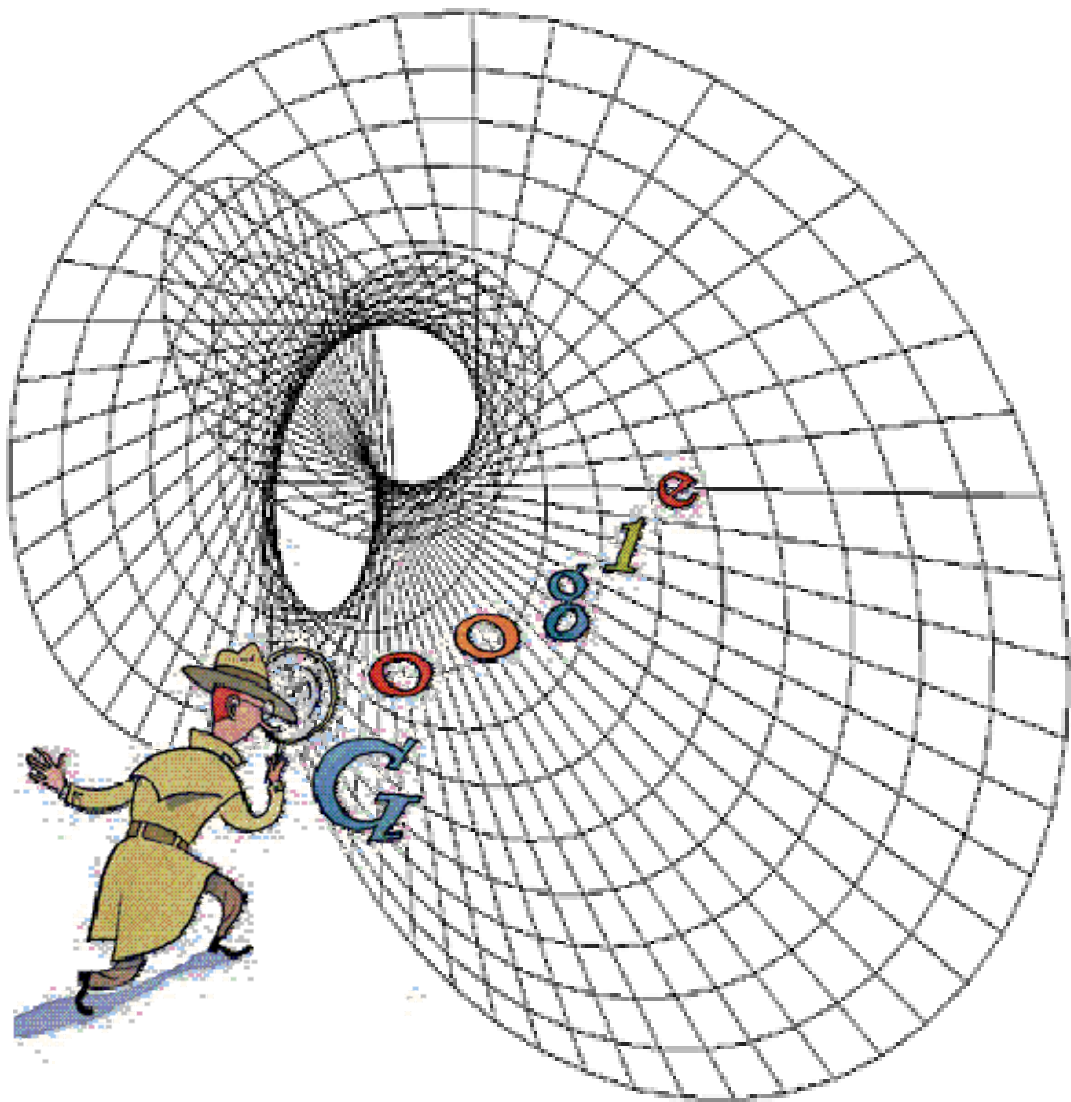


ΕΡΓΑΣΙΑ 1

«Η Μηχανή Αναζήτησης Google»



⁶ Εικόνα εξώφυλλου: Λωρίδα του Möbius. 2005. (Σύνθεση).

Η Μηχανή Αναζήτησης Google⁷

Εισαγωγή



Εικόνα 1.1 Λογότυπο.
Leonardo da Vinci's
Birthday - April 15, 2005.

Η εύρεση χρήσιμων πληροφοριών στον Παγκόσμιο Ιστό (World Wide Web) λόγω της πληθώρας των πληροφοριών γίνεται κατά ένα μεγάλο μέρος με τη χρήση των μηχανών αναζήτησης (search engines). Οι μηχανές αναζήτησης

καλούνται σαν συστήματα ανάκτησης πληροφοριών να επεξεργαστούν τις ανάγκες των χρηστών που εκφράζονται με την μορφή ερωτημάτων (queries).

Μιά μηχανή αναζήτησης, στην βασική της μορφή, αποτελείται από τρία βασικά τμήματα:

- Έναν Crawler, που ξεκινά από ένα αρχικό σύνολο διευθύνσεων (URL) σελίδων Web και συλλέγει σελίδες διατρέχοντας το Web, μέσω των υπερσυνδέσμων (links) του αρχικού συνόλου.
- Έναν Δεικτοδοτητή (Indexer), που επεξεργάζεται τις σελίδες που συνέλεξε ο Crawler για την κατασκευή του ευρετηρίου (Index).
- Έναν επεξεργαστή ερωτήματος (query processor), ο οποίος επεξεργάζεται τα ερωτήματα και επιστρέφει τις σχετικές σελίδες βάσει του μοντέλου ανάκτησης πληροφορίας (Information Retrieval) που χρησιμοποιεί, ταξινομημένες (βάσει ενός αλγόριθμου ταξινόμησης) ή όχι.

Μιά τέτοια μηχανή αναζήτησης είναι και το Google, πιθανότατα η δημοφιλέστερη μηχανή αναζήτησης που υπάρχει αυτή την στιγμή στο διαδίκτυο.



Εικόνα 1.2 Λογότυπο.
Happy Birthday to
Michelangelo - March 6, 2003.

⁷ Βλ. Διαφάνειες 1α & 1β, στο Παράρτημα Ι, σελ. 97-98.

Ιστορική Αναδρομή της Μηχανής Αναζήτησης Google⁸

Το 1995 συναντήθηκαν για πρώτη φορά οι απόφοιτοι του τμήματος υπολογιστών του Πανεπιστημίου Stanford, Larry Page και Sergey Brin. Αρχισαν να προβληματίζονται για ένα βασικό πρόβλημα στο χώρο των υπολογιστών: την ανάκτηση πληροφορίας από ένα μεγάλο όγκο δεδομένων. Ο Ιανουάριος του 1996 βρήκε τον Larry και τον Sergey να συνεργάζονται για μια μηχανή αναζήτησης με το όνομα BackRub. Το όνομα αυτό προερχόταν από την ιδιότητα της μηχανής αυτής να αναλύει "προς τα πίσω" τους συνδέσμους (links) που "έδειχναν" προς κάποιο δεδομένο web site.

Το Σεπτέμβριο του 1998 ξεκίνησε η εταιρεία Google Inc στην Καλιφόρνια των ΗΠΑ. Από τη δοκιμαστική του λειτουργία το Google δεχόταν καθημερινά 10.000 αναζητήσεις κατά μέσο όρο. Η αποδοχή του από τους χρήστες του διαδικτύου ήταν ολοένα και αυξανόμενη. Στις 21 Σεπτεμβρίου του 1999 τέλειωσε η δοκιμαστική φάση λειτουργίας του και το Google έκανε το επίσημο ξεκίνημά του.

Από τότε συνέχισε να αναπτύσσεται τεχνολογικά και να προσελκύει όλο και περισσότερους χρήστες. Ο λόγος ήταν απλός: έβρισκαν αυτό που έψαχναν! Σήμερα εκατομύρια χρήστες το επισκέπτονται και πάνω από 3 δις ιστοσελίδες βρίσκονται στη βάση δεδομένων του, έτοιμες να μας δώσουν τον πλούτο των πληροφοριών τους.



Εικόνα 1.3 Λογότυπο.
*Google celebrated Monet's
birthday on November 14.*

⁸ Η λέξη Google είναι μια παράφραση της λέξης googol με την οποία ο Milton Sirotta, ανηψιός του Αμερικανού Μαθηματικού Edward Kasner, αναφερόταν στον αριθμό 1 ακολουθούμενο με 100 μηδενικά (10 εις στην εκατοστή δύναμη). Ο googol είναι ένας πολύ μεγάλος αριθμός και δεν υπάρχει τίποτα στο Σύμπαν σε αριθμό ίσο με αυτό (π.χ. αστέρια, διαστημική σκόνη κτλ). Ο όρος Google αντικατοπτρίζει το στόχο της μηχανής αναζήτησης να "οργανώσει" τη φαινομενικά "άπειρη" πληροφορία που διατίθεται στο web.

Το Google Σήμερα

Το Google είναι μία ιδιωτική εταιρεία που εστιάζει στην αναζήτηση και ανάκτηση πληροφοριών. Έχει δημιουργήσει το δικό του Web Site στη διεύθυνση www.google.com, μια ισχυρή μηχανή αναζήτησης η οποία είναι γρήγορη, ακριβής και εύκολη στην χρήση της.



Εικόνα 1.4 Λογότυπο.
Google celebrates Vincent van Gogh's Birthday - March 30. 2005.

1. Το Google σε αριθμούς

- Αναζητήσεις που απαντούνται καθημερινά: περισσότερες από 200 εκατομμύρια.
- Σελίδες που αναζητούνται στον ιστό (Indexer): περισσότερες από 3 δισεκατομμύρια. Θεωρείται η μεγαλύτερη βάση δεικτοδοτημένων ιστοσελίδων από όλες τις μηχανές αναζήτησης.
- Τύποι αρχείων που αναζητούνται: πολλοί περιλαμβανομένων των Hyper Text Markup Language (html), Adobe Portable Document format (pdf), Microsoft Excel (xls), Microsoft Word (doc), κ.α.
- Εικόνες (Images): περισσότερες από 425 εκατομμύρια.
- Μηνύματα μέσω δικτύου: περισσότερα από 800 εκατομμύρια.

2. Χρήστες

Το www.google.com είναι ένα από τα 10 πιο δημοφιλέστερα sites στο διαδίκτυο και χρησιμοποιείται παγκοσμίως από εκατομμύρια ανθρώπους.

- Χρήστες ανά μήνα: 73,5 εκατομμύρια.
- Διαθέσιμες γλώσσες-πηγές: 88.
- Διαθέσιμες γλώσσες για παροχή αποτελεσμάτων: 35.

Εικόνα 1.5 Λογότυπο.
Piet Mondrian's Birthday
- March 7. 2002.



Μηχανισμός Λειτουργίας - Τεχνολογικές Καινοτομίες του Google⁹



Εικόνα 1.6
PigeonRank System.

Η μηχανή αναζήτησης Google παραμένει εστιασμένη στον βασικό σκοπό της δημιουργίας της, δηλαδή την ανάπτυξη της “τέλειας μηχανής αναζήτησης”, όπως διευκρινίζεται από τον Larry Page, δηλαδή στο «να καταλάβει ακριβώς τι εννοεί ο χρήστης κατά την αναζήτηση πληροφοριών και να του δώσει πίσω ακριβώς αυτό που θέλει». Η Google είναι μία link-based μηχανή αναζήτησης λαμβάνοντας υπόψη τη σημασία των υπερσυνδέσμων (hyperlink ή link) κατά την ταξινόμηση των αποτελεσμάτων, πέραν του μοντέλου ανάκτησης πληροφορίας (Information Retrieval) που χρησιμοποιεί. Η Google εισήγαγε τεχνολογικές καινοτομίες στον μηχανισμό λειτουργίας της δημιουργώντας τον αλγόριθμο PageRank και άλλαξε τον τρόπο αναζήτησης στο διαδίκτυο (Link-Based Τεχνικές αναζήτησης στο διαδίκτυο).

Από την αρχή οι δημιουργοί του Google αναγνώρισαν ότι η απόδοση γρηγορότερων και όσο το δυνατό ακριβέστερων αποτελεσμάτων απαιτούσε την κατασκευή ενός νέου είδους server. Έτσι, έχοντας στη διάθεσή τους περισσότερους από δέκα χιλιάδες servers, το λογισμικό πίσω από την τεχνολογία της μηχανής αναζήτησης διεξάγει μια σειρά υπολογισμών που απαιτούν λιγότερο από ένα δευτερόλεπτο, χρησιμοποιώντας εξαιρετικά πολύπλοκους αλγόριθμους.



Εικόνα 1.7 Λογότυπο.
Happy Birthday Picasso!
- October 25, 2002.

⁹ Βλ. Διαφάνεια 1β, στο Παράρτημα Ι, σελ. 98.

1. Ο αλγόριθμος ταξινόμησης σελίδων (PageRank)

Το PageRank είναι μιά από τις μεθόδους που χρησιμοποιεί το Google για να προσδιορίσει την καταληλότητα μιάς σελίδας ή το πόσο σημαντική είναι αυτή και βασίζεται σε μιά λογαριθμική μαθηματική φόρμουλα. Το PageRank είναι μέτρο της σημαντικότητας των σελίδων του ιστού, ανεξάρτητο του ερωτήματος (query independent). Πραγματοποιείται μια αντικειμενική μέτρηση της σημαντικότητας των σελίδων του ιστού με την λύση μιας εξίσωσης που περιέχει περισσότερες από 500 εκατομμύρια μεταβλητές (variables) και 2 δισεκατομμύρια όρους (terms). Δεν υπάρχει ανθρώπινη παρέμβαση ή παραποίηση των αποτελεσμάτων ανάκτησης των πληροφοριών και αυτός είναι ένας από τους κύριους λόγους που οι χρήστες εμπιστεύονται το Google ως πηγή αντικειμενικών πληροφοριών.

Δεν μπορούμε να ξέρουμε τις ακριβείς λεπτομέρειες της κλίμακας γιατί το μέγιστο PageRank όλων των σελίδων του ιστού αλλάζει κάθε μήνα όταν το Google επαναλαμβάνει την ευρετηρίαση του (Re-indexing).

Το PageRank είναι μιά “ψήφος” (vote), από όλες τις άλλες σελίδες του ιστού σε σχέση με το πόσο σημαντική είναι μια σελίδα. Ένας σύνδεσμος (link) προς μια σελίδα μετρά σαν μια ψήφος υποστήριξης και εάν δεν υπάρχουν καθόλου σύνδεσμοι δεν υπάρχει καθόλου υποστήριξη. Δηλαδή εάν μιά σελίδα A έχει συνδέσμους προς αξιόλογες σελίδες, τότε η γνώμη της A γίνεται βαρύνουσα: εφόσον η A «δείχνει» την B, με μεγάλη πιθανότητα η B να είναι μιά αξιόλογη σελίδα.

Το PageRank της κάθε σελίδας εξαρτάται από το PageRank των σελίδων οι οποίες «δείχνουν» σε αυτή. Δεν θα ξέρουμε τι PageRank έχουν αυτές οι σελίδες μέχρι που οι σελίδες που «δείχνουν» σε αυτές θα έχουν

υπολογισμένο το PageRank τους κ.ο.κ. Αυτό σημαίνει ότι υπολογίζουμε το PageRank μιας σελίδας χωρίς να ξέρουμε την τελική τιμή του PageRank των άλλων σελίδων. Βασικά κάθε φορά που “τρέχουμε” (Rank) αυτόν τον υπολογισμό υπολογίζουμε με ακριβέστερη προσέγγιση την τελική τιμή. Έτσι αυτό που χρειάζεται να κάνουμε είναι να θυμόμαστε την τιμή κάθε υπολογισμού



Εικόνα 1.8 Λογότυπο.
*Google celebrates MC Escher's
birthday - June 16, 2003.*

και να επαναλαμβάνουμε τους υπολογισμούς πολλές φορές μέχρι οι αριθμοί να μην αλλάζουν πλέον σημαντικά.

Το PageRank μιας σελίδας A θα είναι υψηλό εάν υπάρχουν πολλές σελίδες που «δείχνουν» (link) την A ή εάν υπάρχουν (έστω) λίγες σελίδες με υψηλό PageRank που «δείχνουν» την A. Μία σελίδα με λίγους υπερσυνδέσμους προς εξωτερικές σελίδες, συνεισφέρει υψηλότερο PageRank για τις σελίδες αυτές σε σχέση με σελίδες που έχουν πολλούς υπερσυνδέσμους προς εξωτερικές σελίδες. Έτσι σημαντικές σελίδες λαμβάνουν υψηλότερο PageRank και εμφανίζονται στην κορυφή των αποτελεσμάτων αναζήτησης.

Από την εμφάνιση της πρώτης μορφής του αλγόριθμου τα δεδομένα του Internet άλλαξαν. Ανάλογη όμως υπήρξε και η προσαρμογή του αλγόριθμου, με αποτέλεσμα να έχει ακόμα και σήμερα την ίδια επιτυχία με το παρελθόν. Οι καινούργιες τεχνικές αποτελούν ένα εμπορικό μυστικό. Ο αλγόριθμος μοιάζει να προσαρμόζεται απόλυτα πάνω σε καινούργια δεδομένα που προκύπτουν στο Internet και είναι διαρκώς ενημερωμένος για καινούργιες παραμέτρους που προκύπτουν από την εξέλιξη του Web. Φυσικά μεγάλο ρόλο στη διαρκή ανανέωση των στοιχείων που συγκροτούν τη μηχανή αναζήτησης παίζει και η μηνιαία ανανέωση του ευρετηρίου του Google (Index), που είναι ευρύτερα γνωστό σαν Google-Dance.

Εικόνα 1.9 Λογότυπο.
Salvador Dali May 11th. 2004.

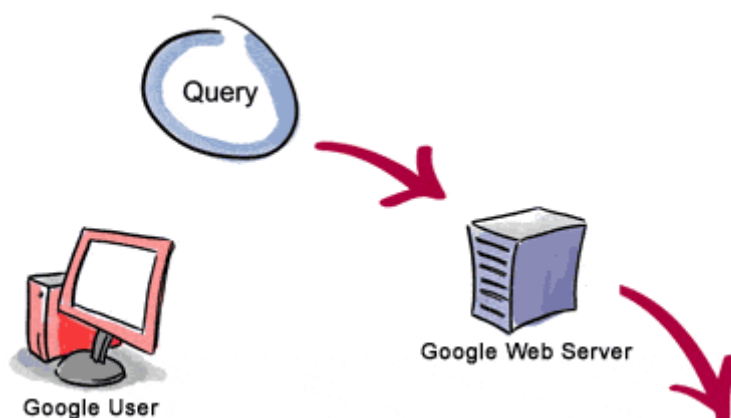


2. Συνδυασμένη Ανάλυση Ιστοσελίδων (Hypertext-Matching-Analysis)

Η μηχανή αναζήτησης Google αναλύει επίσης το περιεχόμενο μιας ιστοσελίδας. Η τεχνολογία του Google αναλύει ολόκληρο το περιεχόμενο κάθε ιστοσελίδας και το τοποθετεί σε υποκατηγορίες και με ακριβή τοποθεσία της κάθε λέξης. Επίσης αναλύει το περιεχόμενο γειτονικών ιστοσελίδων προσφέροντας τη σιγουριά ότι τα αποτελέσματα είναι τα πιο κατάλληλα στις αναζητήσεις του χρήστη.

3. Διάρκεια ζωής των ερωτημάτων (Query)

Η διάρκεια ζωής μιας ερώτησης στο Google διαρκεί κανονικά λιγότερο από μισό δευτερόλεπτο, όμως περιλαμβάνει διάφορα διαφορετικά στάδια που πρέπει να ολοκληρωθούν προτού να μπορέσουν τα αποτελέσματα να παραδοθούν σε κάποιον που αναζητά τις πληροφορίες.



3. Τα αποτελέσματα αναζήτησης επιστρέφονται στο χρήστη σε ένα μέρος ενός δευτερολέπτου.

1. Ο κεντρικός υπολογιστής δικτύου στέλνει την ερώτηση στους κεντρικούς υπολογιστές δεικτών. Το περιεχόμενο μέσα στους κεντρικούς υπολογιστές δεικτών είναι παρόμοιο με το δείκτη στο πίσω μέρος ενός βιβλίου - λέει ποιες σελίδες περιέχουν τις λέξεις που ταιριάζουν με την ερώτηση.



2. Η ερώτηση ταξιδεύει στους κεντρικούς υπολογιστές εγγράφου, οι οποίοι ανακτούν πραγματικά τα αποθηκευμένα έγγραφα. Τα αποκόμματα παράγονται για να περιγράψουν κάθε αποτέλεσμα αναζήτησης.



(πηγή: <http://www.google.com.gr/intl/el/corporate/tech.html>)



Εικόνα 1.10 Λογότυπο.
Salvador Dalí.

φανατικούς χρήστες του διαδικτύου αλλά και για πολλούς οικονομικούς αναλυτές. Είναι μια μηχανή αναζήτησης που ξεκινώντας από ένα ερευνητικό πρόγραμμα (project) στο πανεπιστήμιο Stanford και χάρη στον αλγόριθμο του PageRank ακολούθησε μιά ραγδαία ανοδική πορεία και έγινε η κυρίαρχη μηχανή αναζήτησης και κατέκτησε τον κόσμο.

Το φαινόμενο Google, σύμφωνα με κάποιους αναλυτές οφείλει την επιτυχία του στο γεγονός ότι δεν άλλαξε τον προσανατολισμό του και παρέμεινε μια απλή και εύκολη στην χρήση μηχανή αναζήτησης. Η αποστολή του Google, σύμφωνα με την ίδια την εταιρεία παραμένει αναλλοίωτη, να οργανώνει την παγκόσμια πληροφορία και να την κάνει προσβάσιμη και χρήσιμη για κάθε άνθρωπο.

Επίλογος

Το Google, η δημοφιλέστερη μηχανή αναζήτησης που υπάρχει αυτή την στιγμή στο διαδίκτυο, αποτελεί φαινόμενο και αντικείμενο μελέτης όχι μόνο για τους

Εικόνα 1.11 Λογότυπο.
Andy Warhol's Birthday - August 6, 2002.



Πηγές από το διαδίκτυο

<http://www.7.scu.edu.au/programme/fullpapers/1921/com1921.htm>

http://www.chip.gr/_magazine/viewthema.asp?id-thema=1678

<http://www.enet.gr/online-hprint.jsp?p=google&a=&id=74409396>

<http://www.iprcom.com/papers/pagerank/index.html>

<http://www.google.com.gr/intl/el/corporate/tech.html>

<http://www.google.com/holidaylogos.html>

<http://www.google.com/technology/index.html>

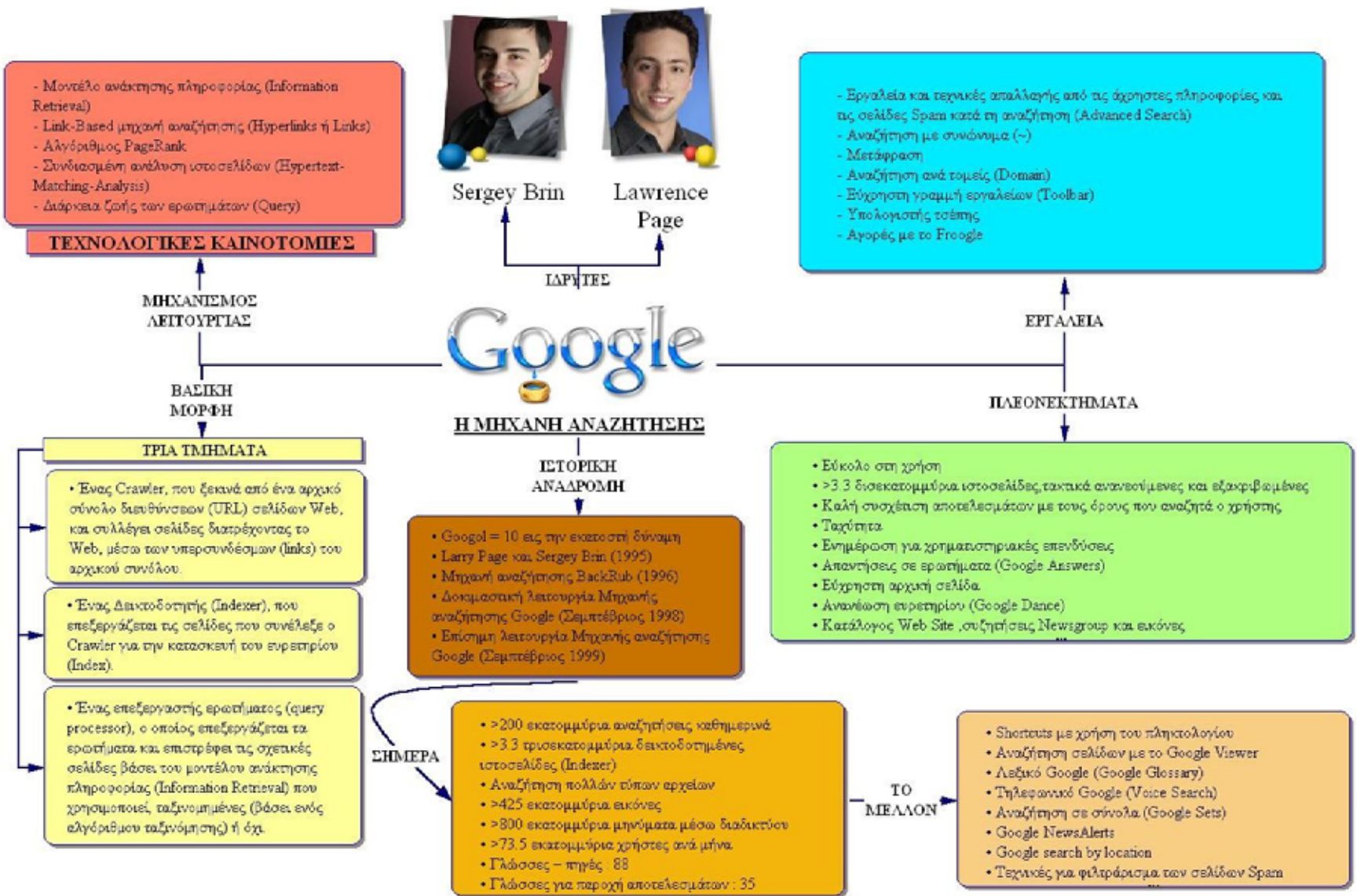
<http://www.google.com/technology/pigeonrank.html>

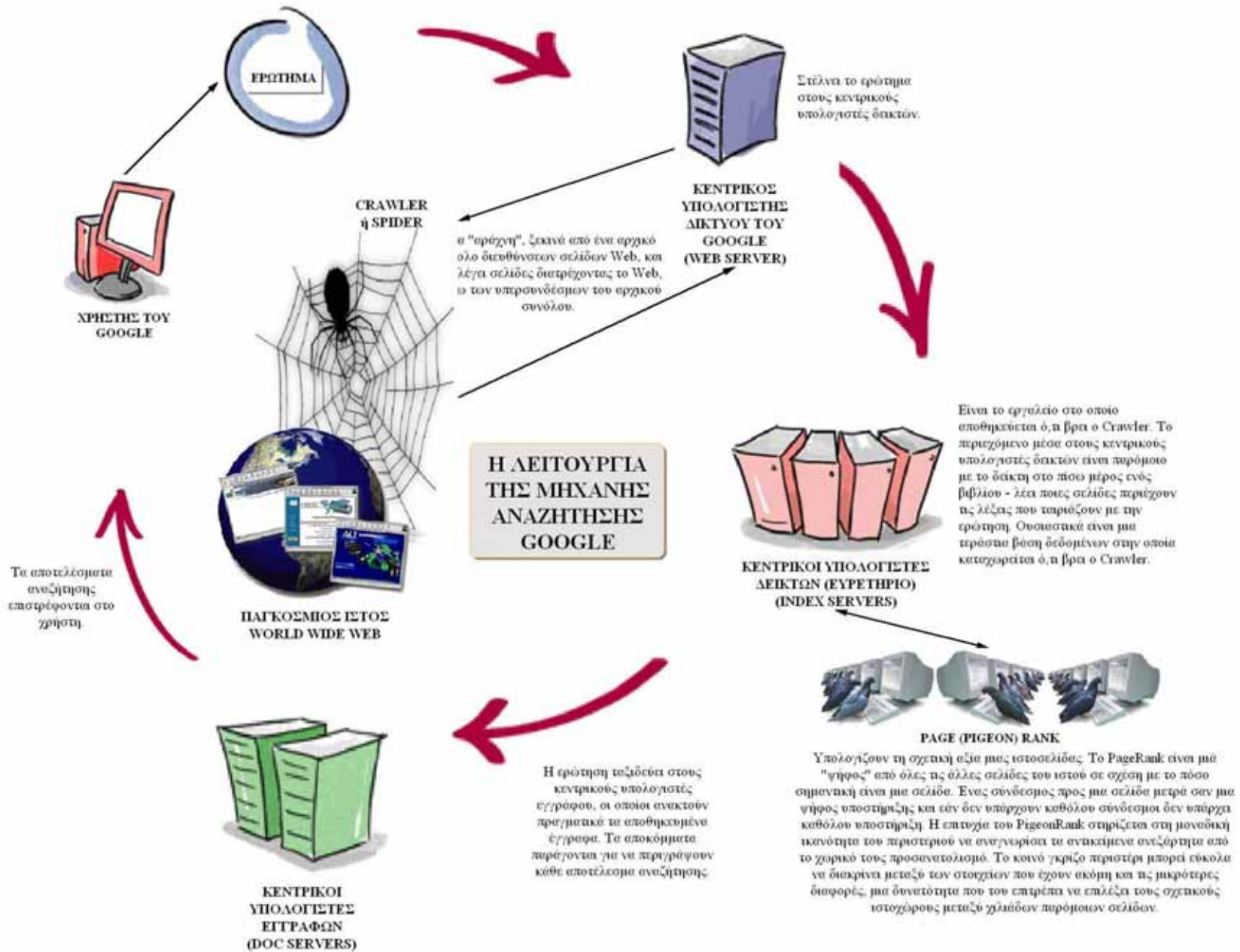
<http://www.gosub.gr/articles/article76.asp>

<http://www.logoogle.com/art.htm>

ΠΑΡΑΡΤΗΜΑΤΑ

(Διαφάνειες Εργασιών)





Διαφάνεια 1β «Η Λειτουργία της Μηχανής Αναζήτησης Google»