

is megérdemli, hogy a szó különös értéket jelző értelmében műhelynek nevez-
zék. A kötetek - jól sikerült - "esztétikai gondozása" (borítólap, belső illusztrációk) Gonda Zoltánra hárult.

FUTALA Tibor

K. Leon MONTGOMERY: Document Retrieval Systems.
Factors Affecting Search Time. Marcel Dekker, Inc.
New York. 1975. VIII. 144 p.

DOKUMENTUMVISSZAKERESŐ RENDSZEREK. A KERESÉSI IDŐT BEFOLYÁSOLÓ TÉNYEZŐK

A visszakereső rendszerek hatékonyságának kérdése napjainkban élő probléma. A Pittsburgh-i Egyetemen (Pennsylvania) kísérletsorozatot folytattak, amelynek eredményeit a jelen - könyv formában publikált - kutatási jelentésben teszik közzé.

A géppel olvasható adattárakban való kereséssel kapcsolatos problémák a hatvanas évek legelején vetődtek fel. Ezek a gépek azonban még meglehetősen kis központi memóriával rendelkeztek, a háttértár gyakran néhány mágnesszalagot vagy egy mágneslemezt jelentett csupán. 1965-ben megjelentek a harmadik generációs gépek, amelyek lehetőségeket nyitottak meg az információ visszakeresés terén is. Három igen lényeges kérdésben hoztak újat, kettőt ebből a hardware, egyet software területén: Hardware tekintetében gyorsabb processzort alkalmaztak és jelentősen megnövelték a belső és külső táruk kapacitását. A harmadik jelentős ujitást az operációs rendszeren hajtották végre.

Már jóideje köztudott, hogy a visszakeresési időt nagy mértékben befolyásolja az adattár nagysága és szervezése. A kísérlet végrehajtói vállalkoztak arra, hogy bemérik a keresési időt befolyásoló összes egyéb paramétert, a költségtényezők függvényében, így megvizsgálták a file szervezési módszereket, a kereső programokat, és a rendelkezésre álló számítógép rendszert. A géppel olvasható adattárak használata csak akkor gazdaságos, ha lehetőleg minél kisebb költséggel minél több kérdést tudnak megválaszolni, minél több felhasználót tudnak kiszolgálni.

Vizsgálataikat a file szervezés tanulmányozásával kezdték. A file szervezésnek két alap formája ismeretes: a lineáris és az invertált. Lineáris file szervezés esetén a dokumentumok rekordjait az azonosító számok sorrendjében tárolják. Minden dokumentumot jellemeznek index kifejezésekkel, és ezeket a dokumentumhoz rendelik. Ez a következőképpen érzékeltethető:

Ha az

1. sz. dokumentumot A, B, F és G index kifejezések jellemeznék és a
2. sz. dokumentumot B, G, és H az adattár szerkezete a következő:
 1. sz. dokumentum A, B, F, G index fogalmak
 2. sz. dokumentum B, G, H index fogalmak.

Invertált file szervezésekor az index fogalmakat rendezik sorba és minden indexfogalom után felsorolják azokat a dokumentumokat, amelyekben a kérdéses indexfogalom előfordul. Az előbbi példánál maradva, az invertált file a következőképpen épül fel:

A indexfogalom	1. sz. Dokumentum
B indexfogalom	1. sz. Dokumentum, 2. sz. Dokumentum
F indexfogalom	1. sz. Dokumentum
G indexfogalom	1. sz. Dokumentum, 2. sz. Dokumentum
H indexfogalom	2. sz. Dokumentum

A lineáris file szervezés lényegesen könnyebben megvalósítható, igen egyszerű az újabb adatok beépítése (update-olás), könnyen követhető a keresés logikája, és nem utolsó sorban batch üzemmódban igen gazdaságos. Valószínűleg ez az oka annak, hogy ez a file szervezési módszer terjedt el szélesebb körben. (NASA Technology Transfer Project, Chemical Abstracts Serv.)

Minél több kérdésre keres egyidőben azonos file-on, annál kisebb az egy kérdésre eső költség.

Megfigyelték, hogy azonos nagyságu file esetén a keresési időt nagyon befolyásolja a keresési profilok száma és a két említett file szervezés módszer, ezért a kísérletek során kidolgoztak egy olyan speciális kereső programot, amellyel a lineáris és invertált szervezésű file-on egyaránt lehet visszakeresni. A programot a Pittsburgh-i Egyetem IBM 360/50-es számítógépére irták.

A munkálatokat a kérdés szakirodalmi felderítésével kezdték, az 1960-as évektől szemleszerűen feldolgozott anyagot a 2. sz. fejezetben ismertetik.

Az irodalomban tárgyalt kérdéseket végül három csoportba rendezték: vizsgálták a keresési időt:

- az egyidejűleg futó kérdések számának függvényében,
- az adatbázis méretének függvényében
- és az alkalmazott file szervezés tekintetében.

Ugy szervezték tehát a kísérleteket, hogy a fenti kérdésekre választ kapjanak, ezért a következő feltételeket kellett biztosítani a kísérletek lefolytatásához:

1. az adatbázis létező, széles körben ismert és használt legyen,
2. aktiv, tényleges, valós kérdésekkel fusson a program,
3. az adatbázis lineáris és invertált file szervezést is reprezentáljon vagy ilyenné átalakítható legyen,
4. és végül, hogy az futtatható legyen az Egyetem adott gépein.

A kísérletek szempontjából két gépi adattár jöhetett szóba a National Aeronautics and Space Administration (NASA) Technology Transfer Project szalagjai vagy a National Science Foundation (NSF) támogatásával a Chemical

Information Center által kiadott Chemical Abstracts Services Condensates adattára. Némi vizsgálódás után végül is a CAS szalagot választották, mivel az Egyetemen ezt az adattárat az IBM 360/50-es gépen futtatják, amely nagyobb kapacitású, mint a NASA adattárral dolgozó másik számítógép. Az IBM gép lehetőséget nyújtott a multiprogramozási lehetőségek tanulmányozására is.

A kísérlet leglényegesebb alkotó elemei: az adatbázis, a kereső profilok, a kereső program és a számítógépes rendszer. Az általuk használt adatbázis, amely 74. volume 2. issue szalagváltozata öt adatszoportot tartalmaz:

1. A dokumentum, a kötetszám és füzetszám,
2. A cím
3. A szerző
4. A szerzővel kapcsolatos adatok
5. Az index fogalmak.

Ezek közül csak három adatszoportot használtak: 1. A dokumentum (kötet és füzetszám); 2. A címet; 3. Az index fogalmakat.

A cím és az indexfogalmak alapján végzik a keresést, és határozzák meg a kötet és füzetszámot. Általában 10-50 fogalom állt rendelkezésre a dokumentumkeresésére és kb. 20 terminussal kerestek minden dokumentumra. Az eredeti szalagon (adatbázison) az adatok lineárisan vannak szervezve. A kérdéseket a Chemical Information Center Project kb. 400 kérdéséből 180-at választottak ki, majd végül 128 kérdéssel végezték a kísérletet, Boole logikai elemek (és, vagy, és nem) alkalmazásával. Az IBM 360/50 operációs rendszerre lineárisan szervezett file-ra íródott, módosíthatók volna ezt is, de végül is egy új, az invertált szervezésű file-on való keresést lehetővé tevő programot írtak.

A kísérlet során két módszert alkalmaztak: először megvizsgálták a keresési idő és a keresett dokumentumok viszonyát, majd a kérdések számának függvényében vizsgálták a keresési időt. Természetesen mindkét variációt: lineáris és invertált szervezésű file-on is kipróbálták, az eredményeket összehasonlították. A dokumentumok számának és a keresési időnek vizsgálatakor 1, 256, 512, 1024, 2048 és 4096 dokumentumot kerestek vissza, lineáris és invertált file-on egyaránt. A kérdések számát a kísérlet második felében pedig a következőképpen változtatták: 1, 16, 32, 64, 96 és 128 kérdést tettek fel. Amikor lehetőség nyílt, háromszor megismételték a futtatást, hogy így átlag értékekhez jussanak.

A kísérlethez ún. bimoduláris programot dolgoztak ki, amely egyaránt használható invertált és lineárisan szervezett adattárak esetében. Olyan programot kellett írni tehát, amelyben lehetőleg azonos szubrutinokat alkalmaznak és csak, ahol feltétlenül szükséges, tér el egymástól a visszakereső program, mert egyébként a különböző program lépésekből, megoldásokból is származnak időeltolódások.

Részletesen ismertetik, a számítógép típusát, lényeges és jellemző adatait, a konfigurációkat, a program minden jellemzőjét, a használt kifejezéseket, az index fogalmakat és minden lényeges adatot.

A kísérleti futtatások során a következő adatokat regisztrálták automatikusan:

- a lineárisan szervezett file-on vagy az invertált file-on kerestek,
- az átnézett dokumentumok száma (1-4096)
- hány kérdést tettek fel (1-128)
- a futtatás pontos ideje és napja
- hány kereső profilt használtak
- a futtatás száma
- a központi memória időigénye
- input és output idők
- teljes számítógépes idő
- a futtatás költsége
- az output sorok száma
- a keresésre előkészített file invertálás ideje
- a 4096 dokumentum szalagról lemezre másolásának ideje
- visszakeresési idő
- a sikeresen visszakeresett dokumentumok száma
- és hogy mennyi idő kell az output perifériák előkészítésére.

Ezeket az adatokat egy 80 oszlopos lyukkártyára vitte fel a gép, minden egyes futtatásról készítettek ilyen kártyát, összesen 321-et, a kísérlet hat hete alatt.

A számadatokból szerkesztett ábrák és diagramok alapján a következő megállapításokat szűrték le: a visszakeresési idő egyenes arányban változik a dokumentumok számával, kezdetben. Kisebb adatbázisnál lineárisan szervezett adattárban érhető el a legrövidebb idő alatt a kérdéses dokumentum. Nagyobb mennyiségű adat kezelésére az invertált szervezés ajánlott. Arra, hogy hol van ez a metszéspont, amelyből megállapítható, hogy hol gazdaságos egyik és honnan már a másik szervezési módszer - az egyidőben futó kereső profilok száma van nagy hatással. Minél nagyobb számú kérdés fut egyidőben, annál kisebb adatbázis esetében az invertált file szervezés eredményesebb. Ha az egyidőben futtatott kereső profilok száma növekszik, akkor már kisebb adatbázis esetében is előnyösebb az invertált file alkalmazása.

A kísérlet során megvizsgálták a dokumentumok számának és a kérdések számának viszonyát is. Adott számú dokumentumbázis esetén az egyidejűleg futó kérdések számának növekedésekor elérkezünk egy olyan kérdésszámhoz, amelynél nagyobb számú kérdés esetén már az invertált szervezésű adattárnál is kisebb az egy-egy kérdésre jutó idő.

A vizsgálatok egyik megállapítása, hogy a visszakeresést célszerű batch üzemmódban végezni, így a legnagyobb a hatékonyság. Ugy találták, hogy 64 kereső profil egyidejű használatakor kapták a legrövidebb egységnyi keresési időt, ha ennél többet kerestek egyidejűleg az egységnyi keresési idő növekedett.

Érdekes megállapítás az is, hogy célszerű az invertált szervezésű file-lal dolgozni, ha 32-nél több kérdés fut egyidőben, és ha az adott adatbázis nagysága meghaladja az 512 dokumentumot.

A kísérletek során számos értékes eredmény és megállapítás született, de egyben újabb kérdések is felvetődtek. Feltétlenül meg kell vizsgálni a közeljövőben a fileszervezés egyéb más megoldásait, lehetőségeit. A munkálatok fő eredményének a bimoduláris STRESS program kialakítása tekinthető.

BOBOKNÉ BELÁNYI Beáta