

Diss. B 720

Makay Árpád

Tezauruszok alkalmazása számítógépes  
információs rendszerekben

S z e g e d

1971.



Diss. B 720



## B e v e z e t é s

A könyvtári és dokumentációs munka három fő területén kínálkozik alkalom a műveletek gépesítésére, automatizálására:

1. A publikálás műveletének automatizálását szolgálják a korszerű sokszorosító berendezések, lyukszalag-, mágnesszalag-szedőgépek. Felhasználhatók számítógépes rendszerek is: írógépből és vizuális megjelenítő készülékből álló, a számítógéppel on-line kapcsolatban lévő terminálokon gépelhető a publikáció. Az így nyert információkat közvetlenül a számítógép szerkeszti, kiadványokat állít össze, esetleg vezérli a sokszorosító berendezést. A kiadványok automatikus előállítása mellett a rendszer másik előnye, hogy a publikáció számítógép számára közvetlenül hozzáférhető marad bármilyen egyéb feldolgozási célokra is.

2. A publikáció előállításának automatizálása azt jelenti, hogy meglévő információk alapján a számítógép készít újabb kiadványokat. Elsősorban bibliográfiák, kötetkatalogusok, névmutatók, tárgyszóindexek előállítására gondolhatunk. Ebbe a kategóriába sorolható az automatikus tartalmi kivonatok készítése is: a publikáció teljes szövegéből egy /általában lényegesen rövidebb/ szólánc előállítása, amely tartalmi szempontból a "lényeg", illetve a lényegre utaló kifejezéseket tartalmazza.

3. A három terület közül az információkeresés (haszná-

latosabb terminológiával élve, az információvisszakeresés) automatizálása a legösszetettebb. Sajátos problémái mellett szoros kapcsolatai vannak az automatikus tartalmi kivonatolással, tárgyszó-indexeléssel. Lényegében arról van szó, hogy publikációk egy halmazából tartalmi vagy formai ismérvek alapján a számítógép válogasson és a kiválasztott publikációkat teljes egészükben vagy csak a rájuk vonatkozó utalások (azonosító szám, lelőhely, stb.) formájában az igénylő rendelkezésére bocsássa.

A tudományos és technikai publikációk mennyiségi növekedése indokolja az információkeresés automatizálásának szükségességét. Másrészt az információkeresés ma már főként tartalmi szempontok alapján elképzelhető, ezért tárolni csak tartalmilag feltárt és tartalmi ismérvek alapján könnyen hozzáférhető dokumentumot érdemes. Az információszolgáltatásnak pedig egyéni igényeket kell tudni kielégíteni, szemben a bibliográfiák, katalógusok, indexek (általában a kiadványok) rendszerével, amelyeket tömegkommunikációs eszközöknek tekinthetünk.

A publikációk tartalmi feltérésének /22/ művelete, sőt magának a fogalomnak a jelentése is igen összetett. A publikáció egésze, mint információhordozó, bármilyen átdolgozás, kivonatolás, rövidítés folyamán veszít információ-tartalmából. Mégis lehetnek olyan szempontjaink, amelyek alapján a publikáció által szolgáltatott információk egy részét lényegtelennek, más részét pedig lényegesnek tekintjük. Az utóbbi rész valamilyen formában történő megadását nevezzük tartalmi feltérésnek.

Információkereső rendszerekben a "lényeg" megadásának szempontjait a következőképpen jellemezhetnénk. Adva van bizonyos típusu publikációknak (pl. számítástechnikával kapcsolatos folyóirat-cikkeknek) egy halmaza. Ugyanakkor elég jól körülhatárolható az az embercsoport (a felhasználói kör), amelynek tagjait érdekli a publikációk tartalma. Vizsgáljuk meg, milyen "szokásaik" vannak a felhasználóknak akkor, amikor valamilyen őket érdeklő téma után kutatnak a publikáció-halmazban. Ma a járható út általában az, hogy néhány folyóiratot időrendben visszafelé lepozva átfutjuk a címeket, összefoglalásokat, így egy idő múlva összegyűlik a számunkra elegendő "olvasnivaló". Ha még a már megtalált publikációk irodalomjegyzékét is figyelembe vesszük (azaz felhasználjuk a cikk írója korábbi irodalom-kutatásának eredményét), meggyorsíthatjuk a munkát. Nem sok a munkára fordítandó időnyereség akkor sem, ha a könyvtár állományába tartozó szakember végzi ezt helyettünk: megjelenik azonban a probléma, hogyan jelöljük meg számára a minket érdeklő témát, hogy valóban azt kapjuk, amire szükségünk van. Ugyanez a nehézség merül fel akkor is, ha a keresést számítógép végzi.

Nyilvánvalóan fel kell sorolni a keresett témára jellemző ismérveket: ezek általában a témával kapcsolatos fogalmakat jelölő szavak, kifejezések (kulcsszavak, deskriptorok). Kérhetjük ezután azokat a publikációkat, amelyek két vagy több fogalom mindegyikével foglalkoznak; vagy egy fogalommal kapcsolatosak, de egy másikkal nem; stb.

A keresés folyamán most már az elsőrendű feladat, hogy megállapítsuk: egy adott publikáció kapcsolatos-e egy adott kifejezéssel megjelölt fogalommal. Nyilvánvalóan minden egyes felhasználó megjelenésekor végigolvasni minden egyes publikációt a publikáció-halmazból, és így megállapítani a megfelelőségét, igen időigényes volna (emberi erővel egyáltalán lehetetlen). A tartalmi feltárás feladata információkereső rendszerben tehát éppen az, hogy hozzárendelje a publikációhoz mindazokat a fogalmakat, amelyek alapján a felhasználói kör bármelyik tagja később keresheti az adott publikációt. Számítógépes rendszerben ez még további feltételeket is szab: a fogalmakat kifejezésekkel írjuk le, így a publikációhoz azokat a kifejezéseket (deskriptorokat) kell rendelni, amelyek alapján a felhasználók keresni fogják. A későbbiekben a tartalmi feltárás eredményének mindig egy kifejezés-sorozatot (deskriptor-sorozatot) tekintünk.

Az alábbiakban információvisszakereső számítógépi rendszerek két aktuális problémájának megoldására szolgáló, lista-eljárásokkal megvalósítható kérdéséről lesz szó. Az egyik a dokumentumok tartalmi feltárásának (legalábbis részben) automatizálását tűzi ki célul, míg a másik az információs rendszerben tárolt, tartalmilag feltárt dokumentum-halmaz igényelt (speciális tartalmu) részének megjelölését (az igény megfogalmazását) célozza.

A dokumentumok tartalmi feldolgozásának teljes automatizálását részben a természetes nyelvekre vonatkozó grammatikai és szemantikai ismereteink jelenlegi hiányosságai, részben a dokumentumok terjedelme (ami a számítógép számára hozzá-

férhető adathordozóra való konvertálásának időigényességét jelenti) akadályozzák. A "tartalmi kivonatolásban" ezért részmegoldásokra szorítkozunk. Még sokáig szakemberre lesz szükségünk, aki a dokumentumot tartalmilag értékeli, jellemzőit összegyűjti. Egyszerűsíthető azonban a forma, amelyben a számítógéppel az ismérveket közli. Az ismertetésre kerülő módszer egyik előnye pl. a következőkben foglalható össze. A dokumentum feldolgozását végző szakember a tartalmat szavakkal, szócsoportokkal (tárgyszavak ill. deskriptorok) írja le. Ezek közismert, de általában az adott fogalomra nem kizárólagosan használt kifejezések: velük párhuzamosan élnek a fogalom egyéb megjelölései is. Ha a dokumentum csupán ezekkel a deskriptorokkal kerül az információs rendszerbe, akkor a felhasználó is csak ezeket használhatja igényének megfogalmazására, másként nem juthat el az adott dokumentumhoz. A tartalmi feltárást végző ill. a felhasználó szakembernek, "egy nyelvet" kell tehát használni. Ezt segíti elő az un. tezaurusz: egy adott szakterület releváns kifejezéseinek gyűjteménye, mindegyik mellett egyebek között jelölve, hogy az illető kifejezés használható-e deskriptorként, illetve ha nem, mit kell helyette használni. Kézenfekvő, hogy a tezaurusz alapján a számítógép hozza közös nevezőre a feldolgozót és a felhasználót, egyszerűsítve a tartalmi feltárás és az igények megfogalmazásának munkáját.

A tezaurusz használatát még egy lényeges szempont indokolja. A tartalmi feltárás során a dokumentumhoz rendelhetőek olyan kifejezések, amelyek tulságosan "szűk", speciális

fogalmakat jelölnek: ekkor az adott kifejezés mellé fel kell venni egy vagy több olyan deszkriptort, amely bővebb kategóriát jelent. A visszakeresés így történhet akár a szűkebb, akár a bővebb kategóriák alapján. Természetesen ehhez az szükséges, hogy a teaurusz egy adott kifejezés mellett felsorolja a genetikus fogalmakat is.

Ez a "homogenizáló" fázis akkor sem hagyható el, ha eljutunk a tartalmi feldolgozás teljes automatizálásáig. A dokumentum szövegében illetve a kivonatban a tartalomra nézve releváns kifejezések (a gépi kivonatolás eddigi tapasztalatai szerint) nagy számát mindenképpen csökkenteni kell egy információvisszakereső rendszerben. A dokumentum és az igény összehasonlításához szükséges idő csökkentésének egyik módja, (különösen nagy dokumentum-file-ok esetében) hogy elkerüljük a dokumentumhoz rendelt deszkriptorokban a redundanciát. Ugyanakkor a tartalmi feldolgozás eredményének leírásában, illetve az igények megfogalmazásában egyre nagyobb szabadságot kell adnunk, ami egyben azt is jelenti, hogy közelíthetünk a természetes nyelvekhez.



## 1. Információs rendszerek

1.1. Az alábbiakban dokumentációs célokat szolgáló információtároló és visszakereső rendszerrel foglalkozunk. A rendszert "adatbanknak" is szokás nevezni, mivel nem szorítkozik a könyvtári gyakorlatban ismert dokumentációra: formai szempontból teljesen azonos rendszer megvalósíthat személynyilvántartást, szabadalmi, műszaki konstrukciók nyilvántartását, stb. Ezért a továbbiakban a dokumentum terminológia használatakor nem szorítkozunk ennek csupán könyv, folyóirat-cikk, publikáció jelentésére, hanem egészen általánosan egy (legtöbbször szövegesen megfogalmazott) információ-hordozó egységet értünk alatta.

A dokumentum tartalmi vagy formai ismérveit deszkriptoroknak nevezzük. Ezek általában szavak, szócsoportok, számok, olyan egységek, amelyek a dokumentumra vonatkozó valamilyen tulajdonságot fejeznek ki. Közismert deszkriptorok pl. a könyvtári gyakorlatban a folyóirat-cikk tartalmi megjelölésére az EFO-számok, egy-egy szakterületre jellemző (tezauruszban lefektetett) tárgyszavak, kifejezések; személynyilvántartásra gondolva a hajszin megjelölése, testmagasság, iskolai végzettség, stb. Egy dokumentumhoz több deszkriptor is tarthat, ezeket egymástól függetlenül rendeljük a dokumentumhoz.

Egy információs rendszerben különböző típusu deszkriptorok szerepelhetnek. Némelyek a tartalom, mások formai jegyek leírását szolgálják. A deszkriptor indikátora jelzi,

hogy az általa képviselt tulajdonság a dokumentumot mely oldalról közelíti meg. Az indikátor lehet valamilyen azonosító (pl. NÉV, ETOSZÁM,...), de gyakrabban számokkal jelezzük (legalábbis a számítógép belső ábrázolási formájában). Az indikátorok rendszere utal az információs rendszer felhasználási körére: a dokumentumokról tárolt ismérvek jellegét illetve a visszakeresés szempontjait határozza meg.

Az alábbi példákban az indikátort a deskriptor szövegétől / jel választja el a C indikátor pedig a tartalmi ismérvet jelöli:

SZERZŐ/Weissmann, G.

C/programozási nyelvek

C/LISP

C/szövegfeldolgozás számítógéppel

C/listák

A formai ismérvek (szerző, kiadó, leltári szám, stb.) alapján történő keresés nem okoz olyan jellegű problémákat, amellyel a dolgozat foglalkozik. A továbbiakban, amikor deskriptorokról beszélünk, gondoljunk mindig a publikáció tartalmi feltárása során nyert kulcsszavakra, tárgyszavakra. A könyvtári és dokumentációs gyakorlatban is ez az általánosan elfogadott jelentése.

Bár a tartalmi feltárás során a dokumentumhoz rendelt deskriptorokat egymástól függetleneknek tekintjük, az információkereső rendszerben használható, illetve használt deskriptorok rendszere strukturális tulajdonságokkal rendelkezhet: elemeit többféle kapcsolat fűzheti össze. A genetikus és specifikus fogalmakat jelölő kifejezések közötti

reláció például a deszkriptor-halmaz parciális rendezését adja; a bevezetőben említett "kitüntetett deszkriptorok" a struktúra ekvivalencia-osztályozását jelentik, stb. A deszkriptorrendszer struktúrája alkalmas arra, hogy 1./ a publikáció tartalmi feltárása során kapott deszkriptor-sorozatban és 2./ a felhasználói igényben szereplő deszkriptorokon a struktúrában megengedett műveleteket végezzünk.

A deszkriptor-halmaz elemeinek felsorolására (véges, de az időben változó, bővülő halmaz) és az elemek közötti kapcsolatok jelölésére alkalmas a tezaurusz. A tezauruszok közös vonása, hogy tartalmazzák a tudományok vagy a technika egy-egy területén előforduló fogalmakra a szakirodalomban általában használt kifejezéseket.

A tezaurusznak, mint halmaznak az elemeit nevezzük kulcsszavaknak. A legelterjedtebb strukturális tulajdonságok, amelyeket tezauruszokban szokásos megadni, a következők:

a.) Az ekvivalens kifejezések (a kifejezésekkel leírt fogalmak azonosak, illetve az információkeresés szempontjából azonos jelentésűek) rendszerének leírása. A tezaurusz minden egyes kulcsszóhoz az ekvivalens kifejezés-osztály egy kitüntetett elemét rendeli. Az osztály kitüntetett elemét megengedett deszkriptornak, az osztály többi elemét tiltott deszkriptornak nevezzük. Szokás deszkriptornak illetve nem-deszkriptornak is nevezni.

b.) A fogalom-hierarchiák jelölése. Ehhez elegendő a kitüntetett deszkriptoroknak megfelelő kulcsszavak mellett felsorolni az általa képviselt fogalom genetikus, illetve

specifikus fogalmait. Szokásos még rokon-fogalmak megjelölése is. Természetesen problémátikus, hogy a teaurusz milyen mértékben általánosítson, illetve specializáljon. Információkereső rendszerben pl. a túl általános fogalom a publikációk közül sokra érvényes, így a visszakeresési eljárás "bőbeszédű" lesz. A fogalom-hierarchia felépítése az adott területen igen nagy szakértelmet követel, ezenkívül figyelembe kell venni a teaurusz felhasználási területét is. Általában csak számítógép képes megállapítani, hogy a hierarchia valóban parciális rendezettség-e, azaz vannak-e a teauruszban ellentmondások.

c.) A teaurusz tartalmazhat a kulcsszavakra vonatkozó statisztikai paramétereket. Elsősorban egy publikációhalmaz elemeinek (automatikus vagy szakemberek által végzett) tartalmi kivonatolása során nyert összes deskriptoron belüli gyakoriság jegyezhető fel. Megállapítható a kulcsszó előfordulásának gyakorisága az információvisszakereső rendszer igényeiben, stb. Ezek a paraméterek segítséget nyújthatnak az automatikus tartalmi kivonatolás egyes fázisaiban, de szükségesek a teaurusz állandó korszerűsítéséhez, bővítéséhez, illetve az adott rendszerhez történő "hozzáidomításához".

A dokumentum egészének (teljes szövegének) tárolása az információs rendszerben általában terjedelme miatt nem megvalósítható, bár pl. mikrofilm-tárolóban nem lehetetlen. Továbbá a dokumentumok keresése a deskriptorok alapján történik, így a visszakeresés szempontjából csupán a dokumentumok azonosítását szolgáló ismérvek szükségesek: leltári

szám, szerző, cím, kiadási adatok, lelőhely. Sokszor ezekhez még egy természetes nyelven megfogalmazott tartalmi kivonat kapcsolódik, amely a felhasználó számára nyújt információkat, még mielőtt a publikációt kézbe venné.

A dokumentum azonosítását szolgáló ismérveket, együtt a tartalmi feltárás során kapott deskriptorsorozattal nevezük dokumentum-rekordnak, rövidebb terminológiával az információs rendszer dokumentumának. Az információs rendszerbe felvett dokumentum-rekordok összességét dokumentum-file-nak nevezük. A dokumentum-file a rendszer üzemelése, használata közben új dokumentum-rekordok belépésével állandóan bővül.

1.2. Az információs rendszereket, céljukat tekintve két csoportba szokták sorolni:

a.) Időszakonként tájékoztatást kíván nyújtani a dokumentum-file tartalmáról. A dokumentumokat egy vagy több szempont szerint sorbarendezve, lista formájában megjelenik a file. Tájékoztató kiadványok, bibliográfiák, review-k, indexek formájában publikálhatók. A deskriptoroknak általában a rendezési elvekben van szerepük. A rendszerben készülő kiadványok felhasználók kisebb-nagyobb csoportja számára készülnek, tömegkommunikációs eszközök.

b.) A több ezer vagy százezer dokumentumból bizonyos tulajdonságaik egy-egy csoportját bármikor szolgáltatni tudja. A tulajdonságok leírására a deskriptorokon kívül általában logikai műveleteket használhatunk. A logikai műveletek (és, vagy, negáció) szerepe az, hogy deskriptorok együttes előfordulását vagy hiányát írassuk elő. A tartalmi ismérvek

megadásának tipikus formáját a következő példa illusztrálja:

Dokumentáció és információvisszakeresés vagy  
automatikus információvisszakeresés vagy  
információvisszakereső számítógépes rendszerek és nem IBM

A tulajdonságok közlését a számítógéppel (felhasználói) igénynek nevezzük, azt az eljárást pedig, amely a dokumentum-file-ból a megfelelő csoportot meghatározza és a felhasználó számára hozzáférhetővé teszi, információvisszakeresésnek.

A rendszer által szolgáltatott információk egyének vagy szűk csoportnak szólnak.

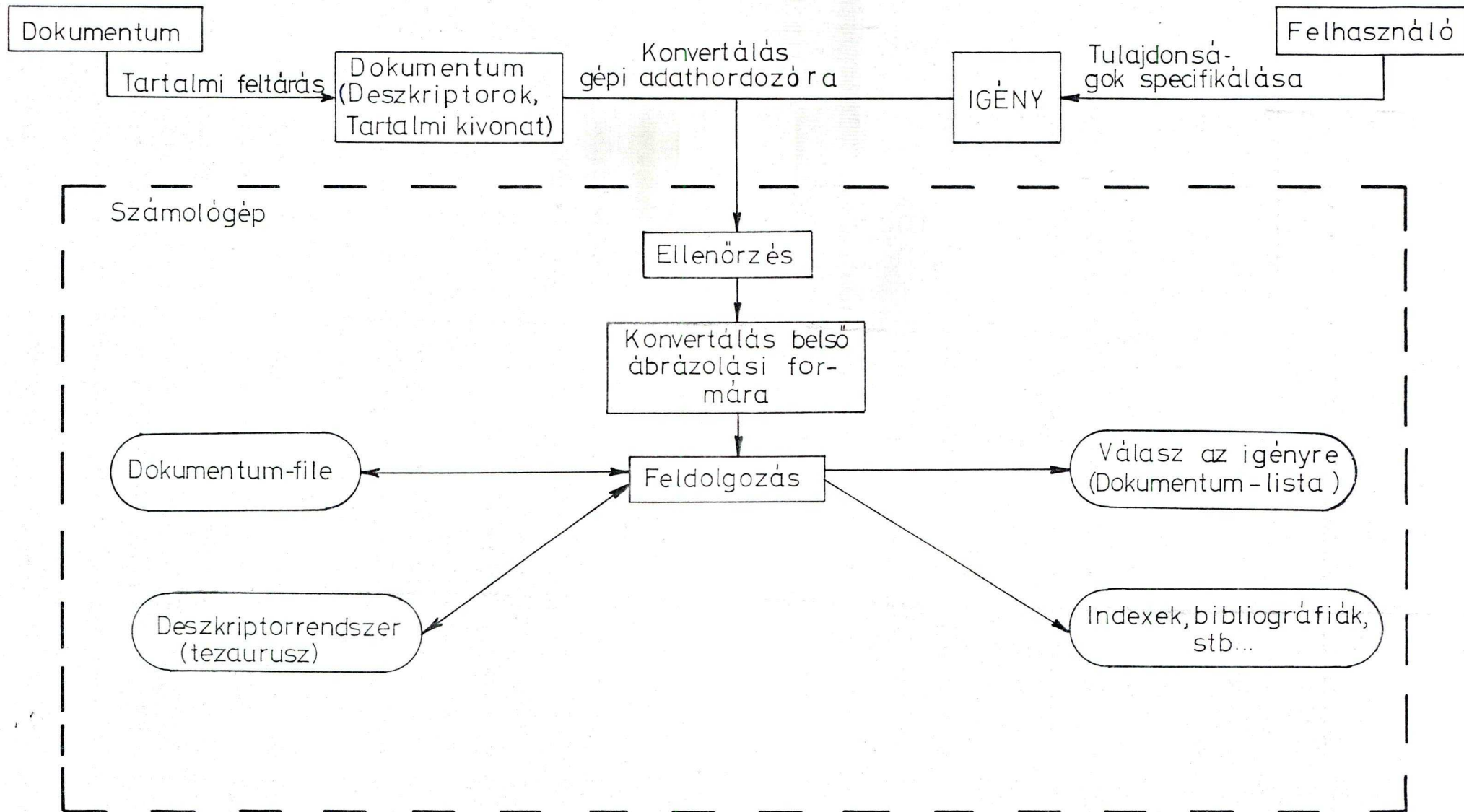
A két cél egymástól mereven nem határolható el, minden esetre a visszakereső rendszert a "kérdés-válasz" jellemzi. Az igények lehetnek hosszabb időszakra nézve állandóak vagy esetenkéntiek, ebből a szempontból a felhasználó (vagy a felhasználó és a gép között közvetítő személy) off-line vagy on-line (konverzációs) kapcsolatban lehet a rendszerrel. Másrészt a válaszadásra fordítható idő is igen változó (a könyvtári kölcsönzésben percek alatt, míg a témafigyelő szolgálatban havonként egyszer kell kielégíteni az igényeket). Mindezek a körülmények befolyásolják

- a dokumentum-rekordok és a dokumentum-file struktúráját
- a dokumentum-file gépi adathordozójának kiválasztását és alkalmazhatóságát (mágnesszalag, mágneslemez, mikro-film, stb.)
- a visszakereső eljárás felépítését
- használható számítógép hardware-strukturáját és konfigurációsükségletét.

1.3. Az 1. ábra az információs rendszerek általános strukturáját szemlélteti. Az egyedi rendszerek a különböző blokkok tényleges tevékenységeiben térnek el egymástól.

A dokumentum, rendszerbe lépésekor először tartalmának megfelelő szakember kezébe kerül, aki a tartalmi feldolgozást végzi. A számítógépi információs rendszerek használatában a leggyakoribb kifogás, hogy ez a fázis (legalábbis jelenleg) elkerülhetetlen. Az automatikus kivonatolással kapcsolatos kutatások (bár széles körben folynak) eddig nem eredményeztek kielégítő megoldást. Ez elsősorban nyelvészeti problémák megoldatlanságának következménye, amihez hazánkban is még egy súlyos kérdés járul: többnyelvűség a szakirodalomban.

A feldolgozást végző szakember a dokumentumhoz tartalmi kivonatot készít, deskriptorokat rendel hozzá, formai és egyéb ismérveit az indikátorrendszernek megfelelően összegyűjti. A számítógép ebben a fázisban is egyre inkább szerephez jut. A dokumentum azonosítását szolgáló ismérvek, a tartalmi kivonat és a tartalmi feltárás eredményeként kapott deskriptorsorozat gépi adathordozóra (lyukszalag, kártya, mágneskártya, stb.) kerül, amely már számítógép számára közvetlenül hozzáférhető. Függetlenül attól, hogy a deskriptorrendszer rendelkezik-e valamilyen strukturával, mindenképpen formai ellenőrzéseket kell első lépésben a számítógépnek végrehajtania. A deskriptor-szótárral (amely a visszakeresési eljárás végrehajtásához is általában szükséges, tehát rendelkezésre áll) összehasonlításokat végez. Ha ezenkívül a deskriptorrendszernek strukturája is van (azaz tezaurusz is ren-



1. ábra. Információs rendszer.



delkezésre áll), a struktúra alapján egy standard eljárást hajt végre a dokumentumhoz rendelt deskriptor-listán, amelynek eredményeképpen elmaradhatnak, megváltozhatnak, illetve bővíthetnek a lista elemei. Ezzel egyidőben jelzéseket ad a rendszer üzemeltetőjének, ami alapján korrigálható a belépő dokumentum-rekord illetve változtatható, bővíthető a tezaurusz struktúrája. A változtatások elsősorban a terminológia egységesítése miatt szükségesek. Nem elhanyagolható másik szempont, hogy az egyes fogalmakat jelölő kifejezések mellé felvegyük az adott fogalomra nézve genetikus fogalmakat jelölő kifejezéseket is. Így biztosítható, hogy a visszakeresési eljárás rátalál a dokumentumra akár a szűkebb, akár a bővebb fogalom alapján.

A specifikus fogalmakkal ugyanez az eljárás már nem ajánlatos: a tágabb kategóriából a szűkebbre segéd-információk nélkül nem következtethetünk.

A változatlanul hagyott vagy módosított dokumentum-rekordot, ezután a számítógép belső ábrázolási formára konvertálja, sorszámozva csatolja a dokumentum-file-hoz. A dokumentum-rekord struktúráját elsősorban a visszakeresési eljárás szempontjai határozzák meg: gyors hozzáférhetőség a deskriptorokhoz, szövegszerű összehasonlítások egyszerű végrehajthatósága. Mind a rekord, mind a rekord mezői változó hosszúságúak, ezért általában szét kell választani az indikátort és a tulajdonságot jelölő karaktersorozatot illetve számot. A indikátor mellett jegyezzük fel a deskriptor szövegére vonatkozó lényeges információkat (hossza, helye

a rekordon belül, stb.), így a visszakeresési eljárásban sorrakerülő összehasonlitások nagy része fix hosszúságu egységeken végezhető. Bibliográfiai rekordokra az USA-beli MARC-II. szabvány rögzít ezeknek a kritériumoknak megfelelő mágnesszalag-tárolási formát /3/.

A felhasználó a dokumentum-file anyagához az igény közlésével fér hozzá. Az igény elemi egysége a deszkriptor, szerepe pedig az, hogy a dokumentum-halmazt két részhalmazra osztja: az egyik csoport dokumentumaihoz hozzá van rendelve, a másik csoport dokumentumaiban nem szerepel az illető deszkriptor. Az igényben a két részhalmaz közül bármelyiket kijelölhetjük, ezt a részhalmazt azután további deszkriptorok segítségével újra részhalmazokra bonthatjuk, ilymódon az igény a dokumentum-file egy meghatározott részét választja ki. A használatos információs rendszerek a vázolt metszet-képzés művelete mellett általában az egyesítést is megengedik, de ennél többet csak akkor adhatunk, ha a deszkriptor-rendszert strukturával látjuk el.

1.4. A visszakeresési eljárás tényleges végrehajtására alapvetően két mód kínálkozik:

a./ A dokumentum-file minden egyes elemét összehasonlitjuk az igénnyel, a megfelelőeket különválasztva a felhasználó rendelkezésére bocsátjuk. A dokumentum-rekord és az igény összehasonlitása időigényes eljárás, nagy dokumentum-file-ok esetében így ez az út nem járható.

b./ A dokumentum-rekordokat deszkriptoraiknak megfelelő sorrendben tároljuk. Ha egy dokumentumhoz több deszkriptor tartozik, a dokumentum-rekord teljes egészében többször is-

méltódik. A deskriptorok segítségével megfogalmazott igény bejelentésekor adott deskriptor mellett megtaláljuk az igényelt dokumentumokat is. Az ismétlődéseken kívül a módszer további hátránya, hogy pl. két deskriptor együttes jelenlétének ellenőrzése igen nehézkes.

A két (szekvenciális és invertált) módszert egyesíti a kombinált információvisszakeresési eljárás /23, 25, 8/. Amikor a dokumentum a rendszerbe lép, egy sorszámot (indexet) kap. A dokumentum-file-ba a rekordok sorszám szerint növekvő rendben épülnek be. Ezzel egyidőjűleg módosul az a deskriptor-file, amelyben minden egyes deskriptor mellett fel vannak sorolva azoknak a dokumentumoknak az indexei, amelyekhez az illető deskriptort hozzárendeltük. Ha most egy igény megjelenik, a keresési eljárás azzal kezdődik, hogy

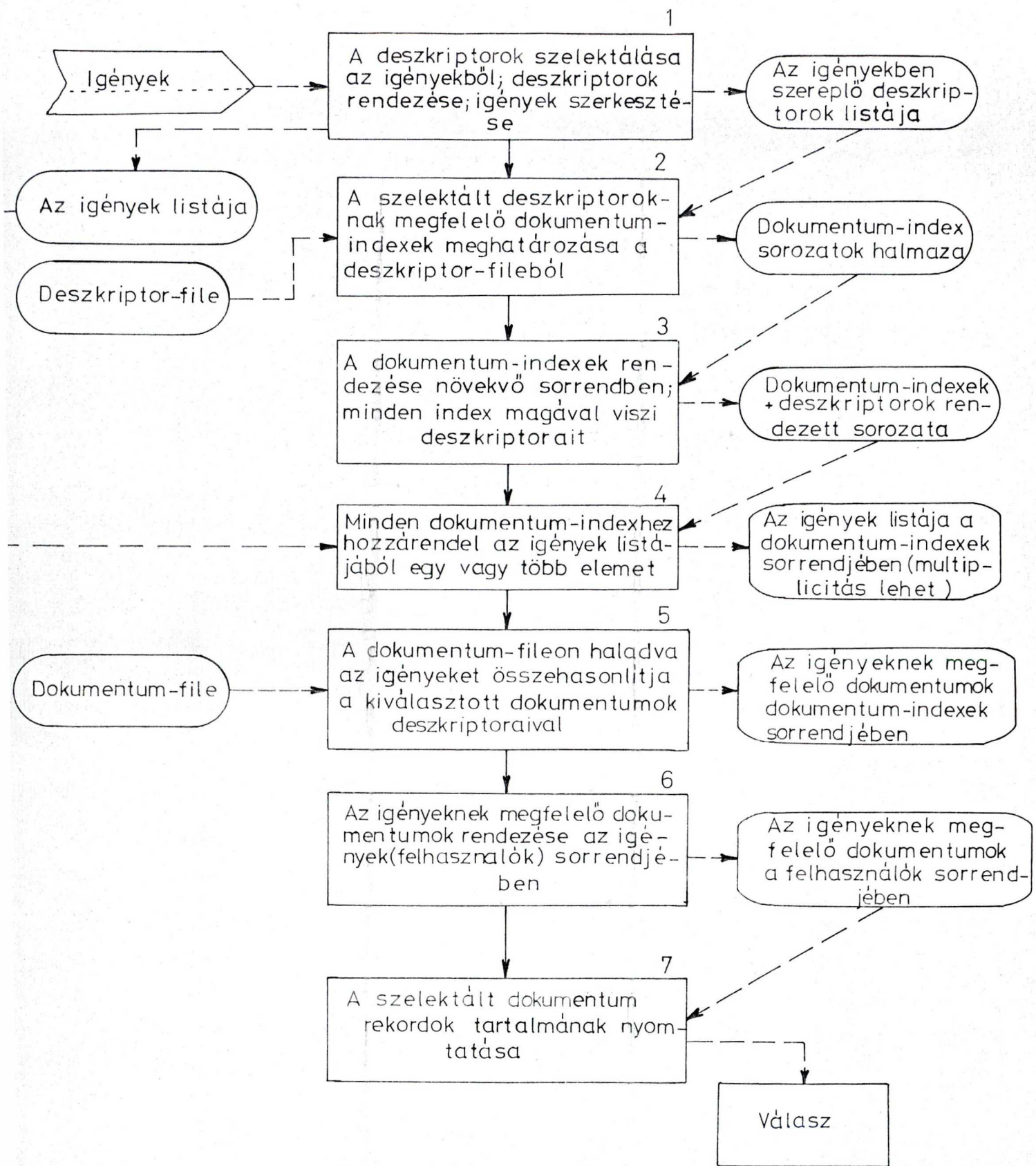
- szelektáljuk az igényből az összes előforduló deskriptort
- a deskriptor-file-ből kiolvassuk mindazokat az indexeket, amelyek a szelektált deskriptorok mellett találhatóak.

Ezután a tényleges összehasonlítást csak azokon a dokumentum-rekordokon hajtjuk végre, amelyek indexeit kigyűjtöttük.

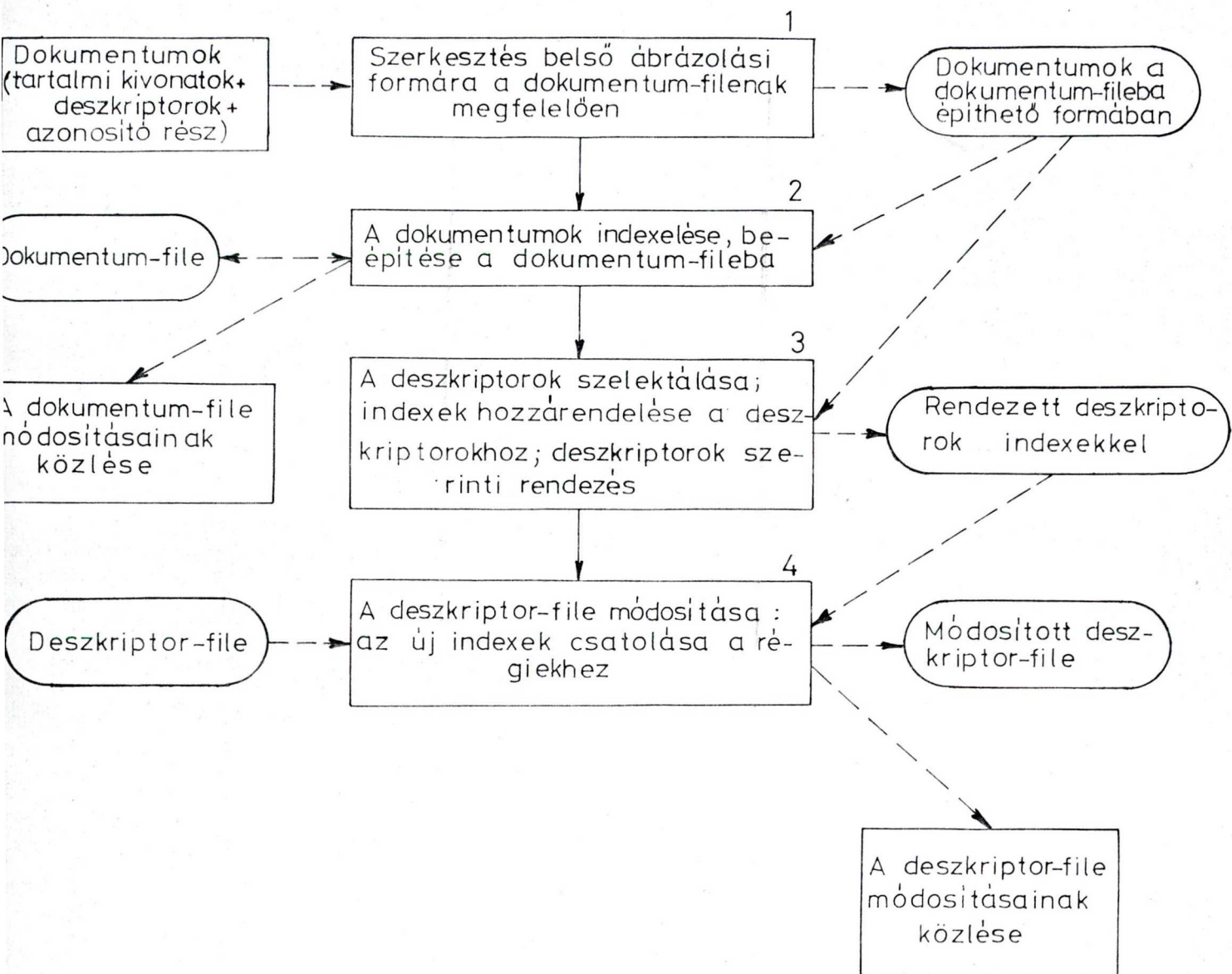
Az egy igényre fordított keresési időt lényegesen csökkenti az igények csoportosítása. A dokumentum-file egyszeri végigfutása közben több felhasználó igénye elégíthető ki. A 2. és 3. ábra kombinált információvisszakeresési eljárást alkalmazó rendszerben szemlélteti a dokumentum belépését illetve a visszakeresés végrehajtását. A téglalapok az egyes

tevékenységeket jelölik, közöttük a folyamatos vonallal jelölt nyilak a tevékenységek sorrendjét jelzik. A többi blokk a kiinduló, a közbűlső adatokat, illetve az eredményeket jelölik. A szaggatott vonallal jelölt nyilak az egyes tevékenységek input-adataiból illetve outputjaiba mutatnak.

A továbbiakban felépítjük a deskriptorrendszer strukturájának, a tezaurusznak lista-modelljét, definiáljuk azokat a műveleteket, amelyek alkalmasak a strukturának birtokában kiválasztani az igénynek megfelelő dokumentum-halmazt.



2. ábra. Kombinált információvisszakeresési eljárás



3. ábra. Az információs rendszer bővítése

## 2. Listák és elemi műveletek

2.1. A számítógépet hagyományosan numerikus jellegű feladatok megoldására használják. Az eközben felhasznált algoritmusok (bizonyos pontossággal megadott) racionális számokon végeznek aritmetikai műveleteket, számokkal képeznek relációkat, számokból építenek fel vektorokat, mátrixokat. Számok tárolására általában egy-egy memóriabeli rekesz (byte-szervezésű gépeken byte-ok egymásutáni sorozata) szolgál, míg a tömböket a memória egy intervalluma képviseli. Ez az ábrázolási forma két speciális tulajdonsággal rendelkezik:

a./ Egy szám ábrázolására a memóriában mindig ugyanannyi, illetve előre meghatározott nagyságu hely szükséges, így a számokból képzett tömbök memóriabeli helyigényének megállapításához elegendő az elemek számának ismerete.

b./ A tömb egy elemére egyértelműen hivatkozhatunk a tömböt tároló memória-intervallum kezdőcímével és az elem tömbön belüli sorszámával.

Mindkét előnyös tulajdonság annak köszönhető, hogy az alaplmenyiségeket (a számokat) ismert hosszúságu memóriaregység tárolja (szóban, byte-ban vagy bit-ben számolva). A nem-numerikus jellegű feladatok éppen ezt az alapfeltételt nem teljesítik: az elemi mennyiségek (karakter-sorozatok, pl. szavak; decimálisan ábrázolt számok) változó hosszúságúak. Másszóval mindaddig nem ismerjük az adat tárolásához szükséges memória méretét, amíg az adat ténylegesen rendel-

kezésre nem áll. Ebből következik új módszerek alkalmazásának szükségessége a memória optimális kihasználása végett.

A probléma megoldásához jelentősen hozzájárult a byte-szervezésű számítógépek megjelenése: a vezérlő illetve aritmetikai egység képes változó (az utasításban megadott vagy szó-határjellel definiált) hosszúságu alapegységeket értelmezni. Másrészt egy-egy feladattípus kapcsán különböző számítástechnikai eljárások (verem-memória, dinamikus cím-kezelés, stb.) alakultak ki.

Az alap-mennyiségek változó hossza az adatok strukturájának számítógépi ábrázolásában is újat jelentett. Ha pl. egy vektor komponenseit továbbra is indexüknek megfelelő sorrendben egymás utáni memória-címen helyezzük el, akkor egy adott elem címe az összes megelőző elem hosszának függvénye. Ezenkívül egy vektorelem helyettesítése a korábbinál hosszabb karaktersorozattal együtt jár azzal, hogy az összes nagyobb indexű komponenseket a memóriában feljebb kell tolni.

Egyébként sem csupán arra van szükség, hogy a meglévő numerikus eljárások adatstrukturáját másoljuk új tartalommal. Ahogy szélesedik a számítógépek felhasználási köre, egyre újabb strukturák gépi ábrázolására van szükség, ami együtt jár az algoritmusok elemi műveleteinek változásával.

2.2. A listák /2, 7, 19, 27/ a gráfelméletből ismert irányított bináris (dichotomikus)gráfok strukturájával írják le az adatokat. A gráf pontjai a cellák illetve atomok. Az atomok képviselik az elemi adattípusokat, amelyek általá-



ban karaktersorozatok. Az atomokból élek nem indulnak ki, míg minden cella a gráf két pontjára mutató élt reprezentál. A listaműveletek egy része új listákat állít elő vagy meglévőket módosít, míg más része egy cellából egy másik cellába vagy atomba vezető utat definiál, ezáltal allistát (speciális esetben atomot) jelölve ki.

Szükségesek, de nem takinthatók listaműveleteknek azok a műveletek, amelyeket atomokon definiálunk és az atomok belső strukturáján alapulnak (aritmetikai, logikai műveletek, az atom hosszának meghatározása karakterben számolva, két elem összehasonlítása, stb.).

A továbbiakban szereplő példákat az Országos Műszaki Könyvtár és Dokumentációs Központ által 1970-ben közreadott Általános műszaki fogalmak teaurusza /1/ szolgáltatja. Az alapanyag előállításában (a rendelkezésre álló angol nyelvű teaurusz átrendezésében a magyarra fordítás után, a három rész rendezésében, a permutációs fejezet kidolgozásában) a Kibernetikai Laboratórium Minszk-22 számítógépe is szerephez jutott. Az Általános műszaki fogalmak teaurusza megkülönböztet deszkriptoroknak és nem-deszkriptoroknak nevezett kulcsszavakat. A nem-deszkriptorok mellett megtalálhatjuk az ugyanazt a fogalmat kifejező deszkriptort. A deszkriptoroknál megjelöli a rokon fogalmakat, a specifikus fogalmakat és a genetikusan fogalmakat reprezentáló deszkriptorokat. Az alábbi lista-modell teaurusz-reprezentációjában az atomoknak a kulcsszavak (deszkriptorok és nem-deszkriptorok) felelnek meg, a listák strukturái pedig a deszkriptorok fogalom-kapcsolatait fejezik ki. Pontosabban szól-

va: minden kulcsszónak megfelel (egyértelműen, de nem kölcsönösen egyértelműen) egy lista, amelynek speciális al-listáiban szereplő atomok képviselik az azonos jelentésű deszkriptort, a speciális, genetikus és rokon fogalmak halmazait.

A műveletek jelölési módjában a LISP /29/ programozási nyelv koncepciója az irányadó.

2.3. A listák elemi egysége az atom. Megjelenési formája a karaktersorozat, hossza nem korlátozott. A lista-műveletekben mint oszthatatlan egység vesz részt. Pl.

Fourier-sorok

integrálok

-512

A4

Speciális (nem lista-) műveletek, amelyeket atomokon definiálunk, tulajdoníthatnak az atomoknak "belső" strukturát. Ez alapján osztályozhatjuk az atomokat: számok, szavak, kifejezések, logikai értékek, kódok, stb. Deszkriptorokra gondolva lehet az indikátor és a tulajdonságot kifejező karaktersorozat egysége /T: nem-deszkriptor vagy tiltott kifejezés, D: deszkriptor/:

T/ harmonikus analízis

D/ Fourier-analízis

Hogy egyetlen atomból is lehessen listát felépíteni, illetve, hogy a cellákból mindig két él fusson ki, bevezetjük a speciális üres atom fogalmát, amelyet 0 jelöl.

2.4. Cellának nevezzük az (x,y) párt, ahol x is és

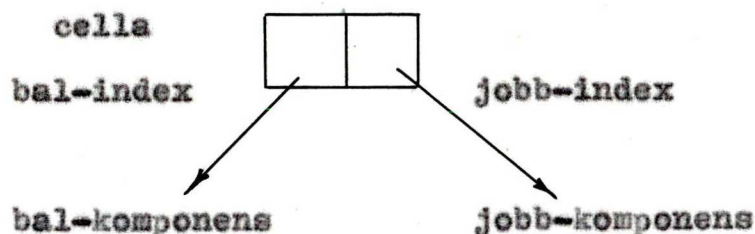
y is atom vagy cella lehet. x a bal-, y a jobb-komponens. Pl:

(T/ harmonikus analízis, D/ Fourier-analízis)

(D/ Fourier-analízis, 0)

(T/ harmonikus analízis, (D/ Fourier-analízis, 0)).

A cellákat írásban a fenti módon zárójelpárral és vesszővel jelöljük. A grafikus ábrázolás részben szemléletes, másrészt a cella számítógépi reprezentációjához is közel áll: a cella bal- és jobb-indexből áll, amelyek rendre a bal- illetve jobb-komponensre mutatnak. Számítógépben az indexek olyan címek, amelyek a komponensek memóriabeli helyére utalnak.



## 2.5. Listának nevezünk

1./ Egy atomot, vagy

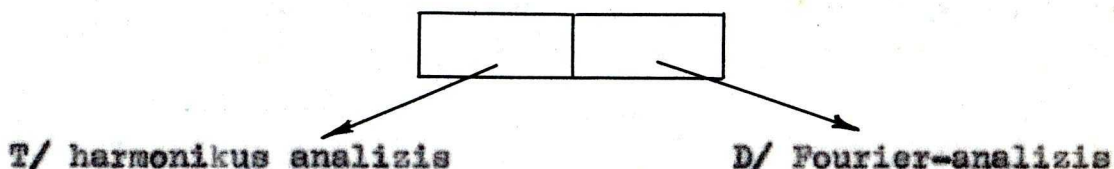
2./ listákból képzett párt (cellát).

Írásban a lista minden egyes cellájára egy zárójelpár utal. A listát egyértelműen megjelölhetjük (azonosíthatjuk) a legkülső zárójelpárnak megfelelő cellával, a lista-fejével. A lista bármelyik cellájához (zárójelpárhoz) illetve atomjához tartozik egy al-lista. Az al-lista fejét (azaz magát az al-listát) a listában definiált uttal adhatjuk meg: az a cellából (a lista feje) a b cellába vagy atomba vezető uton értjük az  $a_1, a_2, \dots, a_n (n > 0)$  sorozatot, ahol  $a_1 = a, a_n = b$

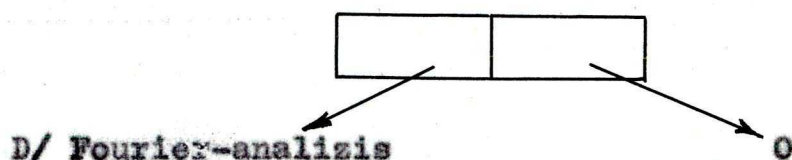
és minden  $a_i$  ( $i = 1, 2, \dots, n-1$ ) olyan cella, amelynek  $a_{i+1}$  vagy bal-, vagy jobb-komponense.

Példák:

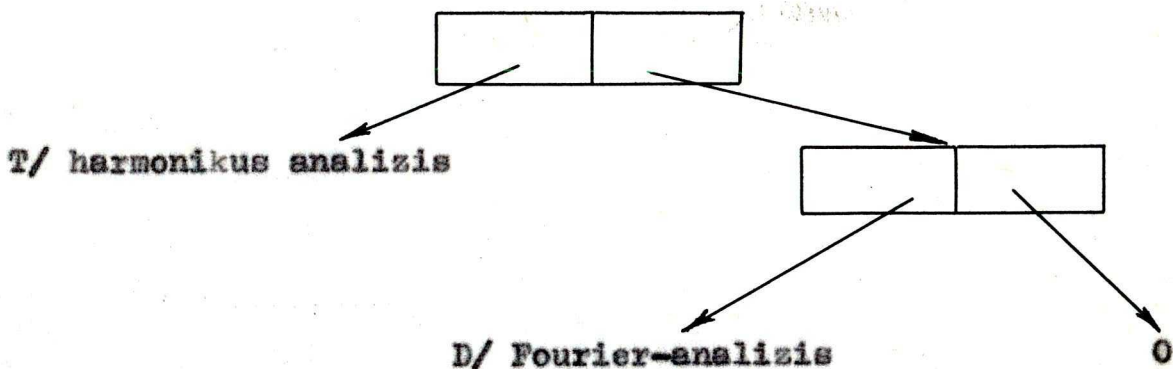
1./ (T/ harmonikus analízis, D/ Fourier-analízis)



2./ D/ Fourier-analízis, 0)



3./ (T/ harmonikus analízis, (D/ Fourier-analízis, 0))



2.6.  $L_1$  és  $L_2$  lista. Az  $(L_1, L_2)$  pár (cella) egy új lista feje, azaz a pár-képzés kétváltozós, listákon értelmezett és lista-értékű művelet. Hagyományos műveleti jelek hiányában függvény-jelölést használunk:

$$(L_1, L_2) = \text{cons}(L_1, L_2)$$

Számítógépben ez egyetlen új cella képzését jelenti, amelynek bal-indexe az  $L_1$ , jobb-indexe az  $L_2$  lista fejére mutat.

A függvény-jelölésben a zárójelpárnak és vesszőnek

argumentum-kijelölő, míg a cella-írásmódban bal- és jobb-komponenset meghatározó szerepe van. Ezt a kettősséget az egyes programozási nyelvek úgy oldják meg, hogy valamelyik jelölésmódot kizárják. A fordító- illetve értelmező-programok szempontjából a függvény-jelölés kedvezőbb: balról haladva az argumentum(ok) olvasásakor a végrehajtandó művelet már rendelkezésre áll.

2.7. Al-lista illetve ut definiálására szolgál a hd és tl művelet. Mindkettő nem egyetlen atomból álló L listán van értelmezve. hd(L) az L lista fejének bal-komponense (cella vagy atom). tl(L) az L lista fejének jobb-komponense (cella vagy atom).

A cons művelettel ellentétben a hd(L) és tl(L) nem képez új cellát, hanem kijelöl egy meglévőt.

2.8. A cons, hd, tl és egyéb elemi műveletekből kompozíció-képzéssel alkothatunk összetett kifejezéseket. Néhány egyszerű összefüggést találunk a definiált műveletek között:

$$\begin{aligned} \text{hd}(\text{cons}(L_1, L_2)) &= L_1 \\ \text{tl}(\text{cons}(L_1, L_2)) &= L_2 \\ \text{hd}(\text{hd}(\text{cons}(L_1, L_2))) &= \text{hd}(L_1) \\ &----- \end{aligned}$$

Viszont

$$f(L) = \text{cons}(\text{hd}(L), \text{tl}(L)) \neq L,$$

mivel a cons művelet mindig egy új cellát definiál, Az  $f(L)$  lista strukturája azonos L strukturájával, de lista-fejük különböző cella.

## 2.9. Legyen

$$H = (L_1, L_2, \dots, L_n)$$

listák egy rendszere. (Rendszert kell mondanunk halmoz helyett, mert elemeit szoros kapcsolatok fűzhetik össze. Különböző listáknak lehetnek közös al-listái; pl. az  $L_i$  és  $L_j = \text{cons}(\text{hd}(L_i), \text{tl}(L_i))$  listák csak a lista-fejben különböznek. Másrészt minden cella illetve atom képviselhet egy listát: közülük csak bizonyosakat soroltunk  $H$ -ba, a többit al-listaként kezeljük.)

A  $\text{cons}$ ,  $\text{hd}$ ,  $\text{tl}$  műveleteket (ismételten) alkalmazva a  $H$  rendszer elemeire, olyan  $H'$  rendszert kapunk, amelynek  $H$  része:  $H \subseteq H'$ . Nem változtattunk meg egyetlen  $L_i$  listát (sőt cellát) sem.

A  $\text{set}$ ,  $\text{sethd}$ ,  $\text{settl}$  műveletek cellákat módosítanak. A  $\text{set}(L_1, L_2)$  függvény alkalmazásának eredményeképpen az  $L_1$  cella helyébe  $L_2$  kerül. A  $\text{sethd}(L_1, L_2)$  illetve a  $\text{settl}(L_1, L_2)$  az  $L_1$  cella bal- illetve jobb-komponense helyébe az  $L_2$ -t helyettesíti.

Ily módon a  $\text{set}$  műveletek azt, vagy azokat a listákat, amelyekből való cellát módosítanak, megváltoztatják.

A listák számítógépi tárolásában a cella-módosító műveletek egy igen lényeges problémát vetnek fel. Az atomoktól eltekintve a listák elemei (a cellák) fix hosszúságú memóriaegységekben (rekeszekben, byte-okban) tárolhatók. Maguk a listák viszont változó számú cellából épülnek fel, így egy  $H$  listarendszer számára szükséges memóriaintervallum a rajta végzett műveletektől (a belőlük képzett listák hely-

igényétől) függ. Listákat ábrázoló programozási nyelvek, illetve programozási rendszerek ezért általában a lehetőségekhez (memóriakapacitás, felhasználó igény-jelzései, stb.) mérten maximális munkamezőt jelölnek ki a program számára, amelyben el kell férni mind a kiinduló adatokat képviselő, mind a generált listáknak. A listákat programozási nyelvekben a változóknak megfelelő azonosítók jelölik, amelyeket a program írója deklarál. Minden azonosító tulajdonképpen egy lista-fej címe a munkamezőben. Az al-listáknak megfelelő celláknak nincs azonosítójuk, ugyanígy nincs a kifejezések feldolgozása során keletkező közbülső listákat reprezentáló celláknak sem.

A munkamezőbe tehát egyre újabb és újabb cellák kerülnek. Ezzel egyidejűleg cellák válhatnak feleslegessé: a műveletek pl. "lekapcsolnak" al-listákat egy-egy listáról. A munkamező optimális kihasználása érdekében így időnként a "felesleges" cellákat fel kell szabadítani, hogy helyüket újra felhasználhassuk.

Tegyük fel, hogy az  $M$  munkamezőben az  $(a_1, a_2, \dots, a_n$  változókhoz rendelt)  $L_1, L_2, \dots, L_n$  listákat tároljuk. A  $c$  cellát az  $M$  mezőben foglaltnak nevezzük, ha az  $L_1, L_2, \dots, L_n$  listák valamelyikének lista-fejéből ut vezet  $c$ -hez.

Az  $M$  mező többi cellája "felesleges", mivel sem azonosítóval, sem al-lista képzéssel nem tudunk rá hivatkozni, így időnkénti (egy-egy kifejezés feldolgozása utáni) regeneráló eljárással felszabadítható további célokra.

A  $c$  cella foglaltsága függvénye az  $M$  mezőhöz rendelt változóknak. Minden  $a_i$  változóhoz (azonosítóhoz) tar-

tozik egy  $L_1$  lista-fej. Asszociációs listának nevezzük az  $M$  mezőben az  $(a_1, L_1), (a_2, L_2), \dots, (a_n, L_n)$  sorozatot.

Az asszociációs listákat általában valóban lista-formában tárolják a programozási rendszerek, mivel a program futása során éppugy változhatnak, mint a kifejezések segítségével felépített listák. (Deklarációk, blokkhatárok, eljáráshívások egyrészt az asszociációs lista hosszát változtatják, másrészt az azonosítókhoz új listákat rendelnek.)

2.10. Feltételes kifejezések képzéséhez logikai értékű függvények definiálására van szükség.

$L, L_1, L_2$  listák. Az atom  $(L)$  függvény értéke igaz, ha  $L$  egyetlen atom, különben hamis. Az egyenlőség értelmezése már problémákat okoz: definiálhatjuk a strukturák azonosságaként vagy egyszerűen a lista-fejek indexeinek megegyezése alapján. Ez utóbbi általában trivialis, így a következő definíció hasznosabb:

a./  $L_1, L_2$  mindegyike atom.  $eq(L_1, L_2)$  igaz, ha az  $L_1, L_2$  atomok megegyeznek (mint karaktersorozatok).

b./ Ha  $L_1$  és  $L_2$  közül az egyik atom, a másik cella, akkor  $eq(L_1, L_2)$  hamis.

c./  $L_1, L_2$  cellák,  $B_1, B_2$  illetve  $J_1, J_2$  a bal- illetve jobb-komponenseik.  $eq(L_1, L_2) = eq(B_1, B_2) \wedge eq(J_1, J_2)$ .



### 3. A tezaurusok listái

3.1. Információkereső rendszerünkben listákat használunk mind a tezaurusz strukturájának leírásában, mind pedig a felhasználói igények megfogalmazásában. Eközben a tezaurusz kulcsszavai (amelyek megegyeznek a dokumentumokhoz rendelhető deskriptorokkal) kétféle szerepük:

a./ A tezaurusz listáinak azonosítói, azaz nevek. Minden egyes kulcsszóhoz tartozik a tezauruszban egy lista (több kulcsszóhoz esetleg ugyanaz a lista): a tezaurusz asszociációs listája felelteti meg a kulcsszavaknak a listafejeket. A felhasználói igény alapegysége a lista-kifejezés, amelynek atomjai a kulcsszavak, szerepük pedig azonosító (változónév) jellegű. Az igény egy kulcsszavának értéke az a lista, amelyet a tezaurusz asszociációs listája a kulcsszóhoz rendel. A felhasználói igényekben megengedhetünk olyan kulcsszavakat is, amelyeket nem tartalmaz az asszociációs lista. A kulcsszónak, mint változónak az értéke legyen ilyenkor az egyetlen atomból (az adott kulcsszóból) álló lista pl. a tulajdonnevek. (IBM, VIDEOTON) nem vehetők fel a tezauruszba, illetve a nevekben nem érhetjük el még a viszonylagosan teljes felsorolást sem.

b./ A tezaurusz listáiban ugyanezek a kulcsszavak (de már csak a deskriptorok) a listák atomjai lesznek, azaz elemi adatok, és nem azonosítók.

Az információkereső rendszer tezaurusza tehát két részből áll: az asszociációs listából és a deskriptorokhoz tar-

tozó listák halmazából. Részletes strukturáját a következőkben adjuk meg.

3.2.  $L$  lista,  $k \geq 0$  egész szám,  $t$  a lista feje. Az  $L$  lista  $k$ -adrendű (jobb) al-listájának nevezzük azt az  $L^k = L \uparrow k$  listát, amelynek lista-feje  $t^k$  (ha  $L^k$  egyetlen atom, akkor  $t^k$  ezt jelöli) és van olyan  $a_0, a_1, \dots, a_k, a_{k+1}$  cellasorozat ( $a_{k+1}$  atom is lehet), hogy  $a_0 = t, a_{k+1} = t^k$ , minden  $i$ -re ( $i=0, 1, \dots, k-1$ )  $a_{i+1}$   $a_i$ -nek jobb-komponense és  $a_{k+1}$   $a_k$ -nak bal-komponense.

Tehát  $L \uparrow 0$  a lista-fej bal-komponense,  $L \uparrow 1$  a lista-fej jobb komponensének bal-komponense, és így tovább.  $L \uparrow k$  értelmezését kiterjeszthetjük olyan  $k$  számokra is, amelyekhez nem található a definícióban szereplő  $a_i$  sorozat (a jobb-indexek mentén haladva atomhoz jutunk  $p \leq k$  lépésben):  $L \uparrow k$  értéke legyen ilyenkor a  $0$  atom.

A  $d$  deskriptorhoz a teaurusz asszociációs listája rendelje az  $L$  listát. Az  $L \uparrow k$  lista fejéből elérhető atomok halmazát  $d$   $k$ -adrendű értékkészletének nevezzük és  $d_k$ -val jelöljük.

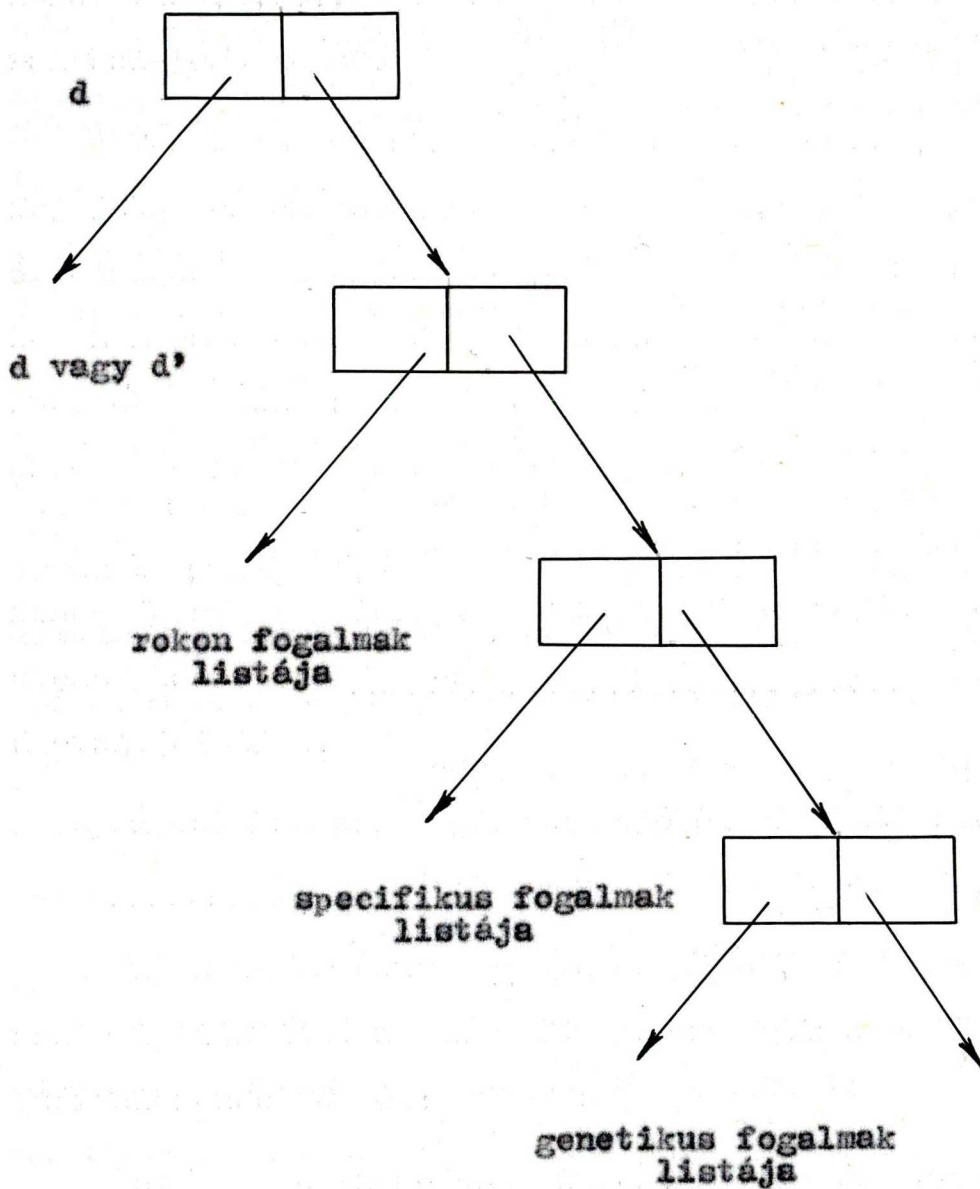
Általában az  $L$  lista fejéből elérhető atomok halmazát  $L$  értékkészletének nevezzük és  $L_H$ -val jelöljük.

A teauruszban a  $d$  kulcsszóhoz, mint azonosítóhoz az asszociációs listán keresztül rendelt lista a következő struktúrával rendelkezék:

a./ Ha  $d$  nem-deszkriptor, akkor  $d_0$  egyetlen deskriptorból álljon, mégpedig a  $d$ -nek fogalmilag megfelelő  $d^0$  kitüntetett deskriptorból. Ha  $d$  deskriptor, akkor

$d_0 = d.$

b./  $d_1, d_2, d_3$  rendre a rokon fogalmak, specifikus fogalmak, genetikus fogalmak deszkriptorainak halmazai legyenek. További fogalom-kategóriák felvételével a struktúrában  $d_4, d_5, \dots$  is kaphat szerepet.





3.4. A tezaurusz számítógépi hordozója közvetlen elérési háttérmemória lehet: mágneslemez, mágnesdob. A szekvenciális elérhetőség (pl. mágnesszalagos tárolás esetén) nem megfelelő még akkor sem, ha maguk a listák egy-egy rekordot képeznek: részben a deskriptorok ismétlődései miatt (több tizszeresére növelné a tezaurusz terjedelmét), másrészt az asszociációs lista kifejezetten a direkt elérési adatszervezés "katalógusának" szerepét tölti be.

A kombinált információvisszakereső rendszerre gondolva a tezaurusz strukturája a következő:

a./ Az asszociációs lista egy eleme

$(d, L, I_1, I_2, \dots, I_n)$

alakú.  $d$  a kulcsszó szövege,  $L$  a kulcsszóhoz rendelt listából képzett rekord közvetlen elérési címe.  $I_1, I_2, \dots, I_n$  ( $n \geq 0$ ) azok a dokumentum-sorszámok, amelyekhez a  $d$  kulcsszót hozzárendeltük tartalmi feltárásuk, vagy a dokumentum-file-ba lépésükkor. Ilymódon az asszociációs lista betölti a kombinált információvisszakereső rendszer deskriptor-file szerepét is.

b./ Minden egyes lista egy rekordot képez a tezaurusz adathordozóján. A cellák jobb-indexei mindig rekordon belüli valamilyen másik cellára vagy a 0 atomra utalnak. A bal-indexek vagy rekordon belüli cellára vagy az asszociációs lista egy elemére (azaz deskriptorra) mutatnak. A cellák tehát nem szimmetrikusak, a jobb-indexek tárolására kevesebb bit is elegendő. Másrészt a cellák ilyen felépítése biztosítja, hogy egy index csak akkor vezet ki a rekordból, ha atomra mutat, így a tényleges lista-műveletekhez szükséges információk a

rekordon belül megtalálhatók. (Még az atomok összehasonlítására is elegendő: minden atom (deszkriptor) csak egyszer szerepel az asszociációs listában, ezért a rájuk mutató indexek egyenlősége az atomok egyenlőségével azonos.)

3.5. Az információvisszakereső rendszerben a (felhasználói) igény megfogalmazására (ld. 1.2.) tezaurusz birtokában a következő nyelv alkalmas:

3.5.1. A változók azonosítóinak szerepét a kulcsszavak töltik be. A változó értéke az asszociációs listán keresztül a kulcsszóhoz rendelt lista.

3.5.2. Feltétlen és feltételes lista-kifejezéseket a változókból, lista-műveletekből (cons, hd, tl, set, sethd, settl, ↑, ...), logikai értékű műveletekből (atom, eq, ...) építhetünk fel. Minden lista-kifejezés eredményeképpen egyetlen lista alakul ki, ez a kifejezés értéke.

3.5.3. Az igény lehet (rekurzív definíció):

- a./ egy lista-kifejezés,
- b./ igény negációja,
- c./ igény kerek zárójelek között,
- d./ két igény közé a konjunkció vagy diszjunkció jelét téve.

3.6. Szemantikailag az  $L$  lista az igényben az  $L$  érték-készletének  $L_H$  halmazát jelöli. Ha az igény egyetlen  $L$  lista-kifejezésből áll, akkor az igényt mindazok a dokumentumok kielégítik, amelyekhez hozzárendeltünk az  $L_H$  halmazból való deszkriptort.

Az  $L^1$  és  $L^2$  lista-kifejezés között a diszjunkció jele

az  $L_H^1$  és  $L_H^2$  halmazok egyesítését jeleníti.

Az  $L^1 \wedge L^2$  igényt mindazok a dokumentumok kielégítik, amelyekhez  $L_H^1$ -ből is és  $L_H^2$ -ből is rendeltünk deszkriptort.

A  $\neg L$  igényt azok a dokumentumok elégítik ki, amelyekhez  $L_H$ -ből való deszkriptor nincs rendelve.

Természetes módon definiáljuk a  $\neg L^1 \wedge L^2$ ,  $\neg L^1 \wedge \neg L^2$ ,  $\neg L^1 \vee L^2$ ,  $\neg L^1 \vee \neg L^2$  igények kielégíthetőségét is. Ezeknél összetettebb igény-kifejezésekre a kielégíthetőséget a konjunktív vagy diszjunktív normál-formára hozott alakon definiáljuk (a változók szerepét a lista-kifejezések illetve az általuk felépített listák értékkészleteinek halmazai töltik be).

3.7. Az információvisszakeresési eljárás a következő lépésekben hajtható végre:

a./ Az igény-kifejezést konjunktív vagy diszjunktív normál-formára hozzuk (a benne szereplő lista-kifejezéseket azonosítókkal helyettesítve).

b./ Az igényben szereplő kulcsszavakhoz az asszociációs lista segítségével kikeressük a teaurusz megfelelő listáit.

c./ Az igényben szereplő lista-kifejezések értékét képviselő deszkriptor-halmazokat előállítjuk.

d./ A deszkriptor-halmazokban szereplő deszkriptorok (asszociációs listában tárolt) dokumentumsorszámait összegyűjtjük és növekvő nagyságrendben egy I sorozatot képezzünk belőlük.

e./ Az I sorozat minden egyes elemének megfelelő dokumentum-rekord deszkriptorsorozatát összehasonlítjuk az igényt képviselő konjunktív vagy diszjunktív normálformával. A normál-

forma azonosítói helyett a nekik megfelelő deskriptor-halmazokkal dolgozunk a 3.6. szabályai szerint.

Példa. Legyen a felhasználói igény-kifejezés a következő alakú:

cons(differenciálegyenletek $\uparrow$ 0, differenciálegyenletek $\uparrow$ 2) $\wedge$   
cons(differenciálanalizátorok $\uparrow$ 0, differenciálanalizátorok $\uparrow$ 1)

Az igényben szereplő kulcsszavak mellett az Általános műszaki fogalmak tezauruszában a következőket találjuk (a lista-formában való megadás áttekinthetlenebb, ezért felsorolás-formában adjuk meg):

differenciálegyenletek

specifikus fogalmak

genetikus fogalmak

rokon fogalmak

lineáris differenciálegyenletek  
parciális differenciálegyenletek  
egyenletek  
Bessel függvények  
differenciálhányadosok  
differenciálok (matematika)  
exponenciális függvények  
Fourier-analízis  
integrálegyenletek  
kezdeti feltételek  
Laplace-transzformáció  
számítási eljárás  
vektor analízis

differenciálanalizátorok $\rightarrow$ analóg számítógépek

analóg számítógépek

genetikus fogalmak

matematikai gépek  
számítógépek



rokon fogalmak    analízátorok (elektromos)  
analóg-digitális átalakítók  
analóg-digitális számítógépek  
differenciálegyenletek  
digitál-analóg konverterek  
digitális számítógép  
elektronikus számítógép  
folytonosság  
integrálók  
logikai áramkörök

A két lista-kifejezést, az igényben való szereplésük sorrendjében jelöljük a-val és b-vel. A teaurusz alapján a-nak és b-nek megfelelő listák fejéből elérhető atomok (deszkriptorok) halmazát jelölje A és B, amelyek a következők:

$A = \{ \text{differenciálegyenletek, lineáris differenciálegyenletek, parciális differenciálegyenletek} \}$

$B = \{ \text{analóg számítógépek, Bessel-függvények, differenciálhányadosok, differenciálok (matematika), exponenciális függvények, Fourier-analízis, integrálegyenletek, kezdeti feltételek, Laplace-transzformáció, számítási eljárás, vektor analízis} \}$

Tegyük fel, hogy a számítógépben tárolt teaurusz asszociációs listájában A és B elemei mellett az  $i_1, i_2, \dots, i_n$  indexeket találjuk meg. Ekkor a dokumentum-file  $i_1, i_2, \dots, i_n$  indexű rekordjaiban található deszkriptorsorozatot hasonlítjuk össze A és B elemeivel. Közülük egy rekord akkor és csak akkor felel meg az igénynek (amely végül is  $A \wedge B$  alakú),

ha tartalmaz A-ból is és B-ből is deszkriptort.

3.8. A dokumentum-rekord információs rendszerbe lépésekor kulcsszavait az asszociációs lista alapján ellenőrizzük. Ha a  $d$  kulcsszó nem-deszkriptor, a tezaurusban hozzárendelt lista  $d_0$  0-adrendű értékészletének megfelelő  $d'$  deszkriptorral helyettesítjük. Ehhez le kell olvasni a lista-rekordok közül is némelyeket. Ha azt akarjuk, hogy elegendő legyen csupán az asszociációs lista alapján beépíteni a dokumentum-rekordot, a nem-deszkriptorokhoz rendelt asszociációs lista-elemet

$(d, L, d')$

formában kell tárolni: a dokumentum-sorszám sorozat helyébe (nem-deszkriptor esetében ez ugyis üres) az ekvivalens jelentésű deszkriptor kerül.

Nem elegendő azonban a dokumentum-rekord beépítéséhez a tezaurusz asszociációs listája, ha nemcsak a nem-deszkriptor kulcsszavakat helyettesítjük az ugyanazt a fogalmat jelölő deszkriptorral, hanem automatikusan a genetikus (esetleg rokon) fogalmak deszkriptorait is hozzá kívánjuk csatolni. Formálisan ez azt jelenti, hogy a  $d$  kulcsszó helyébe a  $d_0$  és  $d_3$  (esetleg  $d_1$ ) deszkriptor-halmazok elemeit tesszük. Az ilyen módon keletkező deszkriptorsorozat azonos tagjai közül természetesen csak egy kerül a dokumentum rekordba.

#### 4. További lehetőségek

4.1. A működő információs rendszerek egy részében a dokumentum-rekordokhoz tartalmi kivonat és nem deskriptorsorozat tartozik. A visszakeresési eljárás sebességének fokozása érdekében a dokumentum-rekord belépésekor a kivonatot szavakra tagolják és a továbbiakban, mint deskriptorokkal dolgozik velük a rendszer. A kivonat tagolására két módszer kínálkozik:

a./ A rendszerbe építünk egy irreleváns szavakból álló szótárt (névelők, kötőszók, stb.). A kivonat minden olyan szava deskriptor lesz, amely a szótárban nem szerepel. Hátránya, hogy a ragozott alakok külön-külön deskriptort képviselnek, amit a felhasználói igények megfogalmazása küszöböl ki. (Csak a szótót jelöli, amihez bármilyen ragot elfogad. Pl. asztal~~sz~~ jelöli az asztal, asztalt, asztalnak, ... szavakat.) Másrészt a tartalomra nézve közömbös szavak is deskriptorok lesznek, ami a terjedelmet növeli, egyúttal a visszakeresést lassítja.

b./ A rendszerbe a releváns szavak (kifejezések) szótárát építjük és a tartalmi kivonat minden olyan kifejezése deskriptor lesz, amely a szótárban (legalábbis fő-alakban) megtalálható. Tezaurusz birtokában ez az út járható: a dokumentum-rekord belépésekor a tartalmi kivonat "értékelését" elvégezve, a visszakeresési eljárás alapját továbbra is a tezaurusz képviseli.

A tartalmi kivonathoz hasonló problémákat vet fel az igények természetes nyelven történő megfogalmazása, azzal a különbséggel, hogy bizonyos grammatikai információkra feltét-

lenül szükség van (és, vagy, nem szerepe, szavak sorrendje, stb.).

4.2. Az információs rendszer felhasználói igényeket értékelő része lényegében lista-kifejezéseket interpretáló program. A 2. és 3. pontban felsorolt lista-műveletek tovább bővíthetők: függvények definiálhatók, rekurzív eljárások írhatók fel, stb. Lényeges szerephez juthatnak a standard módon beépített makró-definíciók: bonyolult kifejezések az igények megfogalmazásakor csupán a makro névére és a paraméterekre hivatkoznak. Pl. téma(d) jelentse a  $\text{cons}(d \uparrow 0, d \uparrow 1)$  kifejezést, ami a visszakeresésben a  $d$  deskriptort és az összes rokon fogalmakat jelenti.

A makro-rendszer bővítései illetve a függvény-írásmód helyettesítése "szabadabb" formákkal, elvezethetnek a természetes nyelveken megfogalmazható igények közelébe.

Irodalomjegyzék

1. Általános műszaki fogalmak tezaurusza.  
Országos Műszaki Könyvtár és Dokumentációs Központ, 1970.
2. An Introduction to COMIT Programming. The Research Laboratory of Electronics and the Computation Center. The M.I.T. Press, 1962.
3. H.D.Avrar, J.F.Knapp, L.J.Rather: The MARC II. Format. A communications Format for Bibliographic Data. Information Systems Office, Library of Congress. Washington, D.C. January, 1968.
4. N.G. de Bruijn: The mathematical language AUTOMATH, its usage, and some of its extensions. Proceedings of the Symposium on Automatic Demonstration. IRIA, Versailles, 1968.
5. R.G.Busacker, T.L.Saaty: Véges gráfok és hálózatok; bevezetés alkalmazásokkal. Műszaki Könyvkiadó, Budapest, 1969.
6. A.Church: The Calculi of Lambda-Conversion. Annals of Mathematics Studies N<sup>o</sup> 6. Princeton University Press, 1941.
7. J.M.Foster: List Processing. Macdonald Computer Monographs, 1968.
8. R.R.Freeman, P.Atherton: File Organization and Search Strategy using the Universal Decimal Classification in Mechanized Reference Retrieval Systems. F.I.D./I.P.I.P. Conference on Mechanized Information Storage, Retrieval and Dissemination. Rome, June, 1967.
9. J.Friant: Les Languages "Contex-Sensitive". Ann. Inst. Henri Poincaré, Vol. III. N<sup>o</sup> 1, 1967. p. 35-120.
10. M.A.Gavrilov, A.D.Zakravskii: LYAPAS, A Programming Language for Logic and Coding Algorithms. Academic Press, 1969.
11. В.М.Плущков: Синтез цифровых автоматов. Государственное Издательство Физико-математической Литературы. Москва, 1962.
12. T.M.Hibbard: Some Combinatorial Properties of Certain Trees with Applications to Searching and Sorting. Journal of the ACM. January 1962.
13. Horváth Tibor: A könyvtárak műszaki fejlesztése. Budapest-Veszprém, 1969. OSzK, Könyvtártudományi és Módszertani Központ.
14. IBIS; System wyszukiwania informacji. Instytut Maszyn Matematycznych. 1968.
15. INBI; System przetwarzania na EMC informacji bibliograficznych Instytut Maszyn Matematycznych. 1968.

16. INFORGA Szimpoziium előadások tézisei. Moszkva, 1965.  
Országos Műszaki Könyvtár és Dokumentációs Központ, 1965.
17. L.Kalmár: Les Calculatrices automatiques comme structures algériques. Gauthier-Villars-Éditeur. 1963.
18. Kalmár László: Matematikai és nyelvi strukturák. Általános Nyelvészeti Tanulmányok II. Akadémiai Kiadó, Budapest, 1964
19. C.A.Kapps: SPRINT - A programming Language with General Structure. The Moor School of Electrical Engineering. Report N<sup>o</sup> 71-18. 1970.
20. Makay Árpád: Tanulmány a kisszámítógépek könyvtári és dokumentációs alkalmazásait illetően. JATE Kibernetikai Laboratóriuma, 1971.Kézirat.
21. Nicht-Numerische Informationverarbeitung. R.Gunzenhäser. Springer Verlag, Wien-New York. 1968.
22. J.O'Connor: Mechanised Indexing Methods and Their Testing. Journal of the ACM. V.11.N<sup>o</sup> 4. October 1964.
23. D.D.Prentice: The Combined File Search System. IBM Corporation, 1964.
24. Scientific, Technical and Economic Information in Poland. Central Institute for scientific, Technical and Economic Information. Warsaw, 1963.
25. Searching Normal Text for Information Retrieval. IBM Data Processing Application, E20-0335-0. 1969.
26. SESAM. Kézirat.
27. C.Weissmann: LISP 1.5. Primer. Dickenson Series in Computer and Information Science.

## Tartalomjegyzék

Bevezetés .....	1
1. Információs rendszerek .....	6
1.1. Alapvető fogalmak .....	6
1.2. Információs rendszerek feladatai .....	11
1.3. Információs rendszerek strukturája .....	13
1.4. Információvisszakeresés .....	16
2. Listák és elemi műveletek .....	21
2.1. Szöveges információk számítógépi ábrázolása .....	21
2.2. A lista fogalma .....	22
2.3. Atomok .....	24
2.4. Cellák .....	24
2.5. Listák .....	25
2.6. Cons .....	26
2.7. hd, tl .....	27
2.8. Összetett kifejezések .....	27
2.9. Listamódosító műveletek, asszociációs listák .....	28
2.10. Feltételes kifejezések .....	30
3. A tezaurusz listái .....	31
3.1. Azonosítók, atomok .....	31
3.2. Egy tezaurusz-elem listája .....	32
3.3. Példák .....	34
3.4. A tezaurusz számítógépi reprezentációja .....	35
3.5. A felhasználói igény megfogalmazása .....	36
3.6. Az igény szemantikája .....	36
3.7. Az információvisszakeresési eljárás .....	37
3.8. A dokumentum beépítése a rendszerbe .....	40
4. További lehetőségek .....	41
4.1. A tartalmi kivonatok feldolgozása .....	41
4.2. Az igények megfogalmazásának további lehetőségei .....	42
Irodalomjegyzék.....	43

