

Universitatea din Craiova

Facultatea de Automatică, Calculatoare și
Electronică

Proiect S.C.D.
*Trasaturile, configuratia si modul de
funtionare pentru Google*

COORDONATOR:

PROF.DR.ING. Matei Vînătoru

STUDENTĂ:

Stănică Roxana
Master - A.S.C

1. Google - Introducere.....	3
Căutarea în WEB.....	3
Câteva sfaturi pentru o căutare eficientă:.....	6
Utilizatori	7
Administrație.....	8
Funcționarea Google – pe scurt.....	8
Hardware	8
Cum manevrează Google cererile de căutare?	10
Sistemul PageRank.....	10
Googlebot-ul.....	11
Sistemul de fișiere Google	11
Scalabilitate:	12
1.1 Motoare de căutare – Evoluție: 1994-2000.....	13
1.2 Google: Evoluție concomitentă cu Web-ul.....	13
1.3 Caracteristici ale configurației.....	14
1.3.1 Capacitate îmbunătățită de căutare	14
1.3.2. O cercetare academică a motorului de căutare	15
2. Trăsăturile sistemului.....	15
2.1 PageRank: Ordonarea Web-ului.....	16
2.1.1 Descrierea unui calcul PageRank.....	16
2.1.2 Explicarea intuitivă	17
2.2 Textul link-ului	17
2.3 Alte trăsături	18
3. Activități similare	18
3.1 Extragerea de informații.....	19
3.2 Diferențe între Web și colecțiile bine controlate	19
4. Configurația sistemului.....	20
4.1 Scurtă descriere a arhitecturii Google.....	21
4.2 Structuri majore de date	21
4.2.1 BigFiles	22
4.2.2 Biblioteca (Repository)	22
4.2.3 Indexul documentelor.....	22
4.2.4 Lexiconul.....	23
4.2.5 Listele de hit-uri (Hit Lists).....	23
4.2.6 Indexul primar (Forward Index)	24
4.2.7 Indexul complementar.....	25
4.4 Indexarea Web-ului.....	25
4.5 Căutarea.....	27
4.5.1 Sistemul de clasificare.....	27
4.5.2 Scurta recapitulare.....	28
5. Rezultate și mod de funcționare	30
5.1 Parametri pentru stocare.....	30
5.2 Modul de funcționare a sistemului	31
5.3 Procesul de căutare	31
6. Concluzii	32
6.1 Activități viitoare.....	32
6.2 Căutare la standarde înalte	32
6.3 Arhitectura scalabilă	33
Referințe.....	34

1. Google - Introducere

Oricât de imparțiali încearcă jurnaliștii să fie, există o tendință clară în presă de a critica tot sau aproape tot ce face Microsoft și de a reda într-o lumină plăcută și optimistă realizările Google.

Este greu de spus când a început această modă și care au fost motivele din spatele ei. Un ochi atent ar observa că până la urmă este vorba de doi giganți din lumea IT care nu sunt atât de diferiți, factorul primar care face diferența între cei doi fiind vârsta.

Inițial, Google a reprezentat un proiect de cercetare conceput în 1996 de Larry Page și Sergey Brin, doi doctoranzi de la Stanford. Cei doi au plecat de la ipoteza că un *motor de căutare* care analizează relațiile dintre website-uri ar duce la obținerea de rezultate mai bune decât cele furnizate de tehnicile folosite atunci, și anume de câte ori se repetă cuvântul cheie pe pagină.

În martie 2001, Eric Schmidt a devenit membru al consiliului de directori, iar în august 2001, a fost numit CEO-ul companiei. Popularitatea noului motor de căutare s-a datorat design-ului aerisit, simplității și acurateții rezultatelor întoarse, iar numărul de utilizatori a crescut exponențial. În timp, după succesul motorului de căutare clasic, Google și-a extins brațele virtuale și peste imagini, video, știri, grupuri, mail, cumpărături, hărți și multe altele. În plus, există zvonuri că Google ar plănuși o rețea Internet proprie în Statele Unite, formată cu ajutorul kilometrilor întregi de fibră optică rămași după boom-ul de la sfârșitul anilor '90. Această fibră optică, care străbate suprafețe întinse ale Statelor Unite poartă numele de "fibră întunecată" (dark fiber).

Căutarea în WEB

Internetul poate fi considerat ca fiind cea mai mare bibliotecă disponibilă și actualizată permanent. De aceea uneori găsirea informației care ne este necesară poate fi dificilă.

World Wide Web (WWW sau Web) reflectă chiar prin numele sau trăsăturile sale esențiale :

- este destinat căutării informației în întregul Internet (world wide = în lumea întreagă)
- folosește hipertextul pentru organizarea informației, ceea ce face ca aceasta să apară ca o pânză de păianjen (web) permițând navigarea cu ușurință de la o pagină la alta.

Pentru regăsirea informației în Web se pot folosi următoarele instrumente:

- **serviciul de navigare** – browserul, care permite accesarea informațiilor prin introducerea adresei de web a paginii (numită URL - Uniform Resource Locator – identificator standard al locului în care se găsește

resursa), sau prin urmărirea unei hiperlegături (*en.* hyperlink) dintr-un alt document (care conține URL-ul paginii respective)

- **serviciul de căutare automată** – prin:
 - instrumente de căutare (*search tool*)
 - instrumente de indexare (*indexing tool*)
 - motor de căutare (*search engine*)

Browserul permite, în general, efectuarea următoarelor operații:

- navigarea printre paginile web și vizualizarea lor.
- urmărirea legăturilor dintre documente care conțin hiperlegături.
- copierea informațiilor din Internet pe propriul calculator.
- căutarea informației în Internet.
- regăsirea rapidă a informațiilor prin folosirea „semnelor de carte” („pagina favorită”) și a istoricului.
- accesarea altor servicii Internet: poșta electronică, mesagerie instantanee etc.

Așadar, browserul integrează accesul la mai multe servicii din Internet printr-o interfață accesibilă și ușor de utilizat.

O alternativă mai rapidă pentru „răsfoirea” paginilor web în căutarea informației dorite este apelarea la un **serviciu de căutare**, adică un site web care conține în general următoarele categorii de informații:

- instrucțiuni care arată modul în care trebuie folosit serviciul
- metoda prin care utilizatorul poate să propună un subiect pentru căutare

Portalurile sunt site-uri specializate care îndeplinesc funcția cataloagelor dintr-o bibliotecă. Aceste site-uri aranjează pe categorii sau domenii diverse site-uri existente în Internet și le ordonează în funcție de anumite criterii în cadrul categoriilor stabilite (pe subiecte, după popularitate, etc.). De multe ori, portalurile oferă și alte servicii pe lângă cel de căutare (poșta electronică, știri etc.).



Motoarele de căutare sunt site-uri care au rolul de a ajuta utilizatorul să găsească mai ușor și mai direct informația în Internet. De cele mai multe ori, în cuprinsul unui site cu rol de căutare, se regăsesc ambele funcționalități – indexare în catalog și meniu de căutare. Motorul de căutare este în fapt o aplicație care „răsfoiește” paginile web din Internet în căutarea cuvintelor sau frazelor cerute de utilizator. Pentru aceasta sunt folosite niște programe automate care alcătuiesc liste de cuvinte din interiorul site-urilor. Rezultatele căutării sunt afișate în funcție de relevanța stabilită de motorul de căutare, utilizând indexarea termenilor din aceste liste.



De aceea, termenii căutați trebuie să fie cât mai definatorii pentru subiectul în cauză (*keywords* – „cuvinte-cheie”).

Pentru a limita aria de căutare, se recomandă utilizarea modului de căutare avansată („*advanced search*”), care permite găsirea mai rapidă a informației datorită criteriilor multiple de căutare.

De exemplu, **Google** afișează doar paginile care includ *toți* termenii căutării. Nu este necesar să includem "and" între termeni. Pentru a rafina căutarea, se adaugă doar alți termeni, iar rezultatele vor conține un subset specific al paginilor întoarse de către cererea originală.

Putem exclude un cuvânt din căutare prin scrierea semnului minus ("-") imediat înaintea termenului pe care vrem să îl evităm. De asemenea, se pot căuta fraze prin includerea acestuia între ghilimele.



Câteva sfaturi pentru o căutare eficientă:

1. Fiți cât mai exact. Printr-o interogare precisă, se obțin mai puține rezultate și conținutul relevant este mai ușor de găsit. De exemplu, dacă am căutat cuvântul *lege*, am obținut aproximativ **3.580.000** de rezultate. Pentru căutarea *lege drept autor* am obținut **12.000** de rezultate, în timp ce pentru fraza exactă "**lege privind dreptul de autor**" am obținut **30** de rezultate.
2. Nu folosiți cuvinte uzuale. Utilizați cuvinte cât mai adecvate subiectului căutat, altfel utilitarul de căutare va returna zeci de pagini web cu informații irelevante pentru dumneavoastră. Vezi diferența între **5.570.000** rezultate pentru **masă** și **4.960** pentru **mijloace comunicare masă**.
3. Învățați să adaptați interogarea. Dacă interogarea dumneavoastră returnează prea multe rezultate, restrângeți aria de căutare. Dacă rezultatele returnate nu sunt suficiente, reformulați-o într-un mod mai general. Nu întotdeauna primele cuvinte -cheie sunt și cele mai bune.
4. Folosiți diferite forme ale cuvintelor. Puteți utiliza diferite cuvinte care se referă la subiectul căutat pentru a obține cât mai multe informații relevante pentru dumneavoastră.
5. Folosiți sinonimele. De exemplu, scrieți și „alergare” și ”jogging”. Dacă folosiți un utilitar de căutare care acceptă combinații de cuvinte - cheie, separați sinonimele prin cuvântul cheie OR
6. Folosiți citate între ghilimele. În cazul când căutați o anumită frază sau un titlu, plasați-le între ghilimele (de ex. „Internet pentru începători”) în formularul utilitarului de căutare.
7. Folosiți majuscule atunci când este necesar. Majoritatea directoarelor și indexurilor de căutare fac diferența între literele mici și literele mari din

șirul de caractere căutat. Dacă textul introdus conține numai litere mici, utilitarul de căutare va identifica numai textul scris fie cu litere mari, fie cu litere mici. În cazul în care scrieți și o litera mare, utilitarul de căutare presupune ca aceasta are o semnificație specială și va afișa numai rezultatele care corespund exact șirului respectiv.

8. Aflați secretele utilitarului de căutare folosit. Unele utilitare de căutare oferă facilități speciale prin care conținutul relevant poate fi găsit mai ușor.

Google este cel mai mare motor de căutare și datorită parteneriatelor sale cu Yahoo!, Netscape și altele este capabil să răspundă la mult mai multe interogări decât orice alt serviciu online asemănător.

- Interogări soluționate în fiecare zi: peste 150 de milioane
- Pagini web parcurse: peste 2,4 miliarde.
- Tipuri de fișiere căutate
- HyperText Markup Language (html)
- Adobe Portable Document Format (pdf)
- Adobe PostScript (ps)
- Lotus 1-2-3 (wk1, wk2, wk3, wk4, wk5, wki, wks,θ wku)
- Lotus WordPro (lwp)θ
- MacWrite (mw)θ
- Microsoft Excel (xls)θ
- Microsoft PowerPoint (ppt)
- Microsoft Word (doc)
- Microsoft Works (wks,θ wps, wdb)
- Microsoft Write (wri)
- Rich Text Format (rtf)
- Text (ans,θ txt)
- Imagini: peste 330 milioane
- Mesaje Usenet: peste 700 de milioane

Utilizatori

Google.com este unul din primele zece dintre cele mai populare site-uri de pe Internet și este folosit de oameni din lumea întreagă

- Media lunară a numărului de utilizatori: 55.6 milioane (Nilesen/NetRatings 5/02)
- Gradul de utilizare atât acasă, cât și la birou: #4 (Nilesen/NetRatings 5/02)
- Limbile pentru care Google asigură o interfață: 81

- Limbile în care Google oferă rezultate: 35
- Utilizatori: mai mult de jumătate din traficul Google.com provine din afara USA.

Administrație

Personalul Google include profesioniști cu experiență în domeniul tehnologiei, iar compania este susținută de fonduri provenite de la două firme ce se ocupă cu investițiile de risc.

- Număr aproximativ de angajați: 400
- Deținători ai titlului de doctor: peste 50

Proiectarea unui motor de căutare este o provocare. Motoarele de căutare indexează zeci sau chiar sute de milioane de pagini web, implicând un număr echivalent de termeni distincți. Acestea raspund la zeci de milioane de întrebări în fiecare zi. Deși importanța acestor motoare de căutare pe web este mare, totuși ele nu au constituit subiectul unei cercetări academice amănunțite. În plus, datorită gradului rapid de avansare a tehnologiei și dezvoltării continue a web-ului, metoda de creare a unui motor de căutare este foarte diferită de cea folosită acum trei ani. Această lucrare oferă o descriere amănunțită a motorului de cautare.

În afara problemelor de măsurare a capacității motoarelor tradiționale de căutare până la a putea suporta o cantitate importantă de date, există noi provocări tehnice corelate cu utilizarea informațiilor adiționale prezente în hipertext, cu scopul producerii de rezultate mai bune. Lucrarea de față pune aceasta întrebare, cum să construiești un sistem practic de măsurare care poate folosi informația adițională din hipertext. De asemenea, luăm în considerare modul de tratare efectiv al colecțiilor de hipertext, care nu pot fi în întregime controlate și unde fiecare este liber să publice ceea ce dorește.

Funcționarea Google – pe scurt

Hardware

Pentru a putea oferi capacitate suficientă de service, structura fizică a Google-lui este alcătuită din *cluster*e și *computere* situate peste tot în lume, cunoscute sub numele de *ferme de servere*. Aceste ferme de servere sunt alcătuite dintr-un număr mare de computere de nivel jos pe care rulează sisteme bazate pe *Linux* care operează cu GFS (Sistemul de fișiere Google), iar cele mai mari dintre aceste ferme au peste 1000 de *noduri* de stocare și peste 300 de *TB* de spațiu de stocare pe disk.

S-a speculat că Google are cel mai mare computer din lume. S-a estimat că Google are în jur de:

- 899 de *rack-uri*
- 79, 112 mașini
- 158,244 Unități centrale de prelucrare (UCP)
- 316,448 GHz putere de procesare
- 158,224 Gb de RAM
- 6,180 Tb de spațiu pe Hard

Clustere-le sunt folosite pentru a face ca 2 sau mai multe calculatoare să apară ca unul singur pentru o comunitate de utilizatori. Custerile sunt folosite pentru a oferi fiabilitate crescută, și/sau pentru a crește performanța de pe un singur calculator.

Un cluster este un grup de calculatoare legate între ele care lucrează împreună ca un computer paralel. Una dintre cele mai populare implementări o reprezintă cluster-ul cu noduri pe care rulează Linux-ul ca și sistem de operare și software-ul Beowulf pentru a implementa paralelismul.

Microsistemele Sun și-au lansat de asemenea un produs gen cluster numit Grid engine.

De asemenea denumită și cluster de servere, fermă de computere sau ranch, o fermă de servere este un grup de servere legate prin rețea care sunt găzduite într-o singură locație. O fermă de servere modernizează procesele interne prin distribuirea lucrului între componentele individuale ale „fermei” și expediază procesele computaționale prin folosirea puterii mai multor servere. „Ferma” se bazează pe software-ul „load-balancing” care realizează sarcini cum ar fi urmărirea cererii pentru procesarea puterii de la diferite mașini, stabilind prioritatea temelor și programându-le și reprogramând-le în funcție de prioritatea și cererea pe care utilizatorii o pun pe rețea. Când un server din „fermă” eșuează, un altul poate interveni ca și back-up.

Un nod este un dispozitiv care este conectat ca și parte a unei rețele de calculatoare. Nodurile pot fi computere, personal digital assistants (PDAuri), telefoane mobile sau variate dispozitive ale rețelei. Pe o rețea IP, un nod este orice dispozitiv cu o adresă IP. Nodurile sunt deseori conectate prin centrale, routere sau printr-un switch de rețea.

Linux-ul este un sistem de operare gratuit de tip Unix, creat cu originalitate de către Linus Torvalds cu ajutorul unor programatori din jurul lumii. Creat sub licența GNU General Public, codul sursă pentru Linux este disponibil gratuit pentru toată lumea.

Un rack de 19 inch este un sistem standard pentru încărcarea unor module electronice variate într-o stivă sau un „grilaj”, cu o lățime de 19 inch.

Cum manevrează Google cererile de căutare?

Când un user introduce o întrebare în căsuța de căutare la Google.com, aceasta este transmisă în mod aleator la unul dintre cluster-le lui Google. Întrebarea va fi apoi manevrată în mod exclusiv de către cluster. Un *echilibrator de încărcare* monitorizează cluster-ul, apoi împrăștie cererea către serverele din cluster ca să asigure că încărcarea de pe hardware este egală.

Echilibratorul de încărcare are nevoie de informații de profil pentru indica cum să re-mapeze procesele și canalele de comunicare pentru a optimiza performanța. Acesta va fi un echilibrator-optimizer static care va lua informația din uneltele de profil și va produce echilibrare de încărcare pentru viitoarele rulări.

Căutarea propriu-zisă are loc în 2 faze:

- În prima fază, cuvintele din cerere sunt verificate cu o listă din index-urile serverelor care conțin detaliile de potrivire ale documentelor. Sistemul PageRank este angajat într-un set relevant de documente, ce sunt identificate de trimeri la alte referințe ale cuvintelor și care vor determina scorul pentru fiecare document, dar la întoarcere afectează poziția documentului pe pagina rezultatelor.
- Rezultatele de la serverele indexate sunt apoi trimise la servere de documente în formularul de identificatori de documente. Rezidentă în serverele document este o copie a lui *World Wide Web*, de unde se ia rezumatul conținutului paginii web apoi le reîntoarce către server. Rezultatul este apoi procesat și este redat înapoi ca un document *HTML* care poate fi afișat pe *browser*-ul user-ului ca *pagină web*.

Sistemul PageRank

PageRank, denumit astfel după Larry Page, care l-a inventat, este unul dintre modurile prin care Google determină importanța paginii, care, la reîntoarcere, va decide unde anume pe lista de rezultate va apărea rezultatul.

Algoritmul exact a lui PageRank, așa cum este extras din “The Anatomy of a Large-Scale Hypertextual Web Search Engine” este în felul următor:

Presupunem că pagina *A* are paginile $T_1..T_n$ care pointează spre ea. Parametrul *d* este un factor, care poate fi setat între 0 și 1. De obicei îl punem pe *d* pe valoarea 0.85. Există mai multe detalii despre *d* în secțiunea următoare. De asemenea, $C(A)$ este definit ca fiind numărul de link-uri care „ies” din pagina *A*. RankPage-ului paginii *A* este redat în felul următor:

$$PR(A) = (1-d) + d (PR(T_1)/C(T_1) + \dots + PR(T_n)/C(T_n))$$

A se nota că PageRank-urile dintr-o probabilitate a distribuției peste paginile web este făcută astfel încât suma PageRank-urilor a tuturor paginilor web va fi una.

Googlebot-ul

Serverele Google sunt populate de către *web crawler*-urile Google-ului, Googlebot, care se mișcă de pe site pe site pe internet, făcând download a unor copii a paginilor web și salvându-le în index-ul Google-ului (cunoscut sub numele de cache) pentru viitoarele refetințe. Googlebot-ul însăși este rezident pe multe calculatoare care accesează simultan pe sute de pagini web. Acesta emulează un browser, așa că, mulți dintre webmasteri își vor da seama că atunci când vizitează, acesta lasă un „marcaj” în secțiunea de browser a protocolului site-ului web mai degrabă decât în secțiunea de „păianjen” decât cele mai multe web crawler-uri care se înregistrează.

Un web crawler (cunoscut și sub numele de web spider) este un program care deschide World Wide Web într-o manieră metodică, automată. Web crawlerii sunt în principal folosiți pentru a crea o copie a tuturor paginilor web vizitate pentru procesări ulterioare de către un motor de căutare, care va indexa paginile download-ate pentru a oferi căutari mai rapide.

Există două metode pe care Googlebot le folosește pentru a găsi o pagină web: ori prin ajungerea la pagina web după ce „s-a târât” prin link-urile acesteia, ori ia pagina după ce a fost înscrisă de către webmaster. Prin înscrierea link-ului principal Googlebot va trece prin toate link-urile care se găsesc în indexul paginii și în fiecare pagină următoare, până când întregul site a fost indexat.

Sistemul de fișiere Google

Sistemul de fișiere Google este un sistem propriu de administrare a fișierelor dezvoltat de către Sanjay Ghemawat, Shun-Tak Leung și Urs Holyle pentru Google, ca un mijloc de a manevra un număr masiv de cereri peste un număr mare de clustere ale serverelor.

Sistemul a fost creat la fel ca majoritatea celorlalte *sisteme distribuite* de fișiere, pentru performanță maximă, pentru a face față unui număr mare de utilizatori, *scalabilitate*, pentru a fi capabil să manevreze expansiunile inevitabile, siguranță, pentru a asigura maximul de *uptime* și *disponibilitate*, pentru a asigura faptul că, toate calculatoarele sunt disponibile pentru a trata întrebările.

Scalabilitate:

1. Cuvânt des folosit care se referă la felul în care un sistem hardware sau software se poate adapta pentru a crește cererile.

2. Se referă la orice la care i se poate schimba mărimea

3. Când este folosit pentru a descrie un sistem al computerului, proprietatea de a rula mai mult de un procesor.

Uptime este o măsură a timpului de când un computer de sistem a fost pornit. A fost creat pentru a descrie antonimul cuvântului downtime, adică timpul cât un sistem a fost non-operațional.

Disponibilitatea este o măsură a cât de mult timp o rețea sau o conexiune a rulat. În general, formula este: Timpul rulării/Timpul măsurat (timpul de rulare împărțit la timpul măsurat). Deși, dacă se măsoară ceva pentru 20 de minute și acesta a fost disponibil doar 19 din acestea, disponibilitatea va fi de 95%.

Din cauza deciziei lui Google de a folosi un număr mare de computere de nivel jos în loc să folosească un număr mai mic de sisteme de tip server, sistemul de fișiere Google a trebuit să fie creat astfel încât să poată trata eșecurile sistemului, să efectueze monitorizarea constantă a sistemelor, detectarea de erori, *toleranța erorilor* și recuperarea automată.

Toleranța erorilor este disponibilitatea unui sistem de a răspunde unui eșec neașteptat a dispozitivelor hardware sau software. Există multe nivele de tolerare a erorilor, cel mai jos fiind amabilitatea de a continua operațiile în cazul unei căderi de tensiune. Multe dintre sistemele de computere care tolerează erorile monitorizează toate operațiile, aceasta este în cazul în care, fiecare operație este efectuată pe două sau mai multe sisteme duplicate, așa că în cazul în care unul dintre acestea eșuează celălalt poate prelua controlul.

Aceasta însemna că, toate clustere-le ar trebui să rețină replici multiple a informației create de web crawler-ii Google-lului. Aceasta este foarte relevant pentru Gmail care a fost recent implementat pe Google, unde email-ul personal al utilizatorilor trebuie să aibă backup pentru a preveni pierderea informațiilor.

Din cauza mărimii *bazei de date* a Google-ului, sistemul a trebuit creat ca să poată trata fișiere enorme de mai mulți giga-bytes totalizând mulți terabytes ca și mărime. A fost creat astfel datorită faptului că stocarea fișierelor în kilobytes ar însemna să existe miliarde de fișiere, care s-ar putea dovedi imposibil de manevrat. O altă metodă angajată să trateze fișiere de mărime enormă este că schimbările (cunoscute și sub numele de mutații) sunt mai degrabă adăugate decât să aibă fișiere peste care s-a rescris, astfel se minimizează timpul de acces la fișiere.

GFS este un sistem propriu așa că aplicațiile software pot fi construite de obicei în jurul lui, asigurând faptul că Google are control maxim asupra sistemului, în același timp permițând sistemului să rămână flexibil.

Web-ul creează noi provocări pentru obținerea de informații. Cantitatea de informații de pe web crește într-un ritm alert, pe măsura numărului de noi

utilizatori lipsiți de experiență în arta căutării pe web. De obicei, oamenii navighează pe web folosind graficul acestuia de link-uri, adeseori începând cu indici superiori calitativ, mentinuți de intervenția umană, cum ar fi Yahoo sau cu motoare de căutare. Listele unde intervine mintea umană acoperă subiecte diverse și populare, dar sunt subiective, costisitoare de întreținut și menținut, greu de îmbunătățit și nu pot acoperi toate subiectele ce țin de domenii specializate. Motoarele de căutare automate, care se bazează pe potrivirea de cuvinte-cheie, oferă, în mod obișnuit, prea multe rezultate neconcludente. Pentru a înrăutăți situația, unii agenți publicitari încearcă să câștige atenția prin diferite metode destinate să inducă în eroare motoarele automate de căutare. Ei au construit un motor de căutare pe scară largă care are răspuns pentru multe din problemele sistemului existent. Acesta utilizează, în special, structura adițională din hipertext pentru a oferi rezultate concludente. Au numit acest sistem Google deoarece este un cuvânt asemănător cu googol sau 10¹⁰⁰ și corespunde cel mai bine intenției lor de a construi motoare de căutare pe scară largă.

1.1 Motoare de căutare – Evoluție: 1994-2000

Tehnologia motoarelor de căutare a fost obligată să se reinventeze în mod constant pentru a putea ține pasul cu expansiunea web-ului. În 1994, unul dintre primele motoare de căutare, World Wide Web Worm (WWW), avea un index de 110,000 pagini web și documente accesibile pe web. În noiembrie 1997, motoarele de căutare cele mai performante realizau indexarea de la 2 milioane (WebCrawler) până la 100 de milioane de documente web (din Search Engine Watch). Este previzibil că, până în anul 2000, un index complet al Web-ului să conțină peste un bilion de documente. În același timp, numărul interogărilor la care pot face față motoarele de căutare a crescut într-un ritm incredibil de rapid. În martie și aprilie 1994, World Wide Web Worm primea în jur de 1500 de întrebări pe zi. În noiembrie 1997, Altavista pretindea că rezolva aproximativ 20 de milioane de interogări pe zi. Odată cu creșterea numărului de utilizatori ai web-ului și sistemelor automate care puneau întrebări motoarelor de căutare, este posibil ca motoarele performante de căutare să poată răspunde la sute de milioane de întrebări pe zi până în anul 2000. Scopul sistemului este să rezolve multe dintre aceste probleme, atât din punct de vedere al calității, cât și al scalabilității, coordonate introduse de tehnologia motoarelor de căutare care au atins cifre extraordinare.

1.2 Google: Evoluție concomitentă cu Web-ul

Crearea unui motor de căutare care să țină pasul cu web-ul zilelor noastre este subordonată multor provocări. Tehnologia rapidă de parcurgere este necesară pentru a ține o evidență a documentelor web și pentru a le actualiza în

permanență. Spațiul de stocare trebuie folosit în mod eficient pentru organizarea indicilor și, opțional, chiar a documentelor. Sistemul de indexare trebuie să proceseze cât mai bine sute de gigabiți de date. Întrebările trebuie să aibă un răspuns într-un interval cât mai scurt, la o rată de sute până la mii pe secundă.

Aceste sarcini devin din ce în ce mai dificile pe măsură ce Web-ul se dezvoltă. Totuși, performanțele hardware și costurile au crescut și ele în mod constant pentru a rezolva parțial această problemă. Există, bineînțeles, câteva excepții importante de la acest progres cum ar fi timpul de căutare pe disc și performanțele sistemului de operare. În proiectarea sistemului Google, s-au luat în considerare atât rata de evoluție a Web-ului, cât și schimbările tehnologice. Google este construit în așa fel încât să parcurgă eficient mari cantități de date. În același timp, Google utilizează eficient spațiul de stocare pentru organizarea indexului. Structurile sale informaționale sunt optimizate pentru accesul rapid și concludent (vezi secțiunea 4.2). Se speră că prețul indexării și stocării de text sau HTML să devină în cele din urmă mai mic în raport cu cantitatea de date ce va fi disponibilă. Acest lucru va genera caracteristici mai performante de organizare pentru sistemele centralizate ca Google.

1.3 Caracteristici ale configurației

1.3.1 Capacitate îmbunătățită de căutare

Scopul principal a fost acela de a îmbunătăți calitatea motoarelor de căutare. În 1994, existau păreri că un index complet de căutare va face posibilă găsirea rapidă a oricărui răspuns. În opinia Best of the Web 1994 – Navigators, “Cel mai bun sistem de navigare ar trebui să faciliteze găsirea unui răspuns la aproape orice întrebare pe Web (din momentul introducerii datelor)”. Totuși, Web-ul anului 1997 este diferit. Oricine a utilizat recent un motor de căutare, poate oricând să spună că nu complexitatea indexului este singurul factor care asigură calitatea rezultatelor căutării. Rezultatele neconcludente (junk results) estompează de obicei relevanța rezultatelor de care un utilizator este într-adevar interesat. De fapt, din noiembrie 1997, numai unul dintre cele 4 motoare de căutare comerciale își menționează numele (returnează pagina proprie de căutare ca răspuns la numele sau în topul primelor 10 rezultate). Una dintre principalele cauze ale acestei probleme este aceea că numărul documentelor din indici a crescut progresiv acoperind diferitele grade de importanță, în timp ce modul de percepere al utilizatorului asupra acestor documente nu s-a schimbat. Oamenii manifestă în continuare interes doar pentru primele 10 rezultate. Din acest motiv, pe măsura ce colecția de date se mărește, este nevoie de instrumente cu un grad ridicat de precizie (adică număr de documente relevante returnate, de exemplu în topul primelor 10 rezultate). Într-adevar, se vrea ca noțiunea de “relevant” să includă caracteristica de cel mai important ignorând chiar căutarea

în sine (numărul total de documente relevante pe care sistemul este capabil să le ofere. Există chiar un sentiment de optimism care se referă la utilizarea mai multor informații hipertextuale ce pot contribui la îmbunătățirea căutării și la alte aplicații. În particular, structura și textul link-ului oferă multe informații pentru efectuarea de analize relevante și filtrări de calitate. Google folosește atât structura link-ului cât și textul acestuia (vezi secțiunile 2.1 și 2.2).

1.3.2. O cercetare academică a motorului de căutare

În afară de o creștere considerabilă, Web-ul a adoptat și o caracteristică comercială de-a lungul timpului. În 1993, 1,5% din serverele web erau pe domeniul .com. Numărul acestora a atins 60 de procente în 1997. În același timp, motoarele de căutare au părăsit academicul în favoarea comercialului. Până în prezent dezvoltarea majorității motoarelor de căutare a fost efectuată de companii care aveau prea puțin de-a face cu detaliile tehnice. Acest fenomen a determinat rămânerea în umbră a tehnologiei motoarelor de căutare și orientarea acestora către comercial. Cu Google, scopul este acela de a orienta dezvoltarea și perceperea acestora către sfera academicului.

O altă trăsătură importantă a acestei configurații era construirea de sisteme pe care mulți oameni să le poată utiliza. Modul de utilizare era important pentru deoarece unele dintre cele mai interesante cercetări implică manipularea unei vaste cantități de informații care sunt disponibile pe sistemele web moderne. De exemplu, zilnic sunt inițiate zeci de milioane de căutări. Totuși, toate aceste date sunt foarte dificil de obținut, în special pentru că se consideră că au valoare comercială.

Scopul final al acestei configurații era realizarea unei arhitecturi care să poată susține activități inedite de căutare pe scară largă a informației pe web. Pentru a sprijini aceste activități de căutare, Google stochează toate documentele în format comprimat pe care le parcurge. Unul dintre scopurile principale ale configurației Google a fost acela de a stabili un mediu unde ceilalți căutători puteau intra repede, de a procesa secțiuni considerabile ale web-ului și oferi rezultate interesante care ar fi fost greu de oferit prin adoptarea unei alte metode. În intervalul scurt de timp în care sistemul a fost construit, existau deja câteva lucrări care foloseau bazele de date oferite de Google și multe altele sunt realizate acum. O altă este aceea de a constitui un mediu asemănător cu un laborator spațial unde cercetătorii sau chiar studenții să poată propune și chiar executa experimente interesante folosind informațiile complexe ale web-ului.

2. Trăsăturile sistemului

Motorul de căutare Google este caracterizat de două trăsături importante care ajută la producerea de rezultate cu un grad ridicat de precizie. În primul

rând, Google se folosește de structura de link-uri a Web-ului pentru a calcula un indice calitativ al fiecărei pagini web. Această estimare a nivelului calitativ se numește *PageRank*. În al doilea rând, Google utilizează link-urile pentru a îmbunătăți rezultatele căutării.

2.1 PageRank: Ordonarea Web-ului

Graficul de link-uri al web-ului este o resursă importantă care a rămas în mare parte neutilizată de motoarele de căutare. S-au realizat hărți conținând nu mai puțin de 518 milioane din aceste hyperlink-uri, o mostră semnificativă a totalului. Aceste hărți permit calcularea rapidă a PageRank-ului unei pagini web, o măsura obiectivă a importanței link-urilor care corespunde cu ideea subiectivă de importanță a oamenilor. Datorită acestei corespondențe, PageRank reprezintă o metodă excelentă de stabilire a gradului de importanță a rezultatelor căutărilor bazate pe cuvinte cheie. Pentru cele mai populare subiecte, un text simplu care se potrivește cu căutarea și care este limitat la titluri ale paginii web este foarte bine reprezentat atunci când PageRank stabilește importanța rezultatelor (demonstrație disponibilă la google.stanford.edu). Pentru căutărilor ce au la bază un text integral în sistemul principal Google, PageRank este, de asemenea, de mare ajutor.

2.1.1 Descrierea unui calcul PageRank

Literatura de specialitate referitoare la link-uri a fost raportată la web, în general prin numerotarea link-urilor sau backlink-urilor unei pagini date. Acest lucru stabilește cu aproximație importanță sau calitatea unei pagini. PageRank extinde această idee nu prin efectuarea unei numerotări a link-urilor din toate paginile, ci prin stabilirea numărului de link-uri dintr-o pagină. PageRank este definit după cum urmează:

Presupunem că pagina A este formată din paginile $T_1 \dots T_n$ care se referă la aceasta (adică sunt link-uri). Parametrul d este un factor de nivelare care se află între 0 și 1. De obicei, stabilim valoarea 0.85 pentru acest factor. Mai multe detalii despre d sunt oferite în secțiunea următoare. De asemenea, $C(A)$ este definit ca un număr de link-uri care nu fac parte din pagina A. PageRank-ul paginii A este după cum urmează:

$$PR(A) = (1-d) + d(PR(T_1)/C(T_1) + \dots + PR(T_n)/C(T_n))$$

Trebuie reținut că PageRank formează o distribuție a probabilității paginilor web, astfel că suma tuturor paginilor web ce țin de PageRank este 1. PageRank sau $PR(A)$ poate fi calculat utilizând un simplu algoritm repetabil și care corespunde principalului vector propriu al matricii link-ului normalizat al web-ului. De asemenea, un PageRank pentru 26 milioane de pagini web poate fi

calculat în câteva ore într-un punct de lucru de mărime medie. Există multe alte detalii care depășesc intenția acestei lucrări.

2.1.2 Explicarea intuitivă

PageRank poate fi considerat un model al comportamentului utilizatorului. Să presupunem că există un navigator oarecare care vizitează o pagină web aleasă la întâmplare și care accesează link-urile, fără a reveni la pagina inițială: în cele din urmă se va plictisi și se va orienta spre altă pagină web aleasă la întâmplare. Probabilitatea că acest navigator să viziteze o pagină este reprezentată de PageRank. Iar d , factorul de nivelare, reprezintă probabilitatea ca navigatorul să se plictisească la fiecare pagină accesată și să continue căutarea paginilor la întâmplare. O variație importantă este aceea de a adăuga doar factorul de nivelare d unei singure pagini sau unui grup de pagini. Acest lucru permite personalizarea și poate face aproape imposibilă inducerea deliberată în eroare a sistemului pentru obținerea unui calificativ superior. Pentru mai multe referiri la PageRank.

O altă explicație intuitivă este că o pagină poate avea un PageRank ridicat dacă există mai multe pagini care fac referire la aceasta sau dacă există câteva pagini care au un PageRank ridicat și care o recomandă. În mod intuitiv, paginile la care se face referire din multe colțuri ale web-ului sunt considerate importante. De asemenea, paginile care probabil au o singură referire de la gazda Yahoo-ului sunt considerate importante. Dacă o pagina nu are un nivel calitativ ridicat sau are un link insuficient, este mai mult decât probabil ca pagina gazdă a Yahoo-ului nu va avea nici un link pentru aceasta. PageRank face față ambelor situații și chiar mai mult de atât prin propagarea recursivă a gradului de importanță în întreaga structură de link-uri a web-ului.

2.2 Textul link-ului

Textul link-ului este tratat într-un mod cu totul special de motorul de căutare. Majoritatea motoarelor de căutare asociază textul link-ului cu pagina de care se leagă. În plus, se asociază cu pagina pe care link-ul respectiv o indică. Acest sistem prezintă mai multe avantaje. În primul rând, link-urile oferă deseori descrieri mai precise ale paginilor web decât o fac paginile respective. În al doilea rând, link-urile pot face referire la documente care nu pot fi indexate de un motor de căutare bazat pe text, cum ar fi: imagini, programe sau baze de date. Acest lucru face posibilă returnarea paginilor web care nici macăr nu au fost parcurse. Trebuie reținut că paginile care nu au fost parcurse pot cauza probleme din moment ce nu le-a fost niciodată verificată validitatea înainte de a fi oferite utilizatorului. În acest caz, motorul de căutare poate oferi o pagină care nu a

existat niciodată cu adevărat, dar care are hyperlink-uri care fac referire la ea. Totuși, este posibil ca rezultatele să fie sortate, astfel ca această problemă apare rareori.

Ideea corelării textului link-ului cu pagina web la care se referă a fost implementată în World Wide Web Worm, în special pentru că ajută la căutarea informației de tip non-text și mărește aria de acoperire a căutării prin numărul mai redus de documente descărcate. Se folosește propagarea de link-uri deoarece textul link-ului poate contribui la oferirea de rezultate mai bune. Utilizarea eficientă a text-ului link-ului este dificilă din punct de vedere tehnic din cauza cantităților mari de date care trebuie procesate. În procesul de parcurgere a 24 milioane de pagini, s-au indexat peste 259 de milioane de link-uri.

2.3 Alte trăsături

În afară de PageRank și de utilizarea textului link-ului, Google are și alte trăsături. Prima este aceea că are informații de bază pentru toate căutările și astfel utilizează, în mod frecvent, proximitatea în procesul de căutare. A doua se referă la faptul că Google are în vedere detaliile vizuale ale prezentării cum ar fi mărimea fonturilor. Cuvintele scrise cu un font mai mare sau cu caractere îngroșate sunt percepute altfel decât celelalte cuvinte. A treia trăsătură este aceea că se ține o evidență a întregului șir al paginilor HTML.

3. Activități similare

Cercetarea navigării pe web are o istorie scurtă și concisă. World Wide Web Worm (WWW) a fost unul dintre primele motoare de căutare. Acesta a

fost ulterior urmat de alte câteva motoare de căutare academice, dintre care multe sunt acum companii publice. Comparativ cu dezvoltarea Web-ului și importanța motoarelor de căutare, există puține documente valoroase referitoare la motoarele de căutare recente. Conform lui Michael Mauldin (cercetător la Lycos Inc) “serviciile variate (inclusiv Lycos) urmăresc îndeaproape detaliile acestor baze de date”. Totuși, există studii despre trăsăturile specifice ale motoarelor de căutare. Un studiu bine realizat este acela care duce la obținerea de rezultate prin procesarea ulterioară a rezultatelor motoarelor de căutare comerciale sau la producerea, la scară redusă, de motoare de căutare individualizate. Cercetări amanunțite au fost făcute pe sistemele care se ocupă cu extragerea de informații, în special pe colecțiile bine controlate. În următoarele două secțiuni, se va discuta despre zone unde această cercetare trebuie extinsă pentru a se putea lucra mult mai bine pe web.

3.1 Extragerea de informații

Cercetarea sistemelor care se ocupă cu extragerea de informații a fost inițiată acum mulți ani și este acum bine dezvoltată. Totuși, majoritatea cercetărilor asupra acestor sisteme se concentrează pe colecțiile omogene și bine întreținute cum ar fi: colecțiile de lucrări științifice sau articole despre un subiect apropiat. Într-adevăr, prototipul extragerii de informații, The Text Retrieval Conference folosește pentru teste colecții mici dar bine controlate. Testul “Very Large Corpus” are doar 20 GB în comparație cu 147 GB ai celor 24 de milioane de pagini web. Tot ce funcționează bine pe TREC nu oferă de obicei rezultate satisfăcătoare pe web. De exemplu, modelul de vector spațial standard încearcă să returneze documentele care se apropie cel mai mult de interogare, ținând seama de faptul că atât interogarea cât și documentul sunt vectori definiți de ocurența cuvintelor lor. Pe web, această strategie oferă de obicei documente scurte care conțin interogarea și câteva cuvinte în plus. De exemplu, s-a observat că un motor important de căutare returnează o pagină conținând doar “Bill Clinton Sucks” și o fotografie dintr-o interogare referitoare la Bill Clinton. Unii sunt de părere că, pe web, utilizatorii ar trebui să specifice cu mai multă acuratețe ceea ce doresc să afle și să adauge mai multe cuvinte întrebării respective. Dacă un utilizator formulează o interogare de tipul “Bill Clinton” ar trebui să primească rezultate concludente din moment ce există o cantitate considerabilă de informații valoroase disponibile pentru acest subiect. Prin aceste exemple procesul standard de extragere de informații trebuie extins pentru a putea utiliza web-ul în mod eficient.

3.2 Diferențe între Web și colecțiile bine controlate

Web-ul reprezintă o colecție vastă de documente eterogene. Documentele de pe web sunt caracterizate printr-o variație internă constantă a documentelor și de asemenea prin meta informație externă care poate să fie disponibilă. De exemplu, documentele diferă din punct de vedere intern prin limbaj (atât uman, cât și de programare), vocabular (adrese de e-mail, link-uri, coduri zip, numere de telefon, numere de producție), tip de format (text, HTML, PDF, imagini, sunete) și pot fi generate chiar de sistem (fișiere jurnal sau extrase din baza de date). Pe de altă parte, definim meta informația externă ca informația care poate fi dedusă despre un document, dar care nu este conținută de acesta. Exemple de meta informație externă includ aspecte precum reputația sursei, frecvența actualizării, calitatea, popularitatea sau uzajul și link-urile. Nu numai că sursele posibile ale meta informației externe sunt variate, dar lucrurile care sunt măsurate variază de asemenea din punct de vedere al gradelor de importanță. De exemplu, dacă se compară informația uzuală de pe o pagină gazdă importantă ca Yahoo, care primește zilnic milioane de accesări de pagini, cu un articol istoric obscur care poate fi accesat o dată la zece ani. În mod evident, aceste două aspecte trebuie tratate în mod diferit de un motor de căutare.

O altă diferență majoră între web și colecțiile tradiționale este aceea că, practic, nu există nici un control asupra a ceea ce pot publica oamenii pe web. Corelând această flexibilitate de a putea publica orice cu influența considerabilă a motoarelor de căutare care pot afecta traficul (și există diverse companii care manipulează deliberat motoarele de căutare pentru a obține profit), ajungem la o problemă serioasă. Această problemă nu a fost pusă în cazul sistemelor tradiționale de extragere de informații. De asemenea, este interesant de observat că eforturile depuse în cazul metadatelor au eșuat în cazul motoarelor de căutare, deoarece orice text de pe pagină care nu este direct oferit utilizatorului este supus manipulării efectuate de motoarele de căutare. Există chiar numeroase companii care s-au specializat în manipularea motoarelor de căutare în vederea obținerii de profit.

4. Configurația sistemului

4.1 Scurtă descriere a arhitecturii Google

În această secțiune va fi o descriere despre modul de funcționare al sistemului. În secțiunile următoare se va discuta despre aplicațiile și structurile de date care nu vor fi menționate în această secțiune. Cea mai mare parte din Google este realizată în C sau C++ pentru eficiență și poate rula atât în Solaris, cât și în Linux.

În Google parcurgerea web-ului (descărcarea de pagini) este făcută de mai multe crawlere diferite. Există un server URL care trimite listele de URL-uri ce trebuie găsite de crawlere. Paginile web care sunt găsite sunt apoi returnate serverului de stocare, care le memorează. Acesta comprimă paginile și le depune într-o bibliotecă. Orice pagină web are un număr de identificare numit docID, care este oferit ori de câte ori un nou URL este analizat și extras dintr-o pagină web. Funcția de indexare este realizată de indexer și de sorter. Indexer-ul îndeplinește o serie de funcții. Citește documentele din bibliotecă, decompromă documentele și le analizează. Fiecare document este convertit într-o serie de asocieri de cuvinte numite hit-uri. Acestea înregistrează cuvântul și poziția sa în document, aproximează dimensiunea fontului și tipurile de litere folosite. Indexer-ul distribuie aceste hit-uri într-o serie de categorii, creând un index parțial dezvoltat de sortare. Indexer-ul mai îndeplinește și o altă funcție importantă. Analizează toate link-urile din fiecare pagina web și stochează informații importante despre acestea într-un fișier de link-uri. Acest fișier conține informații suficiente pentru a stabili unde ne îndreaptă link-ul respectiv, precum și textul link-ului.

Sistemul de analizare a URL-urilor citește fișierul de link-uri și convertește URL-urile relative în URL-uri absolute și, respectiv, în docID-uri. Plasează textul link-ului în indexul inițial care este asociat cu docID-ul la care se referă link-ul. Acesta generează de asemenea o bază de date de link-uri care nu sunt altceva decât corespondentele docID-urilor. Această bază de link-uri este folosită pentru calcularea PageRank-urilor pentru toate documentele.

Sorter-ul preia categoriile care sunt sortate de docID și le clasifică după wordID pentru a forma un index complementar (inverted index). Un program numit DumpLexicon preia această listă împreună cu lexiconul produs de indexer și formează un lexicon nou care poate fi utilizat de searcher. Searcher-ul este rulat de un server și folosește lexiconul construit de DumpLexicon împreună cu indexul complementar și PageRank pentru a răspunde întrebărilor.

4.2 Structuri majore de date

Structurile de date ale Google sunt optimizate astfel încât o colecție amplă de documente poate fi parcursă și indexată cu puțin efort. Deși CPU-urile și majoritatea ratelor de input output s-au îmbunătățit de-a lungul anilor o simplă căutare pe disc tot necesită 10 ms pentru a fi realizată. Google este proiectat să evite acest gen de căutări de câte ori este posibil, iar acest lucru a avut o influență considerabilă asupra formatului structurilor de date.

4.2.1 BigFiles

BigFiles reprezintă fișiere virtuale care parcurg multiple fișiere de sistem și care sunt operaționale cu 64 de biți integrali. Împărțirea între multiplele fișiere de sistem este realizată în mod automat. Aceste fișiere suportă de asemenea alocarea și dealocarea descriptorilor de fișiere, din moment ce sistemele de operare nu oferă destul pentru ceea ce avem noi nevoie. BigFiles suportă și opțiuni simple de comprimare.

4.2.2 Biblioteca (Repository)

Biblioteca conține HTML-ul integral al fiecărei pagini web. Fiecare pagină este comprimată prin folosirea zlib. Optarea pentru o tehnică de compresie reprezintă echilibrul între viteză și proporția comprimării. S-a ales viteza zlib dintr-o serie de îmbunătățiri semnificative aduse comprimării de bzip. Rata compresiei bzip era de aproximativ 4 la 1 în bibliotecă, în comparație cu zlib care oferea o rată de 3 la 1. În bibliotecă, documentele sunt stocate unul după altul și sunt prefixate de docID, precizându-li-se lungimea și URL-ul. Biblioteca nu solicită alte structuri de date care să fie folosite pentru ca aceasta să fie accesată. Acest lucru contribuie la consistența informației ușurând dezvoltarea; se poate reconstrui toate celelalte structuri de date doar din bibliotecă și dintr-un fișier care listează erorile crawler-ului.

4.2.3 Indexul documentelor

Indexul documentelor păstrează informații despre fiecare document. Acesta este un index ISAM (Index sequential access mode) cu o lățime fixă, ordonat de un docID. Informația conținută de fiecare scurtă introducere include statutul curent al documentului, un indicator către bibliotecă, o evidență a documentului și statistici variate. Dacă documentul a fost parcurs atunci conține un indicator către un fișier cu multe variabile numit docinfo și care cuprinde URL-ul și titlul documentului. În caz contrar, indicatorul se îndreaptă către lista URL-urilor care cuprinde numai URL-uri. Această hotărare de design a fost

luată în conformitate cu dorința de a avea o structură compactă de date, precum și cu abilitatea de stabili un record de accesare unică a discului în timpul unei căutări.

În plus, există un fișier care este folosit în convertirea URL-urilor în docID-uri. Acesta conține o lista cu URL-uri împreună cu docID-ul corespunzător și este sortat de sumă de control. Pentru a găsi docID-ul unui anume URL, suma de control a URL-ului este calculată și o căutare binară este realizată pe fișierul de sume de control pentru identificarea docID-ului. URL-urile pot fi convertite în docID-uri luând mai multe simultan prin alipirea la acest fișier. Aceasta este tehnica pe care cel ce soluționează URL-uri o folosește pentru a transforma URL-urile în docID-uri. Această metodă de abordare este importantă pentru că altfel trebuie să efectuam o căutare pentru fiecare link care, ținând cont de disc, ar dura mai mult de o lună pentru baza de 322 milioane de link-uri.

4.2.4 Lexiconul

Lexiconul are mai multe forme diferite. O modificare esențială adusă, spre deosebire de sistemele anterioare, este ca lexiconul să poate ajusta memoria la un preț rezonabil. În implementarea actuală putem păstra lexiconul în memorie într-un sistem cu 256 MB de memorie. Lexiconul actual conține 14 milioane de cuvinte (deși câteva cuvinte rare nu au fost adăugate). Lexiconul are două parti: o listă de cuvinte (înșiruite împreună, dar separate de zerouri) și o serie de indicatori. Pentru a îndeplini diverse funcții, lista de cuvinte are informații auxiliare care se află însă dincolo de scopul acestei lucrări.

4.2.5 Listele de hit-uri (Hit Lists)

O listă de hit-uri corespunde unei liste de apariții ale unui anumit cuvânt într-un document, incluzând informații despre poziția, fontul și tipul de literă folosit. Listele de hit-uri explică cea mai mare parte a spațiului utilizat atât în indicele primar (forward index), cât și în indicele complementar (inverted index).

Din această cauză, este important să fie reprezentate cât mai eficient posibil. Am luat în calcul mai multe alternative pentru poziția de codificare, font și tipul de literă – codificarea simplă (un grup de trei numere inetgrale), codificarea compactă (o serie de biți optimizați manual) și codificarea Huffmann. În final, am ales codificarea compactă optimizată manual deoarece necesită de departe mai puțin spațiu decât codificarea simplă și mult mai puțină manipulare a biților decât codificarea Huffmann. Detalii despre aceste hit-uri sunt arătate.

Codificarea compactă folosește doi biți pentru fiecare hit. Există două tipuri de hit-uri: hit-uri complexe (fancy hits) și hit-uri simple (plain hits). Hit-urile complexe includ apariția hit-urilor într-un URL, titlu, textul link-ului sau meta tag. Hit-urile simple includ restul. Un hit simplu constă dintr-un bit referitor la tipul de literă, mărimea fontului și 12 biți de poziții ale cuvântului într-un document (toate pozițiile ce depășesc 4095 sunt catalogate 4096). Mărimea fontului este reprezentată relativ față de restul documentului utilizând 3 biți (doar 7 valori sunt de fapt folosite deoarece 111 este simbolul care semnaleză apariția unui hit complex). Un hit complex constă într-un bit referitor la tipul de literă, mărimea fontului este setată la 7 pentru a indica că este vorba de un hit complex, 4 biți pentru codificarea tipului de hit complex și 8 biți de poziție. Pentru hit-urile de tip anchor, cei 8 biți ai poziției sunt împărțiți în 4 biți pentru poziție în link și 4 biți pentru conținutul docID-ului în care link-ul apare. Aceasta ne oferă o sintagma redusă de căutare din moment ce nu există multe link-uri pentru un anumit cuvânt. Trebuie să se actualizeze metoda de stocare a hit-urilor anchor pentru permiterea unei rezoluții mai mari în cadrul poziției și câmpurilor de docID-uri. Se va folosi mărimea fontului în legătură cu restul documentului deoarece, atunci când căutam, nu dorim listarea diferită a unor documente identice doar pentru că unul din documente este scris cu un font mai mare.

Lungimea unei liste de hit-uri este stocată înainte chiar de hit-urile în sine. Pentru a economisi spațiu, lungimea listei de hit-uri este combinată cu wordID-ul în indexul primar și cu docID-ul în indexul complementar. Acest lucru o limitează la 8 și respectiv 5 biți (există o serie de trucuri care permit ca 8 biți să fie imprumutați din wordID). Dacă lungimea este mai mare și nu se poate încadra în respectivii biți, atunci un cod de rezervă este folosit în acești biți, iar următorii 2 biți vor conține lungimea actuală.

4.2.6 Indexul primar (Forward Index)

Indexul primar este deja parțial sortat și este stocat într-o serie de categorii (s-au folosit 64). Fiecare categorie conține o serie de wordID-uri. Dacă un document conține cuvinte care țin de un anumit barrel, docID-ul este înregistrat în categorie urmat de o listă de wordID-uri cu liste de hit-uri care corespund cuvintelor respective. Această schemă necesită mai mult spațiu de stocare din cauza docID-urilor duplicate, dar diferența este foarte mică pentru un număr considerabil de categorii și economisește timp și complexitate de codificare în faza finală de indexare făcută de sorter. Mergând mai departe, în loc de a stoca wordID-urile actuale, se stochează fiecare wordID ca o diferență relativă de la wordID-ul minim care se găsește în categoria în care se află și wordID-ul. Astfel, se pot folosi 24 biți pentru wordID-uri în categorii nesortate, lasând 8 biți pentru lungimea listelor de hit-uri.

4.2.7 Indexul complementar

Indexul complementar constă din aceleași categorii ca și indexul primar, cu diferența că acestea au fost procesate de sorter. Pentru fiecare wordID valid, lexiconul conține un indicator către categoria în care wordID-ul este inclus. Acest indicator se referă la o listă de docID-uri luate împreună cu listele de hit-uri corespunzătoare. Această listă reprezintă toate aparițiile aceluși cuvânt în toate documentele.

4.4 Indexarea Web-ului

— *Analiza*: orice sistem de analizare care este destinat să ruleze pe întreg Web-ul trebuie să facă față unei mari cantități de erori posibile. Acestea se pot întinde de la mici greșeli în sintagmele HTML până la kilobiți de zerouri în mijlocul unei sintagme, caractere non-ASCII, sintagme HTML depozitate în interior și o mare varietate de alte greșeli care pot provoca imaginația oricui pentru a-l determina să apară cu altele noi. Pentru viteza maximă, în loc să se folosească YACC pentru a genera un sistem CFG de analiză, se folosește un cablu electric pentru inițierea unui analizator lexical *which* care se dotează cu propria informație. Dezvoltarea acestui sistem de analiză care rulează cu o viteză rezonabilă și care este foarte solid a însemnat un timp îndelungat de muncă continuă.

— *Indexarea documentelor în categorii*: după ce fiecare document este analizat, este codificat într-o serie de categorii. Fiecare cuvânt este convertit într-un wordID utilizând un conținut de tip in-memory (lexiconul). Adăugările care se fac conținutului lexiconului sunt consemnate într-un fișier. Din momentul în care cuvintele sunt convertite în wordID-uri, aparițiile lor în documentul curent sunt traduse în liste de hit-uri și sunt scrise în categoriile primare. Principala dificultate cu simultaneitatea fazei de indexare este aceea că lexiconul trebuie împărțit. În loc de împărțirea lexiconului, se scrie un logaritm pentru toate cuvintele auxiliare care nu se aflau în lexiconul de bază, care s-a fixat la 14 milioane de cuvinte. În acest fel, indexer-ele multiple pot rula în paralel și apoi fișierul jurnal cu logaritmii cuvintelor auxiliare poate fi procesat de un singur indexer final.

— *Sortarea*: pentru a putea genera indexul complementar, sorter-ul preia fiecare categorie primară și o sortează după wordID pentru a produce o categorie complementară pentru titlu și hit-urile de tip anchor și o categorie complementară de text integral. Acest proces se întâmplă cu fiecare categorie în parte, necesitând astfel un spațiu mic de stocare temporară. De asemenea, se compară faza de sortare pentru a utiliza cât mai multe sisteme posibile prin

simpla rulare de sorter-e multiple, care pot procesa categorii diferite în același timp. Din moment ce categoriile nu sunt compatibile cu memoria principală, sorter-ul le subdivide mai departe în subcategorii care se potrivesc cu memoria bazată pe word ID și docID. Apoi, sorter-ul încarcă fiecare subcategorie în memorie, o sortează și îi transcrie conținutul în categorii complementare mici și în categorii complementare detaliate. O altă caracteristică importantă o reprezintă ordinea în care docID-urile trebuie să apară în listă. O soluție simplă este să fie stocate în funcție de docID. Acest lucru permite fuzionarea rapidă a diferitelor liste pentru întrebările cu multe cuvinte. O altă posibilitate este să fie stocate în funcție de apariția cuvântului în fiecare document. Acest lucru implică lipsa de valoare a răspunsurilor la întrebările alcătuite dintr-un singur cuvânt și face posibil faptul că răspunsurile la întrebările cu multe cuvinte să se afle chiar la început. Totuși, fuzionarea este cu mult mai dificilă. De asemenea, acest lucru face ca dezvoltarea să devină mult mai dificilă, în sensul că o schimbare în funcția de ordonare necesită o reconstruire a indexului. S-a ales un compromis între aceste opțiuni, prin păstrarea a două serii de categorii complementare - o serie pentru listele de hit-uri care include titlul sau link-ul și o altă serie pentru toate listele de hit-uri. În acest fel, se verifică primul set de categorii și dacă nu se găsesc destule potriviri în cadrul acestora se verifică și cele mai mari.

Rularea unui crawler web reprezintă o provocare. Există acțiuni riscante și subiecte sigure și, mult mai important, există subiecte sociale. Parcurgerea (crawling) reprezintă cea mai fragilă aplicație din moment ce implică interacțiunea cu sute de mii de servere web și nume de servere diverse care se află dincolo de posibilitatea de control a sistemului.

Pentru a pacurge sute de milioane de pagini web, ***Google are un sistem rapid (fast distributed crawling)***. Un singur server URL oferă liste de URL-uri unui număr de crawlers (în general se folosesc în jur de 3). Atât server-ul URL, cât și crawler-ele sunt realizate în Python. Fiecare crawler ține în jur de 300 de conexiuni (connections) deschise simultan. Acest lucru este necesar pentru regăsirea paginilor web la o viteză suficient de rapidă. La viteze mari sistemul poate să parcurgă peste 100 de pagini pe secundă utilizând 4 crawlere. Acesta se ridică la aproximativ 600K de date pe secundă. O acțiune importantă este reprezentată de verificarea DNS. Fiecare crawler menține un cache DNS propriu, astfel că nu este nevoie să se facă un control DNS înainte de parcurgerea fiecărui document. Fiecare dintre sutele de conexiuni se poate afla în stadii diverse: verificarea DNS, conectarea la gazdă, transmiterea solicitărilor și primirea răspunsurilor. Acești factori fac din crawler o componentă complexă a sistemului. Acesta folosește IO asincron pentru a face față solicitărilor și un număr de secvențe pentru mutarea preluărilor de pagini din secțiune în secțiune. Se adeverește astfel că rularea unui crawler care se conectează la mai mult de jumătate de milion de servere și care generează zeci de milioane de fișiere jurnal implică o cantitate considerabilă de e-mailuri și apeluri telefonice. Datorită numărului mare de persoane care sunt online, există întotdeauna aceia care nu

știu ce este un crawler deoarece acesta este primul pe care îl văd. Aproape în fiecare zi primim un mesaj de genul 'V-ați uitat la foarte multe pagini din site-ul meu. Ce parere aveți despre el?'. Există de asemenea unele persoane care nu au aflat de existența protocolului de excluziune la roboți și se consideră că pagina lor ar trebui protejată împotriva unei indexari de genul 'Această pagina este protejată și nu va fi indexată', care este dificil de înțeles de crawler-ele web. De asemenea, din cauza cantității uriașe de date implicate, lucruri neașteptate se pot întâmpla.

De exemplu, sistemul a încercat să parcurgă un joc online.

Rezultatul a constat într-o mulțime de mesaje iritante chiar în mijlocul jocului. Se pare că aceasta a fost o problemă ușor de rezolvat. Dar această problemă nu aparuse până în momentul în care s-au descărcat zeci de milioane de pagini. Datorită variației ridicate în paginile web și în servere, este practic imposibil să testezi un crawler fără să-l rulezi pe o parte considerabilă a Internetului. Invariabil, apar sute de probleme obscure care se pot ivi pe o singură pagină din tot web-ul și pot cauza distrugerea crawler-ului sau mai rău, poate cauza o reacție imprevizibilă sau incorectă. Sistemele care accesează părți mari din Internet trebuie să fie foarte solide și testate cu multă atenție. Din moment ce sistemele complexe cum sunt crawler-ele vor duce în mod invariabil la apariția problemelor, trebuie să existe resurse semnificative dedicate citirii de e-mail-uri și rezolvării problemelor din momentul în care acestea apar.

4.5 Căutarea

Scopul căutării este acela de a oferi rezultate concludente în timp util. Multe dintre motoarele de căutare comerciale par să fi făcut progrese considerabile din punct de vedere al eficienței. De aceea, se pune accentul mai mult pe calitate.

Pentru marcarea unei limite a timpului de răspuns, odată ce un anumit număr de documente care se potrivesc cu interogarea (40.000 de obicei) este găsit, cel care a inițiat căutarea poate accesa paginile.

Aceasta înseamnă că este posibil ca rezultate neconcludente să fie oferite în schimb. În prezent, se investighează alte metode pentru rezolvarea acestei probleme. În trecut, s-au sortat hit-urile în concordanță cu PageRank, lucru care pare să fi îmbunătățit situația.

4.5.1 Sistemul de clasificare

Google păstrează mult mai multe informații despre documentele web decât motoarele tipice de căutare. Fiecare listă de hit-uri include poziția, fontul și informații despre tipul de literă folosit. În plus, se iau în calcul hit-urile după textul link-ului și PageRank-ul documentului. Combinarea tuturor acestor

informații într-un singur rezultat este dificilă. S-a conceput funcția de ordonare astfel încât nici un factor particular să nu aibă o influență prea mare. Se iamai întâi cazul cel mai simplu - o interogare cu un singur cuvânt.

Pentru afișarea unui document folosind o interogare cu un singur cuvânt, Google parcurge toate listele de hit-uri ale documentului pentru cuvântul respectiv. Google consideră fiecare hit ca aparținând unuia dintre diversele tipuri (titlu, link, URL, fonturi mari și fonturi mici de text simplu etc.), fiecare dintre acestea având grade diferite de importanță în funcție de tipul din care face parte. Aceste grade de importanță formează un vector indexat în funcție de tip. Google numără hit-urile fiecărui tip din lista de hit-uri. Apoi fiecare poziție este reorganizată într-un clasament în funcție de importanță. Gradele de importanță cresc liniar în funcție de primele poziții, dar se reduc repede astfel încât este relevant numai un anumit număr de apariții. Se preia produsul scalar al vectorului de ponderi de apariții împreună cu vectorul de ponderi de tipuri pentru a calcula un scor IR al documentului. În final, scorul IR este combinat cu PageRank pentru a oferi un rezultat final al documentului.

Pentru o interogare alcătuită din mai multe cuvinte, situația este și mai complicată. În acest caz, listele multiple de hit-uri trebuie parcurse simultan astfel încât hit-urile care sunt apropiate într-un document sunt plasate pe poziții superioare față de cele care sunt departate unele de altele. Hit-urile din listele multiple sunt potrivite astfel încât hit-urile apropiate sunt puse împreună. Pentru fiecare set de potriviri de hit-uri, se calculează o apropiere. Această apropiere se bazează pe cât de departate sunt hit-urile în cadrul documentului (sau link-ului), dar este clasificată în 10 clase cu valori diferite, mergând de la o sintagmă apropiată până la 'nu foarte aproape'. Se fac contorizări nu numai pentru fiecare tip de hit, dar și pentru fiecare tip și apropiere. Fiecare pereche de tip și apropiere are o pondere tip-apropiere. Contorizările sunt clasificate în funcție de ponderile de apariții și se preia produsul scalar pentru ponderile de apariții și ponderile de tip-apropiere pentru realizarea unui scor IR. Toate aceste numere și matrice pot fi afișate odată cu rezultatele căutării folosind o metodă specială de corectare. Toate aceste afișări sunt de foarte mare ajutor în dezvoltarea sistemului de ordonare.

4.5.2 Scurta recapitulare

Funcția de ordonare are mulți parametri precum ponderile de tip și ponderile tip-apropiere. Stabilirea valorilor exacte pentru acești parametri ține de domeniul magiei. Pentru a face acest lucru, există un mecanism de feedback al utilizatorului în cadrul motorului de căutare. Un utilizator încrezător poate evalua opțional toate rezultatele care îi sunt returnate. Acest feedback este salvat. Apoi, când se modifică funcția de ordonare, se observă impactul acestei schimbări asupra căutarilor anterioare care au fost ordonate. Deși departe de

perfecțiune, acest lucru ne face sa observăm cum o schimbare în funcția de ordonare poate afecta rezultatele căutării.

5. Rezultate și mod de funcționare

Cea mai importantă caracteristică a unui motor de căutare este calitatea rezultatelor căutării. În timp ce o evaluare completă a utilizatorului se află dincolo de scopul acestei lucrări, experiența cu Google a dovedit că este posibil să se obțină rezultate mult mai bune pentru majoritatea căutărilor decât ale celor mai importante motoare comerciale de căutare. Dacă se caută 'bill clinton' se obțin rezultate ce evidențiază câteva dintre trăsăturile Google. Rezultatele sunt grupate de server. Acest lucru este de mare ajutor atunci când se examinează seturile de rezultate. O serie de rezultate sunt din domeniul whitehouse.gov, care reprezintă rezultatul așteptat pentru o căutare de acest gen. În mod obișnuit, majoritatea motoarelor de căutare comerciale nu oferă nici un rezultat din whitehouse.gov, cu atât mai puțin pe cele corecte. Trebuie observat că nici un titlu nu este afișat pentru primul rezultat, din cauză că nu a fost parcurs. În schimb, Google s-a bazat pe textul link-ului pentru a hotărî dacă acesta era un răspuns potrivit pentru întrebare. În mod asemănător, al cincilea rezultat este o adresă de e-mail, care, bineînțeles, nu poate fi parcursă (crawlable) și care este, de asemenea, un rezultat al textului link-ului.

Toate rezultatele sunt pagini cu un înalt nivel calitativ și, la o ultimă verificare, nici una dintre ele nu este un link irelevant (broken). Acest lucru se întâmplă din cauză că toate au un PageRank ridicat. PageRank-urile sunt marcate prin procentajele de culoare roșie. În concluzie, nu sunt rezultate despre un alt Bill în afara de Clinton sau despre un alt Clinton în afară de Bill și asta pentru că se acordă o importanță majoră aproximării sintagmelor în care apare un cuvânt. Bineînțeles că un adevărat test al calității pentru un motor de căutare ar trebui să implice și un studiu complet despre utilizator sau o analiză a rezultatelor.

5.1 Parametri pentru stocare

În afară de calitatea căutării, Google este proiectat să măsoare capacitatea Web-ului pe măsură ce acesta se dezvoltă. Un aspect al acestui fenomen îl reprezintă modul eficient de stocare. Datorită comprimării, mărimea totală a bibliotecii este de 53 GB, puțin peste o treime din totalul de date pe care le conține. La prețurile actuale ale discurilor biblioteca poate fi considerată o sursă ieftină de date. Ceea ce este mai important - totalul datelor folosite de motorul de căutare necesită un spațiu comparabil de stocare, adică în jur de 55 GB. Mai departe, un răspuns poate fi găsit pentru majoritatea interogarilor prin folosirea indexului complementar.

5.2 Modul de funcționare a sistemului

Este important pentru un motor de căutare să parcurgă și să indexeze eficient. Astfel, informația poate fi permanent actualizată și modificările majore aduse sistemului pot fi testate relativ repede. Pentru Google, operațiunile importante sunt **Crawling (parcurgerea)**, **Indexing (indexarea)** și **Sorting (sortarea)**. Este dificil de măsurat cât a durat crawling-ul în total din cauză că discurile au fost în întregime completate, numele serverelor nu mai sunt funcționale sau din cauza oricărei probleme care putea determina oprirea sistemului. În total a durat cam 9 zile să se descarce cele 26 milioane de pagini (inclusiv erorile). Totuși, din moment ce sistemul rula fluent și mai repede, s-au descărcat ultimele 11 milioane de pagini în doar 63 de ore, în medie 4 milioane de pagini pe zi sau 48,5 pagini pe secundă. S-au rulat indexer și crawler simultan. Indexer-ul a rulat mai repede decât crawler-ele. Aceasta s-a întâmplat din cauză că s-a petrecut destul timp pentru optimizarea indexer-ului astfel încât să nu determine întârzierea dezvoltării sistemului. Aceste optimizări au inclus o serie de actualizări ale indexului documentului și plasarea structurilor precare de date pe un disc local. **Indexer-ul rulează aproximativ 54 de pagini pe secundă. Sorter-ele pot fi rulate complet în paralel, utilizând 4 sisteme, întreg procesul de sortare durând în jur de 24 de ore.**

5.3 Procesul de căutare

Îmbunătățirea procesului de căutare nu a reprezentat punctul principal de interes al cercetării al realizatorilor Google-ului până în acest moment. Versiunea curentă Google răspunde la cele mai multe din interogări între 1 și 10 secunde. Acest interval este dominat în mare parte de IO-ul discului asupra NFS (din moment ce discurile sunt împărțite pe sisteme diferite). Mai departe, Google nu are nici un fel de optimizare cum ar fi cache de căutări, subindici pe termeni comuni și alte optimizări obișnuite. *Intenția lor este de a mări considerabil viteza Google prin distribuție și hardware, software și îmbunătățiri algoritmice.* Scopul este acela de a putea răspunde la mai multe sute de interogări pe secundă.

6. Concluzii

Google este proiectat să fie un motor de căutare scalabil. Scopul principal este acela de a oferi rezultate de calitate pe fondul dezvoltării rapide a World Wide Web. Google folosește o serie de tehnici pentru ameliorarea calității căutării incluzând page rank, textul link-ului și alte informații apropiate. Mai departe, Google reprezintă o arhitectură completă pentru adunarea paginilor web, indexarea lor și efectuarea de interogări asupra lor.

6.1 Activități viitoare

Un motor amplu de căutare este un sistem complex care rămâne deschis completărilor. Scopurile noastre imediate sunt acelea de a îmbunătăți căutarea și de a afișa aproximativ 100 de milioane de pagini web. Câteva îmbunătățiri simple aduse eficienței includ cache de interogări, alocare optimizată a discului și subindici. O altă secțiune care necesita o cercetare amanunțită este reprezentată de updatări. Trebuie să avem algoritmi potriviți pentru a decide ce pagini mai vechi ar trebui reparcurse și ce pagini noi ar trebui parcurse. O zonă promițătoare de cercetare este aceea de a utiliza un *cache proxy* pentru construirea de baze de date cu scopul căutării din moment ce acestea sunt influențate de cerere. Se plănuiește adăugarea de trăsături simple ce sunt suportate de motoarele comerciale de căutare cum ar fi operatorii de tip boolean, negația și implicația. Totuși, alte trăsături cum ar fi feedback-ul relevanței și clustere (Google suportă în mod obișnuit clusterelor bazate pe numele de domeniu) de-abia acum încep să fie explorate. Se intenționează de asemenea susținerea contextul de utilizator (cum ar fi locația utilizatorului) și rezumarea rezultatelor. Apare și preocuparea extinderii utilizării structurii de link-uri și a textului link-ului. Experimentele simple indică faptul că PageRank poate fi personalizat prin amplificarea ponderii homepage-ului unui utilizator sau a bookmark-urilor acestuia. În ceea ce privește textul link-ului, se experimentează utilizarea textului link-urilor învecinate alături de textul link-ului în sine. Un motor de căutare pe Web reprezintă un mediu bogat pentru inițiativele de cercetare.

6.2 Căutare la standarde înalte

Cea mai mare problemă cu care se confruntă astăzi utilizatorii de motoare de căutare o reprezintă calitatea rezultatelor pe care le primesc. Pe când rezultatele sunt deseori amuzante și largesc orizontul utilizatorului, ele pot deveni și frustrante și pot consuma timp pretios. De exemplu, primul rezultat pentru o căutare 'Bill Clinton' pe unul dintre cele mai populare motoare de

căutare comerciale a fost Bill Clinton Joke of the Day: April 14, 1997. Google este destinat să ofere rezultate de o calitate superioară astfel încat Web-ul să continue să se dezvolte rapid, iar informația să poata fi găsită ușor. Pentru a putea realiza acest lucru, Google utilizează frecvent informația hipertextuală ce constă din structura de link-uri și din textul link-urilor. Google folosește de asemenea aproximarea și informația despre fonturi. În timp ce evaluarea unui motor de căutare este dificilă, s-a aflat în mod intuitiv că Google returnează rezultate la un înalt nivel calitativ spre deosebire de motoarele comerciale de căutare. Analizarea structurii de link-uri prin PageRank permite Google să evalueze calitatea paginilor web. Utilizarea textului link-ului ca o descriere a ceea ce indică link-ul contribuie la relevanță și, într-o anumită măsură, la înaltul standard calitativ al rezultatelor. În cele din urmă, utilizarea unor informații asemănătoare ajută la mărirea gradului de relevanță al multor interogări.

6.3 Arhitectura scalabilă

În afară de calitatea căutării, Google este destinat să măsoare. Trebuie să dea dovadă de eficiență atât în spațiu, cât și în timp, iar factorii constanți sunt foarte importanți atunci când vorbim de Web. Prin implementarea Google, s-au observat întârzieri în CPU, accesarea memoriei, capacitatea memoriei, căutările pe disc, conținutul discului, capacitatea discului și IO de rețea. Google a evoluat pentru a depăși o serie din aceste probleme prin intermediul unor operațiuni variate. Structurile majore de date ale Google utilizează în mod eficient spațiul de stocare. Mai departe, operațiunile de parcurgere, indexare și sortare sunt suficient de eficiente pentru a construi un index al unei părți substanțiale a Web-ului – 24 de milioane de pagini în mai puțin de o săptămână.

Referințe

Google Guide (2003) “How Google Works” *Google Guide* Retrieved 8 Sept, 2004 from: googleguide.com

Sullivan, R., (2004) “Google and Googlebot Information” *searchengineposition.com* Retrieved 8 Sept, 2004 from: searchengineposition.com

Brin, S., Page, L. (2000) “The Anatomy of a Large-Scale Hypertextual Web Search Engine” Retrieved 8 Sept, 2004 from: stanford.edu

Barroso, L. A., Dean, J., Hölzle, U. (2003) “Web search for a planet: The Google cluster architecture” *Micro, IEEE*, Volume: 23 , Issue: 2, pp.22-28

Ghemawat, S., Gobioff, H., and Leung, S. T., (2003) “The Google File System” *SOSP'03*, October 19–22

Traducere de pe pagina: http://wiki.media-culture.org.au/index.php/Google-How_It_Works.