X δευτερόλεπτα, κάθε κόμβος επιλέγει τυχαία έναν άλλο κόμβο (από το καθολικό του ευρετήριο) και του ζητάει να του στείλει μια περίληψη το δικού του καθολικού ευρετηρίου.
- Αν διαπιστώσει ότι δεν είναι ενημερωμένος, του ζητάει ό,τι χρειάζεται.
- Purpose: The algorithm allows to avoid the possibility of rumors dying out before reaching everyone

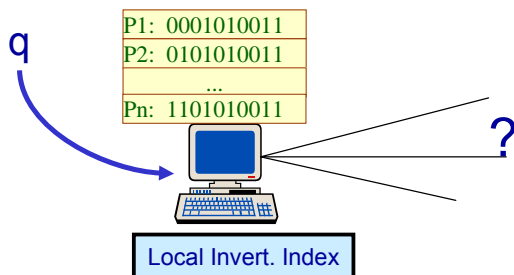




P2P and IR: Το σύστημα PlanetP



PlanetP: Επιλογή Κόμβου



- 1/ **Βαθμολόγηση κόμβων** βάσει της πιθανότητας να έχουν έγγραφα συναφή με την q
- 2/ **Επιλογή** των κόμβων που θα ερωτηθούν και ενοποίηση των αποτελεσμάτων που θα επιστρέψουν

Inverse Peer Frequency (IPF) of a term t =

$$\text{IPF}(t) := |\text{total number of peers}| / |\text{peers that contain the term } t|$$

$$\text{Score}(p_j, q) = \sum \{ \text{IPF}(t) \mid t \in q, t \in \text{Bfilter}(p_j) \}$$

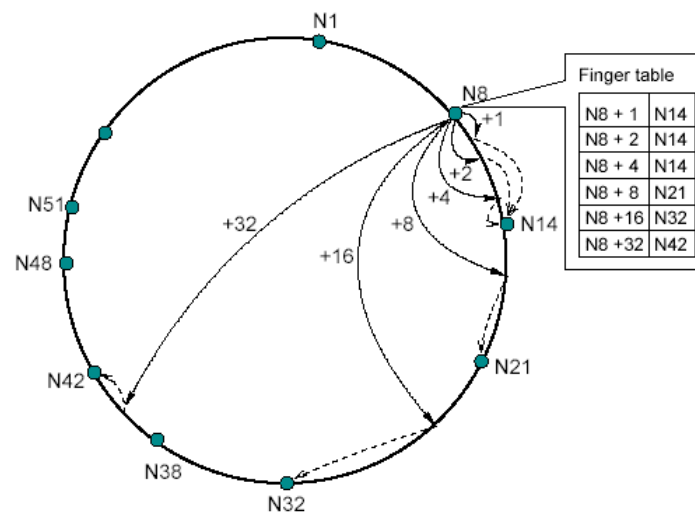


PlanetP: Αποτελεσματικότητα & Επιδόσεις

- Η αποτελεσματικότητα προσεγγίζει αυτήν που θα είχαμε αν κάθε κόμβος είχε ολόκληρο το ευρετήριο
- Τα μηνύματα φτάνουν σε 20%-40% περισσότερους κόμβους σε σχέση με την περίπτωση όπου κάθε κόμβος γνώριζε ακριβώς το καθολικό ευρετήριο
- Gossiping rate 1/second => PlanetP can propagate a Bloom filter containing 1000 terms in less than 40 secs for a community of 1000 peers. This requires an average of 24KB/s per peer.



Ανάκτηση Πληροφοριών σε Δομημένα Ομ. Συσ/τα (Chord-style)



- Ποια είναι εδώ τα κλειδιά ?



Ανάκτηση Πληροφοριών σε Δομημένα Ομ. Συσ/τα (Chord-style)

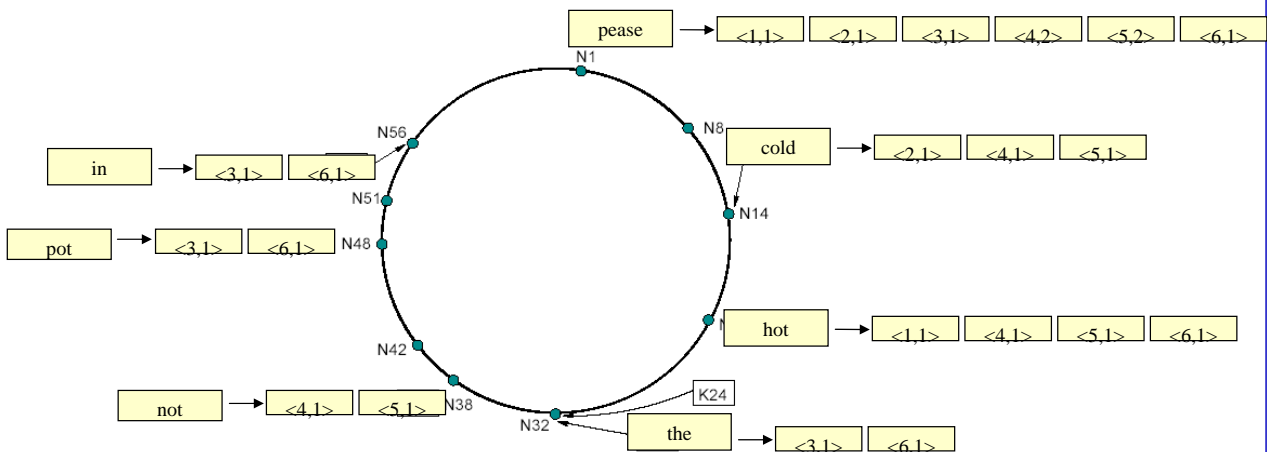
Περίπτωση (I): Κάθε όρος είναι ένα κλειδί

- Το ευρετήριο κατανέμεται βάσει των όρων
 - (άρα έχουμε term-partitioning: θυμηθείτε την παράλληλη Α.Π.)
- Αδυναμία: Η ενημέρωση των ευρετηρίων είναι ακριβή:
 - Εισαγωγή ενός νέου εγγράφου:
 - Για κάθε λέξη του εγγράφου, πρέπει να βρούμε τον κόμβο που είναι υπεύθυνος για αυτήν την λέξη και να του στείλουμε την ανεστραμμένη λίστα

P1	cold	→	<2,1>	<4,1>	<5,1>			
	hot	→	<1,1>	<4,1>	<5,1>	<6,1>		
	in	→	<3,1>	<6,1>				
P2	not	→	<4,1>	<5,1>				
	pease	→	<1,1>	<2,1>	<3,1>	<4,2>	<5,2>	<6,1>
	porridge	→	<1,1>	<2,1>	<3,1>	<4,2>	<5,2>	<6,1>
P3	pot	→	<3,1>	<6,1>				
	the	→	<3,1>	<6,1>				



Ανάκτηση Πληροφοριών σε Δομημένα Ομ. Συσ/τα (Chord-style): Κάθε όρος είναι ένα κλειδί





Ανάκτηση Πληροφοριών σε Δομημένα Ομ. Συσ/τα (Chord-style)

- Αποτίμηση επερωτήσης q
 - βρίσκουμε κάθε κόμβο που έχει τουλάχιστον έναν όρο του q (χρησιμοποιώντας τους πίνακες δρομολόγησης)
 - Σενάριο 1: κάθε ένας από αυτούς τους κόμβους υπολογίζει τα μερικά σκορ και τα στέλνει στον ερωτώντα (αφού του στείλουμε και την επερωτήση)
 - Σενάριο 2: κάθε ένας από αυτούς τους κόμβους επιστρέφει τις ανεστραμμένες λίστες
- [-] Ανταλλαγή πολλών μηνυμάτων για επερωτήσεις με πολλούς όρους

P1	cold	→	<2,1>	<4,1>	<5,1>			
	hot	→	<1,1>	<4,1>	<5,1>	<6,1>		
	in	→	<3,1>	<6,1>				
P2	not	→	<4,1>	<5,1>				
	pease	→	<1,1>	<2,1>	<3,1>	<4,2>	<5,2>	<6,1>
	porridge	→	<1,1>	<2,1>	<3,1>	<4,2>	<5,2>	<6,1>
P3	pot	→	<3,1>	<6,1>				
	the	→	<3,1>	<6,1>				



Ανάκτηση Πληροφοριών σε Δομημένα Ομ. Συσ/τα (Chord-style)

Υπόθεση: Έστω ότι το σύστημα λαμβάνει πολύ συχνά επερωτήσεις με 2 όρους

Περίπτωση (II): Θεωρούμε ως κλειδί κάθε ζευγάρι όρων

- Αν η επερωτήση έχει 2 όρους, τότε ένας μόνο κόμβος θα έχει όλο το κομμάτι του ευρετηρίου που χρειαζόμαστε
- Άρα έτσι έχουμε λίγα μηνύματα
- Π.χ. q= Hotels Crete
 - Ξέρω ότι υπάρχει ένας κόμβος που έχει τις ανεστραμμένες λίστες και των δυο όρων, άρα ο κόμβος αυτός μπορεί να αποτιμήσει πλήρως την επερωτήση
- Αδυναμία: $|V|^*$ ($|V|-1$) κλειδιά, άρα η ανεστραμμένη λίστα κάθε λέξης είναι αποθηκευμένη $|V|-1$ φορές

P1	Hotels	[...inverted list for Hotels ...]
	Crete	[...inverted list for Crete ...]
P1	Hotels	[...inverted list for Hotels ...]
	Cefalonia	[...inverted list for Cefalonia ...]
P3	Crete	[...inverted list for Crete ...]
	Cefalonia	[...inverted list for Cefalonia ...]

Η είσοδος ενός νέου εγγράφου είναι ακόμα πιο ακριβή

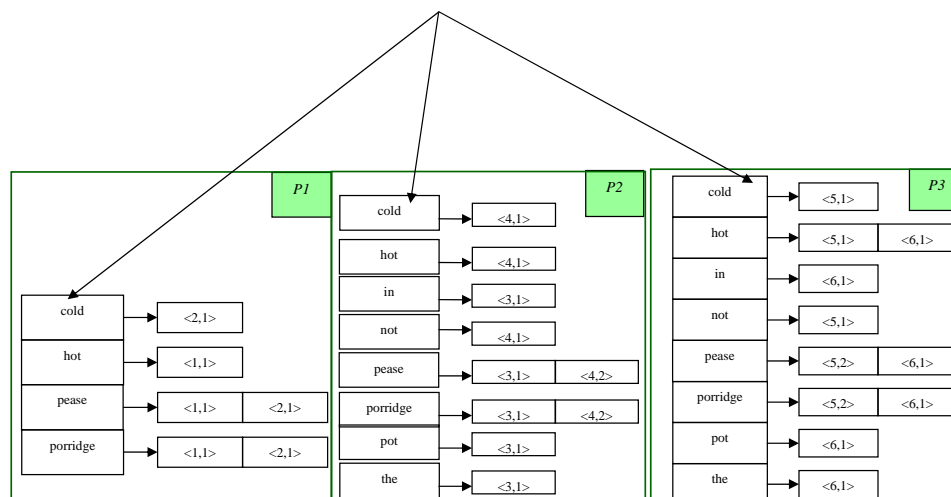


Ανάκτηση Πληροφοριών σε **Δομημένα Ομ. Συσ/τα**

- Περίπτωση (III): Θεωρούμε ως κλειδιά τα διανύσματα των εγγράφων
 - **Ερ**: Ποια προσέγγιση δομημένων συστημάτων είναι κατάλληλη για την παράσταση διανυσμάτων (Chord ή CAN) ;
 - **Απ**: Η προσέγγιση του CAN διότι βλέπει τον χώρο των κλειδιών ως ένα κ-διάστατο χώρο
 - Άρα διαμερίζουμε τα έγγραφα στους κόμβους βάσει των διανυσμάτων τους.
 - (άρα document-partitioning (θυμηθείτε την Παράλληλη Α.Π.))
 - **Ερ**: Τι κερδίζουμε διαμερίζοντας τα έγγραφα όπως το CAN ?
 - **Απ**: Τα κοντινά (ως προς το μέτρο συνημίτονου) έγγραφα τοποθετούνται στον ίδιο ή σε κοντινούς κόμβους.



Document Partitioning: Ο υπολογισμός των καθολικών στατιστικών (IDF) απαιτεί επικοινωνία



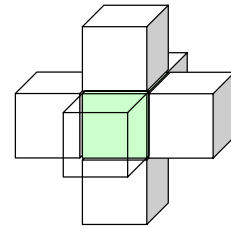
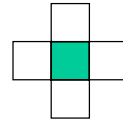


Ανάκτηση Πληροφοριών σε **Δομημένα Ομ. Συσ/τα** (**CAN**-style)

- **Ερ**: Πόσες διαστάσεις έχουν τα διανύσματα των εγγράφων;
- **Απ**: Συνήθως πολλές (π.χ. 10.000)

- **Ερ**: Πόσους γείτονες έχει μια περιοχή k -διάστατου χώρου;
- **Απ**: κατά μέσο όρο $2k$

- Για $k=1$ έχω 2
- Για $k=2$ έχω 4
- Για $k=3$ έχω 6
- Για $k=10.000$ έχω 20.000 !



Ανάκτηση Πληροφοριών σε **Δομημένα Ομ. Συσ/τα** (**CAN**-style). Το σύστημα **pSearch**

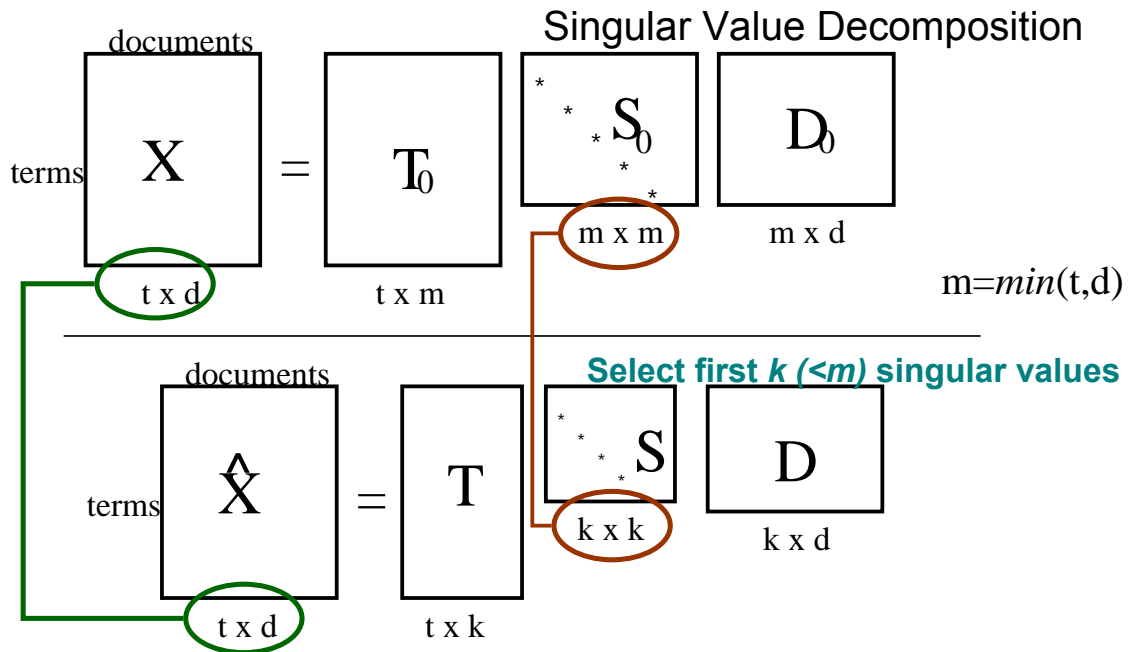
- Μείωση των διαστάσεων των διανυσμάτων για
 - (I) Μείωση του αριθμού των γειτόνων που πρέπει να γνωρίζει (αποθηκεύει) ένας κόμβος.
 - (II) Ομαδοποίηση εγγράφων
 - Αξιοποίηση συνωνύμων, συνεμφανιζόμενων λέξεων, μείωση θορύβου
- Τρόπος μείωσης διαστάσεων: **Latent Semantic Indexing**



Επανάληψη: Latent Semantic Indexing:

t: total number of index terms

d: total number of documents



Επανάληψη LSI: Paper example

Index terms in italics

Titles:

- c1: *Human machine interface* for Lab ABC computer applications
- c2: *A survey of user opinion of computer system response time*
- c3: *The EPS user interface management system*
- c4: *System and human system engineering testing of EPS*
- c5: *Relation of user-perceived response time to error measurement*

- m1: *The generation of random, binary, unordered trees*
- m2: *The intersection graph of paths in trees*
- m3: *Graph minors IV: Widths of trees and well-quasi-ordering*
- m4: *Graph minors: A survey*



Επανάληψη LSI: SVD with minor terms dropped

$$\begin{matrix} T & S & D' \\ \left[\begin{array}{cc} 0.22 & -0.11 \\ 0.20 & -0.07 \\ 0.24 & 0.04 \\ 0.40 & 0.06 \\ 0.64 & -0.17 \\ 0.27 & 0.11 \\ 0.27 & 0.11 \\ 0.30 & -0.14 \\ 0.21 & 0.27 \\ 0.01 & 0.49 \\ 0.04 & 0.62 \\ 0.03 & 0.45 \end{array} \right] & \left[\begin{array}{cc} 3.34 & \\ & 2.54 \end{array} \right] & \left[\begin{array}{cccccccccc} 0.20 & 0.61 & 0.46 & 0.54 & 0.28 & 0.00 & 0.02 & 0.02 & 0.08 \\ -0.06 & 0.17 & -0.13 & -0.23 & 0.11 & 0.19 & 0.44 & 0.62 & 0.53 \end{array} \right] \end{matrix}$$

TS define
coordinates for
documents in latent
space

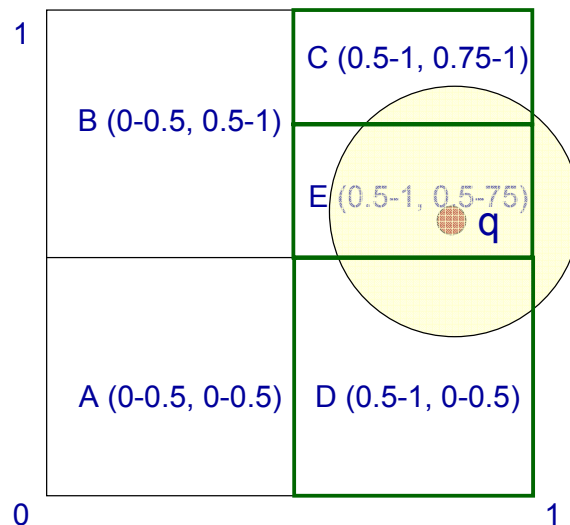


Ανάκτηση Πληροφοριών σε **Δομημένα Ομ. Συσ/τα** (CAN-style). Το σύστημα **pSearch**

- Διάσταση διανυσμάτων μετά την εφαρμογή LSI: **50-350**
- Φτιάχνουμε ένα CAN με διαστάσεις όσες των διανυσμάτων (μετά το LSI).
- Εισαγωγή ενός νέου εγγράφου:
 - Φτιάχνεται το «semantic διάνυσμα» του εγγράφου (βάσει των διαστάσεων που προέκυψαν από την εφαρμογή του LSI) και εισάγεται στον κατάλληλο κόμβο
- Είσοδος μιας νέας επερώτησης
 - Φτιάχνεται το semantic διάνυσμα της επερώτησης και δρομολογείται στον κατάλληλο κόμβο
 - Μόλις φτάσει στον κόμβο, διαδίδεται στους γείτονες σε απόσταση ρ
 - Το ρ μπορεί να δίδεται μαζί με την αρχική επερώτηση



- Επερώτηση

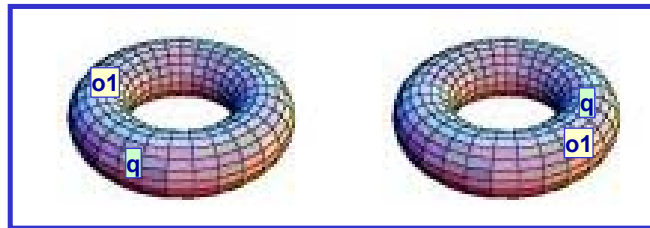


- Ο υπολογισμός του LSI απαιτεί Καθολικά στατιστικά (IDF)
- Επίσης όλοι οι κόμβοι πρέπει να γνωρίζουν την βάση του σημασιολογικού χώρου (για να υπολογίζουμε τα σημασιολογικά διανύσματα των νέων εγγράφων).
- Τα παραπάνω πρέπει να διαδοθούν σε όλους τους κόμβους.
- Το πρόβλημα των διαστάσεων
 - 300 LSI διαστάσεις. Αν έχω λίγους κόμβους τότε η πραγματική διάσταση του CAN είναι μικρότερη γιατί δεν υπάρχουν αρκετοί κόμβοι. Έτσι πολλές διαστάσεις παραμένουν αδιαμέριστες, μεγαλώνοντας έτσι το μήκος του μονοπατιού αναζήτησης.



CAN & Multiple Realities

- Ένας τρόπος αύξησης της ευρωστίας / ανθεκτικότητας είναι να θεωρήσουμε Πολλαπλές Πραγματικότητες (Multiple Realities)
 - Δεν έχουμε 1 αλλά m διαφορετικά συστήματα συντεταγμένων
 - Κάθε κόμβος έχει μια ζώνη για κάθε σύστημα συντεταγμένων
 - Έτσι έχουμε m αντίγραφα ευρετηρίου
 - Μείωση του μήκους του μονοπατιού αναζήτησης (επιλέγεται το σύστημα συντεταγμένων βάσει του οποίου η αναζητούμενη ζώνη είναι εγγύτερα)



Ανάκτηση Πληροφοριών σε **Δομημένα Ομ. Συσ/τα** (CAN-style). Το σύστημα **pSearch**

- Διαμερισμός των διανυσμάτων σε πολλά διανύσματα μικρότερης διάστασης
 - $(x_1, \dots, x_n) \Rightarrow (x_1, \dots, x_{n_1}), (x_{n_1+1}, \dots, x_{n_2}), (x_{n_2+1}, \dots, x_n)$
 - Τα πρώτα διανύσματα αποθηκεύονται σε ένα CAN1
 - Τα δεύτερα σε ένα CAN2, κ.ο.κ
 - Το διάνυσμα μιας επερώτησης επίσης διαμερίζεται σε διανύσματα μικρότερης διάστασης
 - :

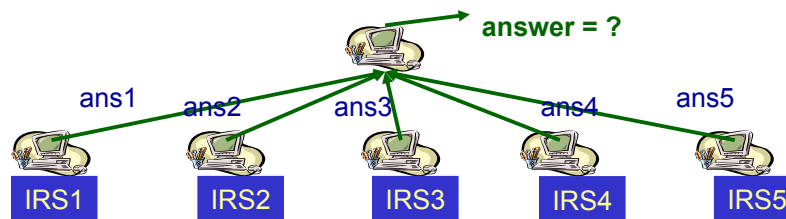


Ανάκτηση Πληροφοριών σε **Δομημένα Ομ. Συσ/τα** (CAN-style). Το σύστημα **pSearch**: Σύνοψη

- Φτιάχνουμε ένα ευρετήριο όπου κάθε έγγραφο δεν περιγράφεται από το διάνυσμα του, αλλά από το διάνυσμα που προκύπτει αν πρώτα εφαρμόσουμε **Latent Semantic Indexing**
 - διανύσματα μικρότερης διάστασης, ομαδοποίηση εγγράφων
- Τα ευρετήριο αυτό διανέμεται στους κόμβους. Το κλειδί του κάθε εγγράφου είναι το διάνυσμα του (μετά την εφαρμογή του LSI). // Αυτό θα τοποθετήσει στον ίδιο κόμβο εννοιολογικά συναφή έγγραφα
- Ο υπολογισμός των διανυσμάτων απαιτεί καθολικά στατιστικά (άρα υπάρχει ανάγκη επικοινωνίας). Επίσης πρέπει να συμφωνηθεί η βάση των διανυσμάτων.
- Μπορεί να χρησιμοποιηθεί και για πολυμέσα (θυμηθείτε **Feature-based Multimedia Indexing**).



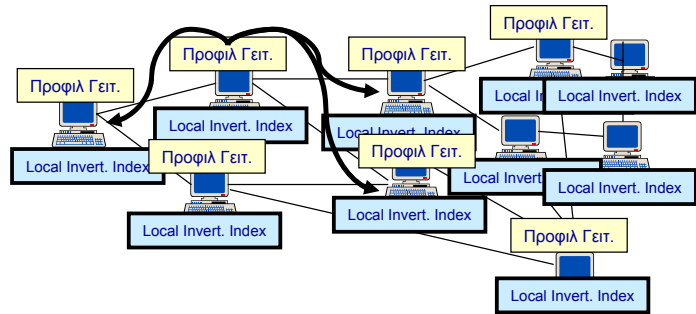
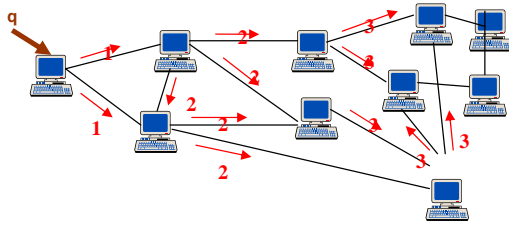
Ενοποίηση Αποτελεσμάτων & Ομότιμα Συστήματα



- **Τεχνικές Ενοποίησης Αποτελεσμάτων**
 - **Round Robin Inter-leaving**
 - **Score-based** (~ merge sort)
 - καλή αν τα σκορ υπολογίζονται βάσει των καθολικών στατιστικών
 - **Weighted-score based**
 - Έστω d_i προερχόμενο από μια πηγή S_j
 - $score(d_i) = score(S_j, d_i) * score(S_j)$
 - Λαμβάνοντας υπόψη μόνο τις διατάξεις και όχι τα σκορ (ενοποίηση διατάξεων)
 - **Borda, Condorcet, Kemeny, Arrow's Impossibility Theorem**



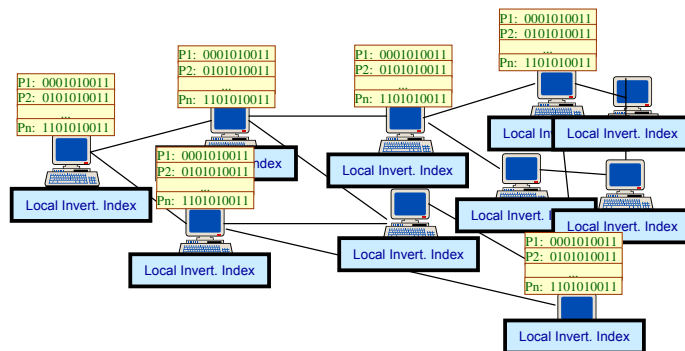
Ενοποίηση Αποτελεσμάτων σε Ομότιμα Συστήματα (I): Οι κόμβοι δεν έχουν στη διάθεση τους **καθολικά** στατιστικά



- **Gnutella-like** systems (document-partitioning):
 - Ενοποίηση: Round-robin interleaving, Score-based, Rank-Aggregation
- Συστήματα βασισμένα σε **προφίλ γειτόνων και >RES**
 - Ενοποίηση: Weighted score-based



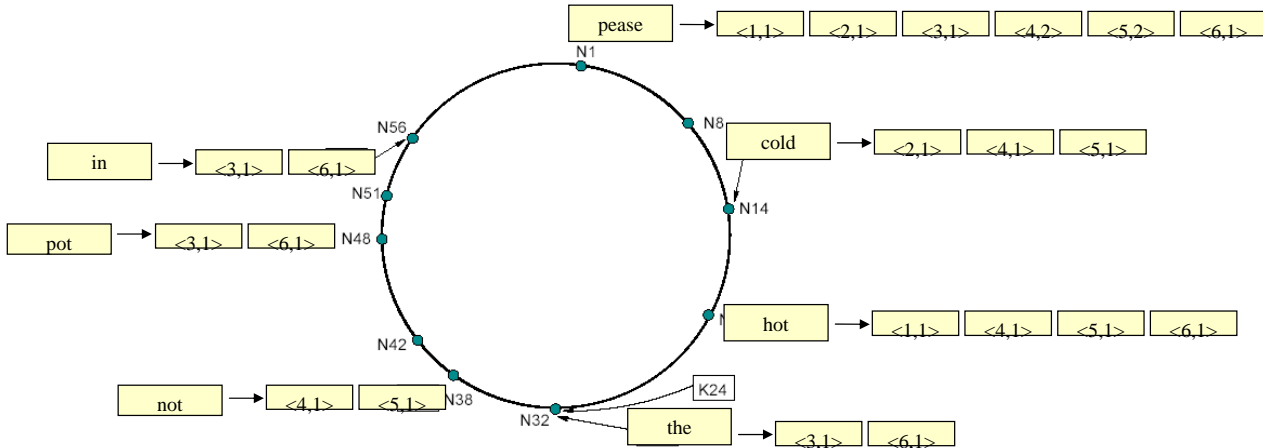
Ενοποίηση Αποτελεσμάτων σε Ομότιμα Συστήματα (II): Οι κόμβοι μπορούν να προσεγγίσουν τα **καθολικά** στατιστικά



- Π.χ. **PlanetP** (κάθε κόμβος μπορεί να προσεγγίσει το καθολικό ευρετήριο)
 - Ενοποίηση: Weighted score-based
 - (καλύτερο από το προφίλ γειτόνων, λιγότερα μηνύματα)



Ενοποίηση Αποτελεσμάτων σε Ομότιμα Συστήματα (III): Οι κόμβοι έχουν στη διάθεση τους τα **καθολικά** στατιστικά



- Π.χ. **Chord-like (term-partitioning)**
 - ο κόμβος που είναι υπεύθυνος για έναν όρο γνωρίζει τις συχνότητες εμφάνισης του καθώς και το πλήθος των κόμβων που έχουν έγγραφα που περιέχουν αυτόν τον όρο
 - Ενοποίηση: απλό **Score-based** είναι μια χαρά
 - κάθε κόμβος υπολογίζει **partial scores**, ο ερωτών τα αθροίζει και παράγει την τελική διάταξη



Ενοποίηση Αποτελεσμάτων σε Ομότιμα Συστήματα (III): Οι κόμβοι έχουν στη διάθεση τους τα **καθολικά** στατιστικά

- Έστω σύστημα όπως το Chord, στο οποίο τα κλειδιά είναι οι όροι και το οποίο συνολικά έχει 100.000 έγγραφα
- Η ανεστραμμένη λίστα ενός όρου έχει το πολύ 100.000 αναφορές σε έγγραφα (έστω ότι κατά μέσο όρο έχει 10.000 αναφορές)
- Έστω ότι ο p λαμβάνει επερώτηση q με 5 όρους. Κάθε όρος της q (μαζί με βάρος του στο q) θα προωθηθεί στον υπεύθυνο κόμβο για τον όρο αυτό
- Κάθε ένας από τους 5 κόμβους θα διατάξει τα έγγραφα βάσει του όρου αυτού και θα επιστρέψει μια λίστα μερικών αποτελεσμάτων
 - το πολύ 100.000 τριάδες (p , docId, score) κατά μέσο όρο 10.000
- Ο p θα λάβει αυτές τις 5 λίστες και θα αθροίσει τα μερικά σκορ
 - $score(doc_i) = score_1(doc_i) + \dots + score_5(doc_i)$
- Άρα $5 * 10.000$ τριάδες ακεραίων πρέπει να μεταφερθούν στο δίκτυο
 - $TotalBytes = 50K * 3 * 4 = 600 KB$
- Ερώτηση: **Αν ο p θέλει να βρει μόνο τα κορυφαία k (π.χ. $k=10$) έγγραφα. Πως μπορούμε να ελαχιστοποιήσουμε την πληροφορία που πρέπει να μεταφέρουμε;**



Top-k Rank Aggregation

- Έχουμε **N** αντικείμενα και τους **βαθμούς** τους βάσει **m** διαφορετικών **κριτηρίων**.
- Έχουμε έναν τρόπο να συνδυάζουμε τα m σκορ κάθε αντικειμένου σε ένα ενοποιημένο σκορ
 - π.χ. min, avg, sum
- Στόχος: Βρες τα **k** αντικείμενα με το υψηλότερο ενοποιημένο σκορ.

Εφαρμογές:

- Υπολογισμός των κορυφαίων-k στοιχείων της απάντησης
 - ενός ΣΑΠ που βασίζεται στο διανυσματικό μοντέλο (τα m κριτήρια είναι οι m όροι της επερώτησης)
 - ενός μεσίτη πάνω από m Συστήματα Ανάκτησης Πληροφοριών
 - μιας επερώτησης σε μια Βάση Πολυμέσων
 - κριτήρια: χρώμα, μορφή, υφή, ...



Άλλο ένα παράδειγμα εφαρμογής

- **Ενοποίηση απαντήσεων σε Μεσολαβητές (middleware)**
 - έστω μια υπηρεσία εύρεσης εστιατορίων βάσει τριών κριτηρίων:
 - τιμή γεύματος
 - απόσταση από ένα σημείο
 - κατάταξη εστιατορίου
 - όπου ο χρήστης μπορεί να ορίσει τον επιθυμητό τρόπο υπολογισμού του ενοποιημένου σκορ ενός εστιατορίου
 - π.χ. $\text{Σκορ} = \text{Τιμή} * 0.5 + \text{Stars} * 0.25 + 0.25 * \text{DistanceFromHome}$
 - η υπηρεσία αυτή υλοποιείται με χρήση τριών απομακρυσμένων υπηρεσιών
 - (a) `getRestaurantsByPrice`
 - (b) `getRestaurantsByStars`
 - (c) `getRestaurantsByDistance`
 - **Πως μπορώ να ελαχιστοποιήσω το πλήθος των στοιχείων που πρέπει να διαβάσω από την απάντηση της κάθε υπηρεσίας, προκειμένου να βρω τα κορυφαία 5 εστιατόρια;**



Εύρεση των κ-κορυφαίων Απλοϊκός Αλγόριθμος

- 1/ Ανέκτησε ολόκληρες τις m λίστες
- 2/ Υπολόγισε το ενοποιημένο σκορ του κάθε αντικειμένου
- 3/ Ταξιλόγησε τα αντικείμενα βάσει του σκορ και επέλεξε τα πρώτα k

Παρατηρήσεις

- Κόστος γραμμικό ως προς το μήκος των λιστών
- Δεν αξιοποιεί το γεγονός ότι οι λίστες είναι ταξινομημένες



Εύρεση των κ-κορυφαίων Παράδειγμα: Απλοϊκός Τρόπος

$S_1 = \langle \mathbf{A} 0.9, \mathbf{C} 0.8, \mathbf{E} 0.7, \mathbf{B} 0.5, \mathbf{F} 0.5, \mathbf{G} 0.5, \mathbf{H} 0.5 \rangle$

$S_2 = \langle \mathbf{B} 1.0, \mathbf{E} 0.8, \mathbf{F} 0.7, \mathbf{A} 0.7, \mathbf{C} 0.5, \mathbf{H} 0.5, \mathbf{G} 0.5 \rangle$

$S_3 = \langle \mathbf{A} 0.8, \mathbf{C} 0.8, \mathbf{E} 0.7, \mathbf{B} 0.5, \mathbf{F} 0.5, \mathbf{G} 0.5, \mathbf{H} 0.5 \rangle$

Ο Απλοϊκός Τρόπος

$$\text{Score}(\mathbf{A}) = 0.9 + 0.7 + 0.8 = 2.4$$

$$\text{Score}(\mathbf{B}) = 0.5 + 1.0 + 0.5 = 2$$

$$\text{Score}(\mathbf{C}) = 0.8 + 0.5 + 0.8 = 2.1$$

$$\text{Score}(\mathbf{E}) = 0.7 + 0.8 + 0.7 = 2.2$$

$$\text{Score}(\mathbf{F}) = 0.5 + 0.7 + 0.5 = 1.7$$

$$\text{Score}(\mathbf{G}) = 0.5 + 0.5 + 0.5 = 1.5$$

$$\text{Score}(\mathbf{H}) = 0.5 + 0.5 + 0.5 = 1.5$$

Τελική διάταξη: $\langle \mathbf{A}, \mathbf{E}, \mathbf{C}, \mathbf{B}, \mathbf{F}, \mathbf{G}, \mathbf{H} \rangle$



Εύρεση των κ-κορυφαίων Πιο Αποδοτικοί Αλγόριθμοι

- Γενική ιδέα: **Αρχισε να διαβάζεις τις διατάξεις από την κορυφή. Προσπάθησε να καταλάβεις πότε πρέπει να σταματήσεις.**
- Αλγόριθμοι
 - **Fagin Algorithm (FA)** [Fagin 1999, J. CSS 58]
 - **Threshold Algorithm (TA)** [Fagin et al., PODS'2001]

Υποθέσεις

- Υποθέτουμε ότι έχουμε στη διάθεση μας 2 τρόπους πρόσβασης στα αποτελέσματα μιας πηγής:
 - **Σειριακή πρόσβαση** στις διατάξεις: φθίνουσα ως προς το σκορ
 - **Τυχαία προσπέλαση**: Δυνατότητα εύρεσης του σκορ ενός αντικειμένου με μία πρόσβαση
- Συναρτήσεις βαθμολόγησης (σκορ)
 - Τα σκορ ανήκουν στο διάστημα $[0,1]$
 - Η συνάρτηση ενοποιημένου σκορ είναι **μονότονη**
 - αν όλα (m) τα σκορ ενός αντικειμένου A είναι μεγαλύτερα ή ίσα των αντίστοιχων σκορ ενός αντικειμένου B , τότε σίγουρα το ενοποιημένο σκορ του A είναι μεγαλύτερο ή ίσο του σκορ του B



Εύρεση των κ-κορυφαίων Ο Αλγόριθμος του Fagin (FA) [1999]

- 1.α/ Κάνε σειριακή ανάκτηση αντικειμένων από κάθε λίστα (αρχίζοντας από την κορυφή), έως ότου η τομή των αντικειμένων από κάθε λίστα να έχει κ αντικείμενα
- 1.β/ Για κάθε αντικείμενο που ανακτήθηκε (στο 1.α) συνέλεξε τα σκορ που λείπουν (με χρήση του μηχανισμού τυχαίας προσπέλασης)
- 2/ Υπολόγισε το ενοποιημένο σκορ του κάθε αντικειμένου
- 3/ Ταξινόμησε τα αντικείμενα βάσει του ενοποιημένου σκορ και επέλεξε τα πρώτα κ

Σχόλια

Αξιοποιεί (α) το γεγονός ότι οι λίστες είναι ταξινομημένες και (β) ότι η συνάρτηση ενοποίησης είναι μονότονη

[–] Το πλήθος των αντικειμένων που θα ανακτηθούν μπορεί να είναι μεγάλο



Εύρεση των κ-κορυφαίων Παράδειγμα: Αλγόριθμος του Fagin (FA)

S1 = < **A** 0.9, **C** 0.8, **E** 0.7, **B** 0.5, **F** 0.5, **G** 0.5, **H** 0.5 >
S2 = < **B** 1.0, **E** 0.8, **F** 0.7, **A** 0.7, **C** 0.5, **H** 0.5, **G** 0.5 >
S3 = < **A** 0.8, **C** 0.8, **E** 0.7, **B** 0.5, **F** 0.5, **G** 0.5, **H** 0.5 >

Έστω ότι θέλω το Top-1

Το E εμφανίζεται σε όλες

(μονοτονία => δεν μπορεί κάποιο δεξιότερο του E να είναι καλύτερο του E

Το E δεν είναι σίγουρα ο νικητής.

Υποψήφιοι νικητές = {A, B, C, E, F}. Κάνουμε τυχαίες προσπελάσεις για να βρούμε τα σκορ που μας λείπουν

getScore(S2,A), getScore(S1,B), getScore(S3,B), getScore(S2,C), ...

Πράγματι, top-1= {A}



Εύρεση των κ-κορυφαίων **Ο Αλγόριθμος TA** (Threshold Algorithm) [Fagin et al. 2001]

Ιδέα:

Υπολόγισε το μέγιστο σκορ που μπορεί να έχει ένα αντικείμενο που δεν έχουμε συναντήσει ακόμα.

- 1/ Κάνε σειριακή ανάκτηση αντικειμένων από κάθε λίστα (αρχίζοντας από την κορυφή) και με χρήση τυχαίας προσπέλασης βρες όλα τα σκορ κάθε αντικειμένου
- 2/ Ταξινόμησε τα αντικείμενα (βάσει του ενοποιημένου σκορ) και κράτησε τα καλύτερα κ
- 3/ Σταμάτησε την σειριακή ανάκτηση όταν τα σκορ των παραπάνω κ αντικειμένων δεν μπορεί να είναι μικρότερα του μέγιστου πιθανού σκορ των απαραίτητων αντικειμένων (threshold).



Εύρεση των κ-κορυφαίων
 Παράδειγμα: Αλγόριθμος του Fagin (FA)

$S1 = \langle \mathbf{A} 0.9, \mathbf{C} 0.8, \mathbf{E} 0.7, \mathbf{B} 0.5, \mathbf{F} 0.5, \mathbf{G} 0.5, \mathbf{H} 0.5 \rangle$
 $S2 = \langle \mathbf{B} 1.0, \mathbf{E} 0.8, \mathbf{F} 0.7, \mathbf{A} 0.7, \mathbf{C} 0.5, \mathbf{H} 0.5, \mathbf{G} 0.5 \rangle$
 $S3 = \langle \mathbf{A} 0.8, \mathbf{C} 0.8, \mathbf{E} 0.7, \mathbf{B} 0.5, \mathbf{F} 0.5, \mathbf{G} 0.5, \mathbf{H} 0.5 \rangle$

Έστω ότι θέλω το Top-2

Το E, B (και το A) εμφανίζονται σε όλες
 (μονοτονία => δεν μπορεί κάποιο δεξιότερο του B να είναι καλύτερο του B)



Εύρεση των κ-κορυφαίων
 Παράδειγμα: Αλγόριθμος TA:

$S1 = \langle \mathbf{A} 0.9, \mathbf{C} 0.8, \mathbf{E} 0.7, \mathbf{B} 0.5, \mathbf{F} 0.5, \mathbf{G} 0.5, \mathbf{H} 0.5 \rangle$
 $S2 = \langle \mathbf{B} 1.0, \mathbf{E} 0.8, \mathbf{F} 0.7, \mathbf{A} 0.7, \mathbf{C} 0.5, \mathbf{H} 0.5, \mathbf{G} 0.5 \rangle$
 $S3 = \langle \mathbf{A} 0.8, \mathbf{C} 0.8, \mathbf{E} 0.7, \mathbf{B} 0.5, \mathbf{F} 0.5, \mathbf{G} 0.5, \mathbf{H} 0.5 \rangle$

Έστω ότι θέλω το Top-1

$\text{Score(A)} = 0.9 + 0.7 + 0.8 = 2.4$
 $\text{Score(B)} = 0.5 + 1.0 + 0.5 = 2$
 $\text{UpperBound} = 0.9 + 1.0 + 0.8 = 2.7$
 αφού $2.7 > 2.4$ συνεχίζω

$\text{Score(C)} = 0.8 + 0.5 + 0.8 = 2.1$
 $\text{Score(E)} = 0.7 + 0.8 + 0.7 = 2.2$
 $\text{UpperBound} = 0.8 + 0.8 + 0.8 = 2.4$
 αφού 2.4 δεν είναι μεγαλύτερο του 2.4 (σکور του A) σταματάω.



Σύγκριση: Fagin vs. TA

- Ο FA ποτέ δεν τερματίζει ενωρίτερα του TA
- Ο TA χρειάζεται μόνο έναν μικρό (k) ενταμιευτή (buffer)
- Ο TA μπορεί όμως να κάνει περισσότερες τυχαίες προσπελάσεις

Ο TA είναι βέλτιστος για όλες τις μονότονες συναρτήσεις σκορ

- Συγκεκριμένα, είναι "instant optimal": είναι καλύτερος πάντα (όχι μόνο στην χειρότερη περίπτωση ή στην μέση περίπτωση)

• Επεκτάσεις

- Αλγόριθμος NRA (Non Random Access)
 - Έκδοση του TA για την περίπτωση που η τυχαία πρόσβαση είναι αδύνατη. Επίσης "instant optimal".
 - Do sequential access until there are k objects whose lower bound no less than the upper bound of all other objects
 - Αλγόριθμος CA (Combined Algorithm)
 - Έκδοση του TA που θεωρεί τις τυχαίες προσπελάσεις ακριβότερες των σειριακών.



Ενοποίηση Αποτελεσμάτων σε Ομότιμα Συστήματα (III): Οι κόμβοι έχουν στη διάθεση τους τα **καθολικά** στατιστικά

- Έστω σύστημα όπως το Chord, στο οποίο τα κλειδιά είναι οι όροι και το οποίο συνολικά έχει 100.000 έγγραφα
- Η ανεστραμμένη λίστα ενός όρου έχει το πολύ 100.000 αναφορές σε έγγραφα (έστω ότι κατά μέσο όρο έχει 10.000 αναφορές)
- Έστω ότι ο p λαμβάνει επερώτηση q με 5 όρους. Κάθε όρος της q (μαζί με βάρος του στο q) θα προωθηθεί στον υπεύθυνο κόμβο για τον όρο αυτό
- Κάθε ένας από τους 5 κόμβους θα διατάξει τα έγγραφα βάσει του όρου αυτού και θα επιστρέψει μια λίστα μερικών αποτελεσμάτων
 - το πολύ 100.000 τριάδες (p , docId, score) κατά μέσο όρο 10.000
- Ο p θα λάβει αυτές τις 5 λίστες και θα αθροίσει τα μερικά σκορ
 - $score(doc_i) = score_1(doc_1) + \dots + score_5(doc_1)$
- Άρα $5 * 10.000$ τριάδες ακεραίων πρέπει να μεταφερθούν στο δίκτυο
 - $TotalBytes = 50K * 3 * 4 = 600 KB$

• Ερώτηση: **Αν ο p θέλει να βρει μόνο τα κορυφαία k (π.χ. $k=10$) έγγραφα. Πως μπορούμε να ελαχιστοποιήσουμε την πληροφορία που πρέπει να μεταφέρουμε;**



Ομότιμα Συστήματα (P2P) και Ανάκτηση Πληροφοριών

- Διαφορές με Κατανεμημένη Ανάκτηση
 - Napster-style
 - Gnutella-style (local inv. Index)
 - Freenet-style
 - $(p, q, |\text{ans}(q)|), > \text{RES}$
 - $(p, q, |\text{ans}(q)|), > \text{RES} * \text{sim}(q)$
 - Hierarchical
 - PlanetP (Bloom filters)
 - Chold-style: key=1 term, a term pair, ...
 - term partitioning
 - CAN-style: key = LSI vector pSearch (LSI + CAN)
 - document partitioning
 - Result aggregation
- Γενικά: P2P & IR = αντικείμενο έρευνας σήμερα



Αναφορές

- D. Zeinalipour-Yazti, Vana Kalogeraki, Dimitrios Gunopulos, Information Retrieval in P2P Networks
- Text-Based Content Search and Retrieval in *ad hoc* P2P Communities, Francisco Matias Cuenca-Acuna and Thu D. Nguyen
- Jie Lu, Jamie Callan, «Federated Search of Text-Based Digital Libraries in Hierarchical Peer-to-Peer Networks», SIGMOD'04 workshop
- Fagin, Lotem, and Naor, Optimal Aggregation Algorithms for Middleware (PODS 2001)