

A mű eredeti címe: Database System Implementation
First Edition by Hector Garcia-Molina, Jeffrey D. Ullman, Jennifer Widom, Copyright © 2000.

All rights reserved.

Published by arrangement with the original publisher Prentice Hall, Inc.,
a Pearson Education Company.

Translation copyright © 2001 by Panem Könyvkiadó Kft.

ISBN 963 545 280 2

A kiadásért felel a Panem Könyvkiadó Kft. ügyvezetője, 2001

Sz. keszlette: dr. Benczúr András

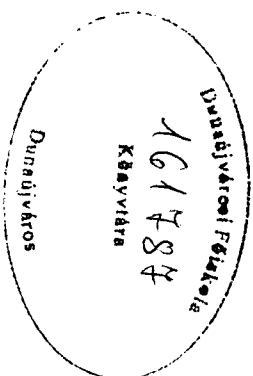
Fordította: Búza Attila, Cserges Enikő, dr. Hajas Csilla, dr. Kiss Attila, Kovács György,

Limbek Réka, Nikovits Tibor, Vincellér Zoltán

Lektorálta: dr. Márkus Tibor

Műszaki szerkesztő: Érdi Júlia

Borítóterv: Érdi Júlia



Tartalomjegyzék

1. Bevezetés az adatbázis-kezelő rendszerek implementálásába	24
Előszó	21
Előszó a magyar kiadáshoz	19
1.1. Bevezetés: a Megatron 2000 adatbázisrendszer	25
1.1.1. A Megatron 2000 implementálásának részletei	26
1.1.2. Hogyan hajítja végre a Megatron 2000 a lekérdezéseket?	27
1.1.3. Mi a baj a Megatron 2000-rel?	29
1.2. Egy adatbázis-kezelő rendszer áttekintése	30
1.2.1. Az adatdefiniációs nyelv parancsai	30
1.2.2. A lekérdezés feldolgozásának áttekintése	30
1.2.3. A központi memória pufferei és a puffervezelő	32
1.2.4. A tranzakció feldolgozása	33
1.2.5. A lekérdezésfeldolgozó	34
1.3. A könyv vázlatos felépítése	35
1.3.1. Előismeretek	35
1.3.2. A tárkezelés áttekintése	36
1.3.3. A lekérdezésfeldolgozás áttekintése	37
1.3.4. A tranzakciófeldolgozó áttekintése	37
1.3.5. Az információintegráció áttekintése	38
1.4. Az adatmodellek és nyelvek áttekintése	38
1.4.1. A relációs modell áttekintése	38
1.4.2. Az SQL áttekintése	39
1.4.3. A relációs és objektumorientált adatok	42

A Panem könyvek megrendelhető az (1) 340-1515 hívószámú telefonon, illetve az 1385 Budapest, Pf. 809 levélcímen.

panem@panem.hu

www.panem.hu

Minden jog fenntartva. Jelen könyvet, illetve annak részeit tilos reprodukálni, adatregisztráció rendszerben tárolni, bármilyen formában vagy eszközzel – elektronikus úton vagy más módon – közölni a kiadók engedélye nélkül.

1.5.	Összefoglalás	44	×2.5.3.	Stabil tárolás	92
1.6.	Irodalomjegyzék	44	2.5.4.	A stabil tárolás hibakezelő képessége	92
2.	Adattárolás	46	2.5.5.	Feladatok	93
2.1.	A memóriahierarchia	47	×2.6.	Lemezhiba helyreállítása	94
2.1.1.	Cache	47	2.6.1.	A lemezek meghibásodási modelljei	94
2.1.2.	A központi memória	48	2.6.2.	A túlközés mint redundanciatechnika	95
2.1.3.	Virtuális memória	49	2.6.3.	Paritásblokkok	96
2.1.4.	Másodlagos tárolás	51	2.6.4.	Egy továbbfejlesztés: az 5. szintű RAID	100
2.1.5.	Harmadlagos tárolás	52	2.6.5.	Mi a teendő, ha több lemez is tönkrement?	101
2.1.6.	Felejítő és nem felejítő tárolás	54	2.6.6.	Feladatok	104
2.1.7.	Feladatok	55	2.7.	Összefoglalás	107
× 2.2.	Lemezek	55	2.8.	Irodalomjegyzék	109
2.2.1.	A lemezek mechanikája	56	3.	Adatelemek ábrázolása	111
2.2.2.	A lemezvezérlő	57	×3.1.	Adatelemek és mezők	111
2.2.3.	A lemeztárolók jellemzői	58	3.1.1.	Relációs adatbázisiselemek ábrázolása	112
2.2.4.	A lemezhozzáférés jellemzői	60	3.1.2.	Objektumok ábrázolása	113
2.2.5.	Blokkok írása	64	3.1.3.	Adatelemek ábrázolása	114
2.2.6.	Blokkok módosítása	64	×3.2.	Rekordok	119
2.2.7.	Feladatok	65	3.2.1.	Rögzített hosszú rekordok építése	119
2.3.	A másodlagos tárolók hatékony használata	66	3.2.2.	Rekordfejlécek	121
2.3.1.	A számtítás I/O-modellje	66	3.2.3.	Rögzített hosszú rekordok blokkokba pakolása	123
2.3.2.	Adatok rendezése a másodlagos tárolóban	67	3.2.4.	Feladatok	124
2.3.3.	Az összefésült rendezés (Merge-Sort)	68	3.3.	Blokkcímek és rekordcímek ábrázolása	125
2.3.4.	Kétfázisú, többitas, összefésült rendezés	70	3.3.1.	Kliens-szerver rendszerek	125
2.3.5.	A többitas összefésültés kiterjesztése nagyobb relációkra	72	3.3.2.	Logikai és strukturált címek	127
2.3.6.	Feladatok	74	3.3.3.	Mutatók helyreigazítása	129
2.4.	A másodlagos tároló hozzáférési idejének javítása	75	3.3.4.	Blokkok visszairása a lemezre	133
2.4.1.	Az adatok cilindres szervezése	77	3.3.5.	Feltöltött rekordok és blokkok	134
2.4.2.	Több lemez használata	78	3.3.6.	Feladatok	135
2.4.3.	Lemezek túlközése	79	3.4.	Változó hosszú adatok és rekordok	137
2.4.4.	A lemez ütemezése és a lift algoritmus	80	3.4.1.	Változó hosszú mezőket tartalmazó rekordok	138
2.4.5.	Korai beolvasás és nagy léptékű pufferezés	84	3.4.2.	Ismétlődő mezőket tartalmazó rekordok	139
2.4.6.	A stratégiák előnyeinek és hátrányainak összegzése	86	3.4.3.	Változó formátumú rekordok	141
2.4.7.	Feladatok	87	3.4.4.	Olyan rekordok, amelyek nem férnek el egy blokkban	142
2.5.	Lemezhibák	89	3.4.5.	Bináris, nagy objektumok (BLOB-ok)	143
2.5.1.	Ideiglenes meghibásodás	90	3.4.6.	Feladatok	144
2.5.2.	Ellenőrző összegek	90			

3.5.	Rekordmódosítások	146
3.5.1.	Beszűrés	146
3.5.2.	Törlés	148
3.5.3.	Módosítás	149
3.5.4.	Feladatok	150
3.6.	Összefoglalás	151
3.7.	Irodalomjegyzék	152

4. Indexstruktúrák

4.1.	Indexek szekvenciális fájlokon	154
4.1.1.	Szekvenciális fájlok	155
4.1.2.	Sűfű indexek	155
4.1.3.	Ritka indexek	158
4.1.4.	Több szintű indexelés	159
4.1.5.	Indexelés ismétlődő kereséskulcs-érték esetén	161
4.1.6.	Indexek kezelése adatmódosításkor	164
4.1.7.	Feladatok	170
× 4.2.	Másodlagos indexek	171
4.2.1.	Másodlagos indexek tervezése	172
4.2.2.	Másodlagos indexek alkalmazása	173
4.2.3.	Közvetett másodlagos indexek	175
4.2.4.	Dokumentumok visszakeresése és az invertált indexek	178
4.2.5.	Feladatok	181
4.3.	B-fák	184
4.3.1.	B-fák szerkezete	184
4.3.2.	B-fák alkalmazása	187
4.3.3.	Keresés B-fában	189
4.3.4.	Tartományra vonatkozó lekérdezések	190
4.3.5.	Beszűrés B-fában	191
4.3.6.	Törlés B-fában	194
4.3.7.	B-fák hatékonysága	197
4.3.8.	Feladatok	197
4.4.	Tördelőtáblázatok	200
4.4.1.	Másodlagos tárolón tárolt tördelőtáblázatok	201
4.4.2.	Beszűrés tördelőtáblázatra	202
4.4.3.	Törlés tördelőtáblázaton	202
4.4.4.	Tördelőtáblázat-indexek hatékonysága	203
4.4.5.	Kiterjeszhető tördelőtáblázatok	204

4.4.6.	Beszűrés kiterjeszhető tördelőtáblázatokba	205
4.4.7.	Lineáris tördelőtáblázatok	207
4.4.8.	Beszűrés lineáris tördelőtáblázatokba	209
4.4.9.	Feladatok	211

4.5.	Összefoglalás	213
4.6.	Irodalomjegyzék	214

5. Többdimenziós indexek

5.1.	Többdimenziós alkalmazások	217
5.1.1.	Térinformaikai rendszerek	217
5.1.2.	Adatcockák	218
5.1.3.	Többdimenziós lekérdezések SQL-ben	219
5.1.4.	Tartománylekérdezések végrehajtása hagyományos indexekkel	221
5.1.5.	Legközelebbi szomszéd-lekérdezések végrehajtása hagyományos indexekkel	222
5.1.6.	A hagyományos indexek további korlátjai	224
5.1.7.	A többdimenziós indexstruktúrák áttekintése	224
5.1.8.	Feladatok	225
5.2.	Tördelésen alapuló struktúrák többdimenziós adatokhoz	226
5.2.1.	Rácsos állományok	227
5.2.2.	Keresés rácsos állományban	227
5.2.3.	Beszűrés rácsos állományba	229
5.2.4.	A rácsos állományok hatékonysága	230
5.2.5.	Particionált tördelőfüggvények	233
5.2.6.	A rácsos állományok és a particionált tördelés összehasonlítása	234
5.2.7.	Feladatok	235
5.3.	Faszertű struktúrák többdimenziós adatokhoz	238
5.3.1.	Többkulcsos indexek	238
5.3.2.	A többkulcsos indexek hatékonysága	240
5.3.3.	kd-fák	241
5.3.4.	Műveletek a kd-fákon	243
5.3.5.	A kd-fák alkalmazása másodlagos tárolók esetén	245
5.3.6.	Quad-fák	246
5.3.7.	R-fák	248
5.3.8.	Műveletek az R-fákon	249
5.3.9.	Feladatok	251
5.4.	Bitképindexek	253
5.4.1.	Indítékok a bitképindexekhez	254

5.4.2.	Tömörített bitértképek	256
5.4.3.	Műveletek szakszösszkódolt bitvektorokon	258
5.4.4.	Bitértképzindexek kezelése	259
5.4.5.	Feladatok	260
5.5.	Összefoglalás	261
5.6.	Irodalomjegyzék	263
6.	Lekérdezések végrehajtása	265
× 6.1.	Algebrai megközelítés	267
6.1.1.	Egyesítés, metszet és különbség	269
6.1.2.	Kiválasztás	270
6.1.3.	Vetítés	272
6.1.4.	Relációk szorzata	273
6.1.5.	Összekapcsolások	274
6.1.6.	Ismétlődések kiküszöbölése	276
6.1.7.	Csoportosítás és összesítés	277
6.1.8.	Rendezés	279
6.1.9.	Kifejezések	280
6.1.10.	Feladatok	282
6.2.	Bevezetés a fizikai lekérdezésv- operátorok világába	285
6.2.1.	Táblák átvizsgálása	285
6.2.2.	Rendezés a táblák átvizsgálásakor	286
6.2.3.	A fizikai operátorok kiszámításának modellje	287
6.2.4.	A költségbecslés paraméterei	287
6.2.5.	Az átvizsgáló operátorok I/O-költsége	289
6.2.6.	Fizikai operátorok megvalósításához használatos iterátorok	290
> 6.3.	Adatbázis- műveletek egyemenetes algoritmusai	293
6.3.1.	Soronkénti műveletek egyemenetes algoritmusai	294
6.3.2.	Unáris, teljes relációs műveletek egyemenetes algoritmusai	295
6.3.3.	Bináris műveletek egyemenetes algoritmusai	299
6.3.4.	Feladatok	302
× 6.4.	Beágyazott ciklusú összekapcsolások	303
6.4.1.	Sor alapú beágyazott ciklusú összekapcsolás	304
6.4.2.	Egy iterátor a sor alapú beágyazott ciklusú összekapcsoláshoz	304
6.4.3.	Egy algoritmus a blokk alapú beágyazott ciklusú összekapcsoláshoz	304
6.4.4.	A beágyazott ciklusú összekapcsolás elemzése	307
6.4.5.	Az eddigi algoritmusok összefoglalása	307
6.4.6.	Feladatok	308
< 6.5.	Rendezésen alapuló kétemenetes algoritmusok	308
6.5.1.	Ismétlődések kiküszöbölése rendezés segítségével	309
6.5.2.	Csoportosítás és összesítés rendezés segítségével	312
6.5.3.	Az egyesítés egy rendezésen alapuló algoritmusa	312
6.5.4.	A metszet és a különbség rendezésen alapuló algoritmusa	313
6.5.5.	Egy egyszerű rendezésen alapuló összekapcsolási algoritmus	315
6.5.6.	Az egyszerű rendezés összekapcsolás elemzése	317
6.5.7.	Egy hatékonyabb rendezésen alapuló összekapcsolás	317
6.5.8.	A rendezésen alapuló algoritmusok összefoglalása	319
6.5.9.	Feladatok	319
6.6.	Tördelésen alapuló kétemenetes algoritmusok	321
6.6.1.	Relációk particionálása tördeléssel	321
6.6.2.	Egy tördelésen alapuló algoritmus az ismétlődések kiküszöbölésére	322
6.6.3.	Egy tördelésen alapuló algoritmus a csoportosításra és az összesítésre	323
6.6.4.	Az egyesítés, a metszet és a különbség tördelésen alapuló algoritmusa	323
6.6.5.	A tördeléses összekapcsolási algoritmus	324
6.6.6.	Lemez I/O-műveletek megtakarítása	325
6.6.7.	A tördelésen alapuló algoritmusok összefoglalása	327
6.6.8.	Feladatok	328
< 6.7.	Index alapú algoritmusok	329
6.7.1.	Nyaláboltt és nem nyáláboltt indexek	329
6.7.2.	Index alapú kiválasztás	330
6.7.3.	Összekapcsolás index segítségével	333
6.7.4.	Összekapcsolások rendezett index segítségével	334
6.7.5.	Feladatok	336
6.8.	Pufferkezelés	337
6.8.1.	A pufferkezelő működése	338
6.8.2.	Pufferkezelő stratégiák	338
6.8.3.	Kapcsolat a fizikai operátor kiválasztása és a pufferkezelés között	340
6.8.4.	Feladatok	342
6.9.	Több mint kétemenetes algoritmusok	343
6.9.1.	Többmenetes, rendezésen alapuló algoritmusok	343
6.9.2.	Többmenetes, rendezésen alapuló algoritmusok műveletigénye	344
6.9.3.	Többmenetes, tördelésen alapuló algoritmusok	345
6.9.4.	Többmenetes, tördelésen alapuló algoritmusok műveletigénye	345
6.9.5.	Feladatok	346
6.10.	Párhuzamos algoritmusok relációs műveletekre	347
6.10.1.	A párhuzamosság modelljei	347
6.10.2.	Soronkénti műveletek párhuzamos megvalósítása	350
6.10.3.	Teljes relációs műveletek párhuzamos algoritmusai	351

6.10.4.	A párhuzamos algoritmusok hatékonysága	352
6.10.5.	Feladatok	355
6.11.	Összefoglalás	356
6.12.	Irodalomjegyzék	358
7.	A lekérdezésfordító	359
7.1.	Elemzés	360
7.1.1.	Szimakritikus elemzés és elemzőfák	360
7.1.2.	Egy leegyszerűsített SQL-részletet leíró nyelven	361
7.1.3.	Az előfeldolgozó	366
7.1.4.	Feladatok	367
7.2.	Algebrai szabályok lekérdezéstervek javítására	367
7.2.1.	Kommutatív és asszociatív szabályok	368
7.2.2.	Kiválasztással kapcsolatos szabályok	371
7.2.3.	Kiválasztások tologatása	374
7.2.4.	Vetítéssel kapcsolatos szabályok	375
7.2.5.	Összekapcsolásra és szorzatra vonatkozó szabályok	379
7.2.6.	Ismétlődések elhagyására vonatkozó szabályok	379
7.2.7.	Csoportosításra és összesítésre vonatkozó szabályok	380
7.2.8.	Feladatok	382
7.3.	Elemzőfák átalakítása logikai lekérdezéstervekké	384
7.3.1.	Átfordítási relációs algebra	384
7.3.2.	Alkérdesések elhárulása feltételekből	386
7.3.3.	Logikai lekérdezéstervek javítása	391
7.3.4.	Asszociatív/kommutatív operátorok csoportosítása	393
7.3.5.	Feladatok	394
7.4.	Műveletek költségének becslése	395
7.4.1.	Közbülső relációk méretének becslése	396
7.4.2.	Vetítés méretének becslése	397
7.4.3.	Kiválasztás méretének becslése	398
7.4.4.	Összekapcsolás méretének becslése	401
7.4.5.	Természetes összekapcsolás több összekapcsolási attribútummal	403
7.4.6.	Sok reláció összekapcsolása	405
7.4.7.	Egyéb műveletek méretének becslése	406
7.4.8.	Feladatok	409
7.5.	Bevezetés a költség alapú tervválasztásba	410
7.5.1.	Méretre vonatkozó paraméterek becslése	411

7.5.2.	Statisztikák növekményes kiszámítása	414
7.5.3.	Logikai lekérdezéstervek költségének csökkentésére irányuló heurisztikák	416
7.5.4.	Fizikai tervek felsorolásának lehetőségei	418
7.5.5.	Feladatok	421
7.6.	Összekapcsolások sorrendjének megválasztása	423
7.6.1.	Összekapcsolások bal és jobb oldali argumentumának jelentősége	423
7.6.2.	Összekapcsolási fák	424
7.6.3.	Bal-mély összekapcsolási fák	425
7.6.4.	Dinamikus programozás az összekapcsolási sorrend és csoportosítás megválasztására	428
7.6.5.	Dinamikus programozás részletesebb költségfüggvényekkel	433
7.6.6.	Egy mohó algoritmus az összekapcsolási sorrend kiválasztására	434
7.6.7.	Feladatok	435
7.7.	A fizikai lekérdezésterv kiválasztásának befejezése	437
7.7.1.	Kiválasztási eljárás megválasztása	437
7.7.2.	Összekapcsolási eljárás megválasztása	440
7.7.3.	Futószalagosítás és materializáció	441
7.7.4.	Unáris műveletek futószalagosítása	442
7.7.5.	Bináris műveletek futószalagosítása	443
7.7.6.	Fizikai lekérdezéstervekkel kapcsolatos jelölések	445
7.7.7.	Fizikai operátorok sorrendbe állítása	449
7.7.8.	Feladatok	449
7.8.	Összefoglalás	451
7.9.	Irodalomjegyzék	452
8.	A rendszerhibák kezelése	454
8.1.	A helyreállítható beavatkozások példái és modelljei	454
8.1.1.	A hibák fajtái	455
8.1.2.	Részletesebben a tranzakciókról	457
8.1.3.	A tranzakciók korrekt végrehajtása	458
8.1.4.	A tranzakciók alapvetékenységei	460
8.1.5.	Feladatok	463
8.2.	Semmisségi (undo) naplózás	463
8.2.1.	Naplóbejegyzések	465
8.2.2.	A semmisségi naplózás szabályai	466
8.2.3.	Helyreállítás a semmisségi naplózás használatával	468
8.2.4.	Az ellenőrzőpont-képzés	471

8.2.5.	Ellenőrzőpont-képzés a rendszer működése közben	473
8.2.6.	Feladatok	476
8.3.	Helyrehozó naplózás (redo logging)	477
8.3.1.	A helyrehozó naplózás szabályai	478
8.3.2.	Helyreállítás a helyrehozó naplózás használatával	479
8.3.3.	Helyrehozó naplózás ellenőrzőpont-képzés használatával	480
8.3.4.	Visszaállítás az ellenőrzőponttal kiegészített helyrehozó típusú naplózással	482
8.3.5.	Feladatok	483
8.4.	A semmisségi/helyrehozó (undo/redo) naplózás	484
8.4.1.	A semmisségi/helyrehozó (undo/redo) naplózás szabályai	484
8.4.2.	Helyreállítás a semmisségi/helyrehozó (undo/redo) naplózás használatakor	485
8.4.3.	Semmisségi/helyrehozó naplózás ellenőrzőpont-képzéssel	487
8.4.4.	Feladatok	489
8.5.	Az eszközök meghibásodása elleni védekezés	490
8.5.1.	Az archivmentés	490
8.5.2.	Archiválás működés közben	491
8.5.3.	Helyreállítás az archivmentés és a napló használatával	494
8.5.4.	Feladatok	495
8.6.	Összefoglalás	495
8.7.	Irodalomjegyzék	497
9.	Konkurenciavezérlés	498
9.1.	Soros és sorba rendezhető ütemezések	499
9.1.1.	Ütemezések	499
9.1.2.	Soros ütemezések	500
9.1.3.	Sorba rendezhető ütemezések	501
9.1.4.	A tranzakció szemantikájának hatása	503
9.1.5.	A tranzakciók és ütemezések jelölése	504
9.1.6.	Feladatok	505
9.2.	Konfliktus-sorbarendeazhetőség	505
9.2.1.	Konfliktusok	506
9.2.2.	Megelőzési gráfok és teszi a konfliktus-sorbarendeazhetőségre	507
9.2.3.	Miért működik a megelőzési gráfon alapuló tesztelés?	510
9.2.4.	Feladatok	511

9.3.	A sorbarendeazhetőség biztosítása zárákkal	513
9.3.1.	Zárak	514
9.3.2.	A zárolási ütemező	516
9.3.3.	A kétfázisú zárolás	517
9.3.4.	Miért működik a kétfázisú zárolás?	517
9.3.5.	Feladatok	519
9.4.	Különböző zármódú zárolási rendszerek	521
9.4.1.	Osztott és kizárólagos zárok	521
9.4.2.	Kompatibilitási mátrixok	523
9.4.3.	Zárak felminősítése	524
9.4.4.	Módosítási zárok	526
9.4.5.	Növelési zárok	527
9.4.6.	Feladatok	529
9.5.	A zárolási ütemező felépítése	532
9.5.1.	Zárolási műveleteket beszűrő ütemező	532
9.5.2.	A zártábla	534
9.5.3.	Feladatok	537
9.6.	Adatbáziselemekből álló hierarchiák kezelése	538
9.6.1.	Többszörös szemcsézettű zárok	538
9.6.2.	A figyelmeztető zárok	539
9.6.3.	Fantomok és a beszűrások helyes kezelése	542
9.6.4.	Feladatok	543
9.7.	Faprotokoll	544
9.7.1.	Fa alapú zárolások idtékeai	544
9.7.2.	Faszervezetű adatok hozzáférési szabályai	545
9.7.3.	Miért működik a faprotokoll?	546
9.7.4.	Feladatok	549
9.8.	Konkurenciavezérlés időbélyegzőkkel	550
9.8.1.	Időbélyegzők	551
9.8.2.	Fizikailag nem megvalósítható viselkedések	552
9.8.3.	Piszkos adatok problémái	553
9.8.4.	Az időbélyegzőn alapuló ütemezések szabályai	554
9.8.5.	Többváltozatú időbélyegzők	556
9.8.6.	Az időbélyegzők és zárolások	558
9.8.7.	Feladatok	559
9.9.	Konkurenciavezérlés érvényesítéssel	560
9.9.1.	Érvényesítésen alapuló ütemező felépítése	560
9.9.2.	Az érvényesítési szabályok	561
9.9.3.	Három konkurenciavezérlés működésének összehasonlítása	564

9.9.4.	Feladatok	565
9.10.	Összefoglalás	565
9.11.	Irodalomjegyzék	568
10.	Bővebben a tranzakciókezelésről	569
10.1.	Tranzakciók, melyek nem véglegesített adatokat olvasnak	569
10.1.1.	A piszkos adat probléma	570
10.1.2.	Továbbgyűrűző visszagörgetés	572
10.1.3.	A visszagörgetés kezelése	573
10.1.4.	Csoportos véglegesítés	574
10.1.5.	Logikai naplózás	576
10.1.6.	Feladatok	579
10.2.	Nézet-sorbarendeázhetőség	580
10.2.1.	Nézetekvivalencia	580
10.2.2.	Poligráfok és nézet-sorbarendeázhetőségi teszt	582
10.2.3.	A nézet-sorbarendeázhetőség tesztelése	585
10.2.4.	Feladatok	585
10.3.	Holtpontkezelés	586
10.3.1.	Holtpontkezelés időkorlátal	586
10.3.2.	A várakozási gráf	587
10.3.3.	Holtpontmegelőzés az elemek sorbarendeázásával	589
10.3.4.	Holtpontkezelés időbélyegzővel	591
10.3.5.	A holtpontkezelő módszerek összehasonlítása	593
10.3.6.	Feladatok	594
10.4.	Osztott adatbázisok	595
10.4.1.	Osztott adatok	596
10.4.2.	Osztott tranzakciók	597
10.4.3.	Adatöbbszörözés	598
10.4.4.	Osztott lekérdeázésopimnalizálás	599
10.4.5.	Feladatok	600
10.5.	Osztott véglegesítés	606
10.5.1.	Az osztott atomosság támogatása	601
10.5.2.	Kétfázisú véglegesítés	601
10.5.3.	Az osztott tranzakciók helyreállítás	604
10.5.4.	Feladatok	606

10.6.	Osztott zárolás	607
10.6.1.	Központosított zárolási rendszerek	607
10.6.2.	Költségmodell az osztott zárolási algoritmusokhoz	608
10.6.3.	Többszörözött elemek zárolása	609
10.6.4.	Az elsődleges példány zárolása	610
10.6.5.	A lokális zárolás a globálisig	611
10.6.6.	Feladatok	612
10.7.	Hosszú tranzakciók	613
10.7.1.	A hosszú tranzakciók problémái	613
10.7.2.	Regék	616
10.7.3.	Kiegyenlítő tranzakciók	617
10.7.4.	Miért működnek jól a kiegyenlítő tranzakciók?	619
10.7.5.	Feladatok	619
10.8.	Összefoglalás	620
10.9.	Irodalomjegyzék	622
11.	Információk egyesítése	624
11.1.	Az információegyesítés módjai	624
11.1.1.	Az egyesítés problémái	625
11.1.2.	Adatbázis-szövetség	627
11.1.3.	Adattárházak	628
11.1.4.	Adatközvetítő	631
11.1.5.	Feladatok	633
11.2.	Borítékolók a közvetítő alapú rendszerekben	635
11.2.1.	Sablonok lekérdeázési formákhoz	635
11.2.2.	Borítékoló generátor	636
11.2.3.	Szűrők	637
11.2.4.	A borítékoló más műveletei	639
11.2.5.	Feladatok	641
11.3.	On-line analitikus feldolgozás	641
11.3.1.	OLAP-alkalmazások	643
11.3.2.	OLAP-adatok többdimenziós nézete	644
11.3.3.	A csillag séma	645
11.3.4.	Szeletelés és kockázás	647
11.3.5.	Feladatok	650
11.4.	Adatkocka	651
11.4.1.	A kockaművelet	652

11.4.2.	Kockaimplementáció megvalósított nézetablakkal	655
11.4.3.	Nézetabló	658
11.4.4.	Feladatok	660
11.5.	Adatbányászat	662
11.5.1.	Adatbányászati alkalmazások	663
11.5.2.	Társtíási szabály bányászat	666
11.5.3.	Az előzetes algoritmus	667
11.6.	Összefoglalás	670
11.7.	Irodalomjegyzék	671

Index	675
--------------------	-----

Előszó a magyar kiadáshoz

A Stanford Egyetem három neves számítástudosa újabb tankönyvvel jelentkezett az adatbázis-kezelő rendszerek témakörében. Az első, a magyar fordításban *Adatbázis-rendszerek. Alapvetés* címen megjelent könyvet most a szerzőknek az adatbázisrendszerek megvalósítási kérdéseiről írt könyve követi. Az előző könyv fordítása során szerzett igen jó véleményünk és a könyv sikere alapján figyelemmel kísértük az előre jelzett új könyv megjelenését, és ahogy megjelent, a Panem–Prentice–Hall Kiadóval azonnal felvettük a kapcsolatot a magyar fordítás kiadásának érdekében. Jólllehet most nem kaptunk támogatást az Oktatási Minisztérium felsőoktatási tankönyv pályázatán, mégis mi, mint a tárgy oktatói fontosnak tartottuk a könyv magyar megjelenését, és a Kiadó legnagyobb örömindkre vállalkozott rá.

Ahogy az előző könyv is hiánypótló volt, ez is az, bár ez már elsősorban azoknak a szakembereknek szól, akik nagy adatbázisok működtetéséért felelősek, megvalósításán, hatékonyabbá tételén dolgoznak.

Folytatva az előző könyv magyar kiadásához írt előszó jéghegy hasonlatát, az adatbázis-kezelő rendszerekről a jéghegy csúcsa után most az alap bemutatására kerül sor. A jéghegy mélyén a megvalósítás technikái húzódnak. Amint a könyv bemutat, az közel 50 év fejlődése során a közvetlen elérésű háttértárolók megjelenését követő fejlődés folyamataiban leisztult megoldások sora. Erőteljes hangszólyt kap ebben a relációs adatbázis-kezelők tipikus lekérdézőoptimalizálási világa, de a legújabb, az internet lehetőségeire alapuló hálózati adatbázis-kezelés elemei is megjelennek már.

A tárolási adatszerkezetek bemutatása során a lemezelérés és adatehelyezés részletei, a hibajavítás és a RAID (független lemezek redundáns tömbje) lemezezségek, az adatelemelek tárolási módjai, az indexelési technikák a legfejlettebb többdimenziós indexszerkezetekig bezárólag kerülnek bemutatásra.

A lekérdezések feldolgozásának megvalósítását az elemi relációs algebrai műveletek kiértékelésének részletes elemzése vezeti be, majd az összetett lekérdezések kiértékelésének elemző, optimalizáló módszerei következnek. A kérdés átfutása hatékonyabban kiértékelhető ekvivalens relációs algebrai kifejezéssé, a logikai lekérdezési terv összeállítás, a műveletek költségbecslése, majd a költségbecslésre alapuló kiértékelési tervváltásztás módszere, az összekapcsolások optimális sorrendje és végül a fizikai kiértékelési terv megadása található a módszerek között.

A rendszerhibák és a kivételükre szolgáló naplózási technikák bemutatását, a konkurencia ellenőrzését, a tranzakciók kezelését részletesen leíró két fejezet követi. Az ütemezések sorbarendehezhetősége, ennek ellenőrzését és garantálását biztosító zárolási és időbélyegzős technikák bemutatása, az osztott adatbázisok tranzakciókezelési módszerei a legfontosabb témái ezeknek a fejezeteknek.

Az utolsó, 11. fejezet a legújabb, nagy teljesítményű rendszerek megvalósításairól ad ízelítőt. A nagy adatbázisok, a különböző forrású adatok közös kezelését biztosító információintegrációs módok sorában bemutatja a szövetséges adatbázisok, az adattárházak és a közvetítők (mediátorok) módszerét. A nagy erőforrás-igényű on-line döntésmozgató rendszerek (OLAP) működését elősegítő adatszerkezetek, például a csillag sémák és adatkockák bemutatása után a könyvet az adathányászati elemek zárják.

A könyv igen jól példázza azt, ahogy az információs technológiák fejlődnek: sok apró, néha egyszerűnek tűnő részlet, ami közel 50 év fejlődése során tisztult le a háttérben tudományos elemzésekkel. A fejlődés során többszörös rétegekben rakódnak egymásra a technológia elemei, az alkalmazásoknál közvetlenül látható réteg mögött ezekről gyakran megfeledkezünk. Az adatbázisrendszerek esetében ez igen veszélyes lehet, amire a könyv 2.1. részében, a memóriahierarchiák ismertetésénél idézett Moore törvénye hívja fel a figyelmet: míg a processzor sebessége, a memória és a háttértároló sűrűsége 18 hónaponként megkétszereződik, és az egységnyi teljesítmény ára feleződik, addig a memória elérésének és a lemezek forgásának sebessége alig növekszik! A relatív távolság a processzor, a memória és a háttértárolók között növekszik. A háttértárolókkal való gazdálkodás, ami az adatbázisrendszer megvalósításának központi kérdésköre, továbbra is meghatározó tényezője lesz a rendszerfejlesztéseknek.

A könyv szakembereknek szólóan dolgozza fel a fentiekben vázolt témákat. Részletesen, példákon keresztül mutatja be az általában aprólétkos lépésekből álló módszereket, eljárásokat. Széles spektrumát fogja át az ismereteknek és a technológiának. Az oktatásban alaposan kipróbált tárgyalásmód igen jól olvasható szakirodalmat jelent mind a hallgatóknak, mind a szakembereknek. Az adatbázis-tervezők, -felhasználók és -alkalmazási programozók számára a három neves szakember igen sok gyakorlati tanácsot nyújt a korszerű adatbázisrendszerek megvalósításához és használatához.

Azok számára, akik az alapozó általános ismereteken túl kívánnak adatbázisokkal foglalkozni, feltétlenül ajánlatos a könyv áttanulmányozása. A hazai felsőoktatásban a programtervező matematikus szakon az adatbázis-előadások eddig is sok mindent lefedtek a könyv témáiból, most végre magyar nyelvű tankönyvként is rendelkezésre áll egy igen kiforrott tananyag. Javasolható a könyv a műszaki informatikus képzés fel-sőbb évfolyamaiban is haladó adatbáziskurzus alapkönyveként.

A magyar nyelvű változat szinte egy éven belül követi a könyv megjelenését, amiért köszönet illeti a Panem Könyvkiadó gyors döntését a kiadás vállalásában, s leginkább a fordítók és a lektor áldozatvállalását, hogy nyári szabadságuk jelentős részének feláldozásával tartani tudták az igen rövid határidőt, és igen jó minőséggel végezték munkájukat. A fordítást végző csapat most is az Eötvös Loránd Tudományegyetem Információs Rendszerek Tanszékének oktatóiból, jelenlegi és volt doktorandus hallgatóiból alakult.

Dr. Benczúr András

Előszó

Ezt a könyvet a CS245 kurzus számára terveztük, amely a Stanford Egyetemen az adatbázisok sorozatban a második kurzus. Az első adatbáziskurzus itt a CS145, amely az adatbázis-tervezés és -programozás témaköröket fedi le, és Jeffrey D. Ullman és Jennifer Widom ehhez írta a Prentice-Hall Kiadónál 1997-ben megjelent *A First Course in Database Systems*¹ című könyvet. A CS245 kurzus ezt követően az ABKR-ek megvalósítását veszi át, nevezetesen a tárolási szerkezeteket, a kérdésfeldolgozást és a tranzakciókezelést.

A könyv használata

Stanfordban negyedéves, quarter rendszerünk van, ezért a könyvet használó alapkurzus, a CS245, mindössze 10 héig tart. 1999 telén Hector Garcia-Molina a könyv „béta” verzióját használta, és a következő részeket adta le: a 2.1.–2.4. részeket, a teljes 3. és 4. fejezetet, az 5.1. és 5.2., a 6.1.–6.7. és a 7.1.–7.4. részeket, a teljes 8. és 9. fejezetet, kivéve a 9.8. részt, a 10.1.–10.3., a 11.1. és a 11.5. részeket.

A 6. és 7. fejezetek többi része (lekérdezések optimalizálása) egy haladó kurzusban, a CS346-ban szerepel, melynek során a hallgatók saját ABKR-t készítenek. A könyv további olyan részei, amelyek nem szerepelnek a CS245-ben, egyéb haladó kurzusban kaphatnak helyet, mint például a CS347, amely az osztott adatbázisokat és a fejlett tranzakciókezelést tárgyalja.

Olyan intézmények, amelyek féléves, szemeszterrendszerben oktatnak, élhetnek azzal a lehetőséggel, hogy kombinálják ezt a könyvet az elődjével. A *First Course in Database Systems* könyvvel. Azt javasoljuk, hogy azt a könyvet az első szemeszterben használják egy adatbázis-alkalmazási projekthez társítva. A második szemeszter lefedheti ennek a könyvnek legnagyobb részét. Annak előnye, hogy az adatbázisok tár-

¹ Szerkesztői megjegyzés: A könyv megjelent magyar fordításban, *Adatbázisrendszerek. Alapvetés* címen a Panem–Prentice-Hall kiadók gondozásában 1998-ban.

gyat két kurzusra bontjuk az, hogy azok a hallgatók, akik nem szándékoznak adatbázisokra specializálódni, beérhessék csak az első kurzus választásával, és képesek legyenek az adatbázisok használatára, bármilyen számfűtádományi területre lépjenek is.

Előismeretek

Azt a kurzust, amely erre a könyvre alapul, ritkán választják a negyedik (senior) év előtt, ezért elvárjuk, hogy a hallgató elég széles háttérrel rendelkezzen a számfűtádomány tradicionális területein. Feltételezzük, hogy az olvasó már tanult valamennyit az adatbázis-programozásról, speciálisan az SQL-ről. Sokat segít a relációs algebra ismerete, valamint az alapvető adatstruktúrákkal való ismeretség. Hasonlóan a fájl-rendszernek és az operációs rendszerek bizonyos mértékű ismerete is hasznos.

Feladatok

A könyv terjedelmes feladatrendszer tartalmaz, szinte minden szakaszhoz tartozik néhány feladat. A nehezebb feladatokat vagy feladatrészeket felkiáltójel jelöli. A legnehezebb feladatok két felkiáltójelre kaptak.

Néhány feladat vagy feladatrész csillag megjelölést kapott. Arra törekszünk, hogy ezeknek a megoldásait a könyv weblapján hozzáférhetővé tegyük. Ezek a megoldások nyilvánosan elérhetőek, és segítenek az önellenőrzésben. Felhívjuk még a figyelmet arra, hogy vannak olyan esetek, amikor egy *B* feladat igényli az olvasó által egy korábbi *A* feladatra adott megoldás módosítását vagy felhasználását. Amennyiben *A* valamely részére weben publikált megoldás található, fellehető, hogy a *B* megfellelő részére is elérhető a megoldás.

Támogatás a World Wide Weben

A könyv honlapja:

<http://www-db.stanford.edu/~ullman/dbsi.html>

Itt megtalálhatók a csillaggal jelölt feladatok megoldásai, a hibajavítások, ahogy értesülünk róluk, valamint segédanyagok. Reményeink szerint elérhetővé tesszük minden meghirdetett CS245 kurzusunkhoz a kiosztott jegyzeteket, és az egyéb adatbáziskurzusok lényeges anyagait úgy, ahogy, tanítjuk őket, beleértve a házi feladatokat, zárthelyiket és megoldásokat.

Köszönetnyilvánítások

Köszönettel tartozunk Brad Adelberg, Karen Butler, Ed Chang, Surajit Chaudhuri, Rada Chirkova, Tom Dienstbier, Xavier Faz, Tracy Fujieda, Luis Gravano, Ben Holzman, Fabien Modoux, Peter Mork, Ken Ross, Mema Roussopolous és Jonathan Ullman segítségéért, amelyet az anyag összegyűjtésében és/vagy a mű korábbi változatban lévő hibák felfedezésében nyújtottak. A megmaradt hibák természetesen a miénk.

Hector Garcia-Molina

Jeffrey D. Ullman

Jennifer Widom

Stanford, CA

Bevetés az adatbázis-kezelő rendszerek implementálásába

Az adatbázisok manapság az üzleti élet minden területén alapvető fontosságot kapnak. Éppúgy használatosak saját feljegyzések kezelésére, mint a világhálós (World Wide Web) kereskedelemben, ahol az ügyfelek és vásárlók számára kell adatokat szolgáltatni. Emellett természetesen sok egyéb kereskedelmi, pénzügyi folyamatot is szolgáltatók bázisokkal segítenek. Hasonló módon adatbázisokat találunk sok tudományos vizsgálat mélyén is. Ezek az adatbázisok olyan adatokat reprezentálnak, amelyeket sok tudós gyűjtött össze, például csillagászok, genetikusok vagy éppen a fehérjék gyógyhatású tulajdonságait kutató biokémikusok.

Az adatbázisok ereje a több évtizeden keresztül kifejlődött tudásanyagának és technológiájának köszönhető. Ez a tudás azokban a speciális szoftverekben ölt testet, amelyeket *adatbázis-kezelő rendszereknek* (DBMS – database management system) vagy röviden „adatbázisrendszereknek” hívunk. Egy adatbázisrendszer hathatós eszköz arra, hogy hatékonyan készíthessünk, kezelhessünk nagy mennyiségű adatot, és lehetővé teszi azt is, hogy ezeket az adatokat hosszú ideig biztonságosan megőrízhessük. Ezek a rendszerek a jelenleg rendelkezésre álló legösszetettebb, legbonyolultabb programtermékek közé sorolhatók.

Nézünk meg, hogy milyen lehetőségeket nyújt egy adatbázisrendszer a felhasználónak:

1. *Maradandó tárolás* (persistent storage). Hasonlóan a fájlrendszerhez, egy adatbázisrendszer is támogatja a nagyon nagy mennyiségű adatok tárolását, és ez a tárolás nem csak az adatokat felhasználó folyamatok futása alatt tart, hanem ezektől függetlenül is létezik. Az adatbázisrendszer azonban tovább megy a fájlrendszerrel abban a tekintetben, hogy olyan adatszerkezetekről is gondoskodik, amelyek segítségével ez a nagyon sok adat hatékonyan, gyorsan érhető el.
2. *Programozási felület* (programming interface). Az adatbázisrendszer megengedi a felhasználónak, hogy az adatokat egy olyan lekérdezőnyelven keresztül érje el, illetve módosíthassa, amelynek elég nagy a kifejezőereje. Az adatbázisrendszer előnye egy fájlrendszerrel szemben itt is a bővebb lehetőségekből adódik, ugyanis a tártolt adatok a fájlírás/olvasásnál sokkal összetettebb módon is kezelhetők.
3. *Tranzakciókezelés* (transaction management). Az adatbázisrendszer támogatja az adatok konkurens elérését, vagyis azt, hogy több különböző folyamat (tranzakció)

A legfontosabb szakkifejezések áttekintése

Ez a könyv azok számára készült, akik felhasználói (például SQL-programozás) szemszögből már ismerik az adatbázisrendszereket legalább az Ullman–Widom: *Adatbázisrendszerek. Alapvetés* (Panem Könyvkiadó Kft., Budapest, 1998.) könyv szintjén. A következő szakkifejezéseket ezért ismerteknek tételezzük fel.

- *Adat*: tetszőleges információ, amelyet érdemes valamilyen (leginkább elektronikus) formában megőrizni.
- *Adatbázis*: hosszú ideig megőrzendő adathalmaz, mely a hozzáféréshez, módosításhoz szükséges szervezettséggel is rendelkezik.
- *Lekérdezés*: olyan művelet, mely meghatározott adatokat gyűjt ki az adatbázisból.
- *Reláció*: az adatoknak olyan kétdimenziójú táblába történő szervezése, ahol a sorok (relációsorok) valamilyen alaptényeket vagy alapegyedeket jelképeznek, és az oszlopok (attribútumok) pedig ezeknek az egyedeknek a tulajdonságait képviselik.
- *Séma*: az adatbázisban szereplő adatok szerkezetének leírása, melyet gyakran „metaadatoknak” is hívunk.

egyszerre tudja elérni az adatokat. Ahhoz, hogy az egyidejű hozzáférés nem kívánatos következményeit elkerüljük, az adatbázisrendszer támogatja az *elkülönítést* (isolation), az *atomosságot* (atomicity) és a *helyreállíthatóságot* (resiliency). Az *elkülönítés* azt jelenti, hogy a tranzakciók között látszólag csak egyet hajt végre egy időben a rendszer. Az atomosság azt a követelményt takarja, hogy egy tranzakciót vagy teljesen végrehajtsunk, vagy egyáltalán nem hajtsunk végre.

Végül a helyreállíthatóságon azt értjük, hogy sokféle rendszerhiba vagy más hiba esetén is legyen meg a lehetőség a rendszer helyreállítására.

1.1. Bevezetés: a Megatron 2000 adatbázisrendszer

Ha valaki már használt adatbázisrendszert, például egy olyan, amely támogatja a megszokott SQL lekérdezőnyelvet, akkor lehet, hogy azt képpzeli, hogy egy ilyen rendszer megvalósítása² (implementálása) nem is olyan nehéz. Képzeljünk el például egy kita-

¹ A relációsor vagy másképpen *n-es* (tuple) egy *n* komponensű vektort jelent. A *fordító megjegyzése*.

² Az angol *implementation* – megvalósítás szó helyett a magyarban leggyakrabban az *implementálás* szót használják. A *fordító megjegyzése*.

lált rendszernek az implementálást, mondjuk a Megatron Rendszernek Rt. képzeletbeli kínálatából vegyük a Megatron 2000 adatbázis-kezelő rendszert. Tegyük fel, hogy ennek a rendszernek UNIX és más operációs rendszerek alatti verziója is kapható, továbbá támogatja a relációs megközelítést és az SQL lekérdezőnyelvet.

1.1.1. A Megatron 2000 implementálásának részletei

Kezdjük azzal, hogy a Megatron 2000 a fájlrendszert használja a relációinak eltarolásához, például a Hallgatók(név, azonosító, tanszék) relációt a /usr/db/Hallgatók fájlban tárolja. A Hallgatók fájlban egy sora felel meg minden egyes relációsornak. Egy relációsor komponenseinek értékeit egymástól speciális karakterrel (például #) elválasztott karakterláncokként (string) tároljuk. Például a /usr/db/Hallgatók fájl kinézhet az alábbi módon:

```
Smith#123#IT
Johnson#522#VM
...
```

ahol IT az Informatika Tanszék, VM a Villamosmérnök Tanszék jelölje.³

Az adatbázissemánta egy /usr/db/séma nevű speciális fájlban tároljuk. A séma nevű fájlban minden relációhoz van egy olyan sora, amely a reláció nevével kezdődik, és amelyben attribútumnevek és típusok váltakozva követik egymást. Például a séma tartalmazhatja a következő sorokat:

```
Hallgatók#név#STR#azonosító#INT#tanszék#STR
Tanszék#név#STR#iroda#STR
...
```

Ezzel leírtuk a Hallgatók(név, azonosító, tanszék) relációt; a név és tanszék attribútumok típusa karakterlánc (string), míg az azonosító egész (integer) típusú. A másik sor egy Tanszék(név, iroda) sémájú relációhoz tartozik.

1.1. példa: Nézzünk meg egy példát a Megatron 2000 adatbázis-kezelő rendszer használatáról. Egy dbhost nevű gépen dolgozunk, amelyen az adatbázis-kezelő rendszert a megatron2000 UNIX-szintű parancs segítségével hívjuk meg.

```
dbhost> megatron2000
```

³ Az USA-ban az egyetemi hallgató választhat egy (esetleg több) tanszék, amely által meghirdetett, egymásra épülő tárgyakat szeretné hallgani. Így a hallgatók ehhez az egy (vagy több) tanszékhez tartoznak. Ez a tanszékhez rendelés hasonlít ahhoz, amit a magyar egyetemeken a szakok jelentenek. A fordító megjegyzi.

Erre a következő választ kapjuk:

```
ÜDVÖZÖL A MEGATRON 2000!
```

Ezután a Megatron 2000 felhasználói felületével kezdünk beszélgetést, amely számára beépíthetünk SQL-lekérdezéseket⁴, és ezzel válaszolunk a Megatron rendszer felhívására (prompt), melyet & jelöl. Egy lekérdezést a # jel zár le. Például az

```
& SELECT *
FROM Hallgatók #
```

válaszul a következő táblát adja meg:

Név	Azonosító	Tanszék
Smith	123	IT
Johnson	522	VM

A Megatron 2000 azt is megengedi, hogy végrehajtsunk egy lekérdezést, és az eredményt egy új fájlban tároljuk el. Ehhez a lekérdezést egy függőleges vonallal és a fájlnevével kell befejezni. Például az

```
& SELECT *
FROM Hallgatók
WHERE azonosító >= 500 | NagyAzonosító #
```

utasítás egy /usr/db/NagyAzonosító fájl fog készíteni, amelybe csak a következő egyetlen sor kerül:

```
Johnson#522#VM □
```

1.1.2. Hogyan hajlja végre a Megatron 2000 a lekérdezéseket?

Tekintsük az alábbi általános formájú SQL-lekérdezést:

```
SELECT * FROM R WHERE <Feltevel>
```

A Megatron 2000 erre a következőket teszi:

1. Beolvassa a séma fájl, hogy meghatározza az R reláció attribútumait és az attribútumokhoz tartozó típusokat.
2. Ellenőrzi, hogy a <Feltevel> szemantikusan érvényes-e az R relációra.

⁴ Az SQL rövid áttekintését az 1.4.2. fejezetben találjuk.

3. Minden egyes attribútumot egy-egy oszlopnak a fejléceként jelenít meg, és húz egy vonalat.
4. Beolvassa az *R* nevű fájlt, és minden egyes sora:
 - a) ellenőrzi a feltételeit, és
 - b) megjeleníti a sort, ha a feltétel igaz.

Abhoz, hogy a Megatron 2000 a

```
SELECT * FROM R WHERE <Feltétel> | T
```

utasítást végrehajtsa, a következőt fogja tenni:

1. Az előbb leírt módon végrehajlja a lekérdezést azzal a különbséggel, hogy kimarad a 3. lépés, amely az oszlopok fejléceit és a fejléceket a soroktól elválasztó vonalat generálja.
2. Kiríja az eredményt egy /usr/db/T nevű új fájlba.
3. Készít egy bejegyzést a /usr/db/séma fájlban a *T* számára, amely ugyanúgy néz ki, mint az *R*-hez tartozó bejegyzés, vagyis a *T* sémája ugyanolyan, mint az *R* sémája, azzal az eléréssel, hogy a reláció neve nem *R*, hanem *T*.

1.2. példa: Most nézzünk egy bonyolultabb lekérdezést, nevezetesen egy olyat, amelyben a Hallgatók és Tanszék miniareléciók összekapcsolására (join) van szükség:

```
SELECT iroda
FROM Hallgatók, Tanszék
WHERE Hallgatók.név = 'Smith' AND
      Hallgatók.tanszék = Tanszék.név #
```

Ez a lekérdezés azt igényli a Megatron 2000-tól, hogy kapcsolja össze a Hallgatók és Tanszék relációkat, vagyis a rendszernek egymás után vennie kell az összes sorpárt, ahol a pár egyik eleme az egyik, a másik eleme a másik relációnak sora, és meg kell határoznia, hogy vajon

- a) a sorpár sorai ugyanazt a tanszékot reprezentálják-e, és
- b) a hallgató neve Smith-e.

Az algoritmust – nem törekedve a teljes formalizálásra – a következőképpen írhatjuk le:

```
minden(Hallgatókhoz tartozó s sorra)
minden(Tanszékhez tartozó d sorra)
  ha(s és d kielégíti a WHERE-feltételt)
    jelenítse meg a Tanszékéből az iroda értékét;
```

1.1.3. Mi a baj a Megatron 2000-rel?

Lehet, hogy nem meglepő, de egy adatbázisrendszert nem szokás úgy implementálni, mint a mi elképzelt Megatron 2000 rendszerünket. Számos oka van annak, hogy az itt leírt megvalósítás alkalmatlan azokra az alkalmazásokra, amelyek jelentős mennyiségű adattal dolgoznak, vagy éppen használják az adatokat. A következőkben a teljesesség igénye nélkül felsoroljuk a legfontosabb problémákat:

- A sorok elrendezése a lemezen nem megfelelő, mivel nem nyújtja azt a rugalmasságot, amire az adatbázis módosításakor szükség lenne. Például ha kicseréljük a VM-et GAZD-ra (Gazdaságtani Tanszék) egy Hallgatókhoz tartozó sorban, akkor az egész fájl újra kell írni, mivel a VM-et követő összes karaktert két hellyel lejjebb kell mozgatni a fájlban.
- A keresés nagyon költséges. Míndig el kell olvasnunk a teljes relációt, még akkor is, ha a lekérdezés olyan értéket (vagy értékeket) használ, amely csak egy sort eredményezne, mint az 1.2. példában, ahol végig kellett néznünk a teljes Hallgatók relációt, annak ellenére, hogy egyedül a Smith nevű hallgatóra voltunk kíváncsiak.
- A lekérdezés végrehajtása úgynevezett „nyers erő” típusú, pedig az összekapcsoláshoz hasonló műveletek végrehajtására sokkal ügyesebb módszerek is léteznek. Például majd látni fogjuk, hogy az 1.2. példához hasonló lekérdezés esetén nem szükséges megnézni az összes sorpárt (ahol a pár egyik eleme az egyik, a másik eleme a másik relációnak a sora) még akkor sem, ha nem írtuk elő egy hallgató (Smith) nevét a lekérdezésben.
- Nincs mód arra, hogy a hasznos adatokat a központi memóriában eláróljuk (puffereeljük): az összes adatot a lemezeről szedjük le minden egyes esetben.
- Nincs konkurenciakézelés. Egyszerre több felhasználó is módosíthat egy fájlt, és emiatt az eredményt nem lehet előre tudni.
- Nem beszélhetünk megbízhatóságról, ugyanis adatokat veszíthetünk el, ha összeomlik a rendszer, vagy műveleteket hagyunk félbe.
- Kicsi a biztonság. Lehet, hogy az alapul szolgáló operációs rendszer valamilyen durva módon felügyeli a hozzáférést, például a különböző felhasználóknak megvan engedve, vagy meg van tiltva, hogy hozzáférjenek egy adott relációt tartalmazó fájlhoz, de nem lehet, hogy valaki mondjunk egy reláció bizonyos attribútumaihoz férhessen csak hozzá, és máshoz nem.

Ezzel a könyvvel az a szándékunk, hogy bevezessük az olvasót abba, hogyan kell ügyesebben felépíteni egy adatbázis-kezelő rendszert. Reméljük, hogy ez a tananyag mindenkinek tetszeni fog.

1.2. Egy adatbázis-kezelő rendszer áttekintése

Az 1.1. ábrán egy teljes adatbázis-kezelő rendszer vázát látjuk. Az egyvonalas dobozok a rendszer alkotórészeit jelentik, míg a dupla dobozok memóriabeli adatszerkezeteket reprezentálnak. A folytonos vonalak jelölik az olyan vezérlésátdást, ahol adatok is áramlanak, a szaggatott vonalak pedig csak az adatmozgást jelölik. Mivel az ábra bonyolult, ezért a részleteket fokozatosan tekintjük át. Először is azt javasoljuk, hogy a legfelső részen az adatbázis-kezelőhöz intézett parancsoknak két forrását különböztesük meg:

1. Szokásos felhasználói és alkalmazói programok, melyek adatokat kérnek vagy módosítják az adatokat.
2. Adatbázis-adminisztrátor (database administrator – DBA) vagy másképpen *adatbázis-rendszergazda*: egy vagy több olyan személy, akik az adatbázissémáért, illetve -struktúráért felelősek.

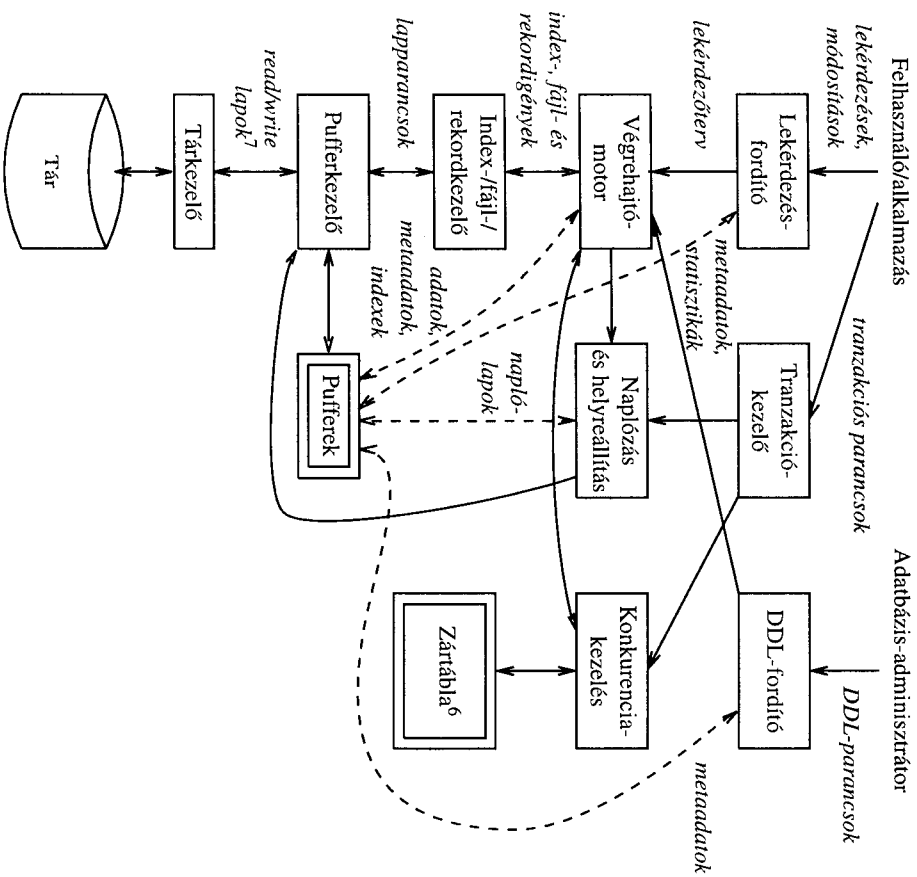
1.2.1. Az adatdefiniációs nyelv parancsai

A második fajta parancsot egyszerűbb feldolgozni. Megmutatjuk, hogyan követhető ez nyomon az 1.1. ábra jobb felső sarkából indulva. Például az adatbázis-rendszergazda elhatározza, hogy egy egyetemi nyilvántartás adatbázisában fel kellene venni egy olyan táblát vagy relációt, amelynek egyik oszlopa a hallgatói jelenti, másik oszlopa a tantárgyat, amit a hallgató felvett az indexébe, és egy harmadik oszlop pedig azt a jegyet tartalmazza, amilyen jegyet kapott a hallgató ebből a tárgyból. A DBA úgy dönt, hogy csak A, B, C, D és F jegyeket lehet megengedni.⁵ Ez az információ a struktúráról és az értékekre vonatkozó megszorításról (constraint) mind az adatbázisséma része. Az 1.1. ábrán látható, ahogy megadja ezt az adatbázis-rendszergazda, aki nek természetesen speciális felhatalmazásra van szüksége ahhoz, hogy végrehajthasson sémamódosító parancsokat, mivel ezek a parancsok alapvető hatással vannak az adatbázisra. Ezeket a sémamódosító *DDL-parancsokat* (ahol a „DDL” az adatdefiniációs nyelv angol megfelelőjének, a „data-definition language”-nek a rövidítése) a DDL-feldolgozó (DDL-processor) elemzi, majd továbbítja a végrehajtomotornak (execution engine), amely továbbadja az index-/fájl-/rekordkezelőnek, hogy az változtassa meg a metaadatokat, vagyis az adatbázis sémainformációját.

1.2.2. A lekérdezés feldolgozásának áttekintése

Az adatbázis-kezelő rendszerrel kapcsolatos kölcsönhatások döntő többsége az 1.1. ábra bal oldalán lévő útvonalat követi. A felhasználó vagy az alkalmazói program olyan működést indít el, amelynek nincs hatása az adatbázissémára, viszont hatással

⁵ Az USA-ban általában betűkkel jelölik az osztályzatokat, + jellel a féljegyet, F jelenti az elégtelent az angol *failure* – bukás szó miatt, és A vagy A+ a legjobb jegy. A *fordító megjegyzése*.



1.1. ábra. Az adatbázis-kezelő rendszer alkotórészei

lehet az adatbázis tartalmára (módosító parancs esetében), illetve adatokat gyűjthet ki az adatbázisból (lekérdezés esetében). Két olyan útvonal van, amely mentén a felhasználó cselekménye hatást gyakorol az adatbázisra:

1. A *lekérdezés megválaszolása*. A lekérdezésfordító elemzi és optimalizálja a lekérdezést. Az eredményül kapott lekérdezés-végrehajlási tervet (röviden *lekérdezéster-vel*), vagy a lekérdezés megválaszolásához szükséges tevékenységek sorozatát továbbítja a *végrehajtomotornak*. A végrehajtomotor kisebb adattarabokra (tipikusan rekordokra vagy egy reláció soraira) vonatkozó kérések sorozatát adja át az erőforrás-kezelőnek (resource manager). Az erőforrás-kezelő ismeri a relációkat tar-

⁶ A szakirodalomban igen elterjedt a lock table kifejezés is. A *fordító megjegyzése*.

⁷ Azok az adategységek, amelyeket beolvassunk (*read*), illetve kiírunk (*write*). A *fordító megjegyzése*.

talnázó *adatfájlokat*, a fájlok rekordjainak formáumát, méretét és az *indexfájlokat* is. Az indexfájlok segítenek abban, hogy az adatfájlok elemeit gyorsan meg lehessen találni. Az adatkéréseket az erőforrás-kezelő lefordítja lapokra (page), és ezeket a kéréseket továbbítja a *pufferkezelőnek* (buffer manager). A pufferkezelő szerepét az 1.2.3. részben fogjuk megvizsgálni, most röviden csak annyit, hogy a másodlagos adattárolón, általában lemezen tartjuk folyamatosan az adatokat, és a pufferkezelő feladata, hogy innen az adatot megfelelő részét hozza be a központi memória puffereibe. A pufferek és a lemez közti adatátvitel egysége általában egy lap vagy egy lemezblokk. A pufferkezelő információt cserél a tárkezelővel (storage manager), hogy megkapja az adatokat a lemeztől. Megőrtíthet, hogy a tárkezelő az operációs rendszer parancsait is igénybe veszi, de tipikusabb, hogy az adatbázis-kezelő a parancsait közvetlenül a lemezvezérlőhöz intézi.

2. *A tranzakció feldolgozása.* A lekérdezéseket és más tevékenységeket tranzakciókba csoportosíthatjuk. A tranzakciók olyan egységek, amelyeket atomosan és elkülönítve kell végrehajtani, ahogy ezt a fejezet bevezetőjében megbeszéltük. Gyakran minden egyes lekérdezés vagy módosítás önmagában is egy tranzakció. Ezenkívül a tranzakció végrehajtásának *tartósnak* (durable) kell lennie, ami azt jelenti, hogy bármelyik befejezett tranzakció hatásai még akkor is meg kell tudni őrizni, ha a rendszer összeomlik a tranzakció befejezése utáni pillanatban. A tranzakciófeldolgozói két fő részre oszthatók:

- a) *Egy konkurenciavezérlés-kezelő* vagy *ütemező* (scheduler) felelős a tranzakciók elkülönítésének és atomosságának biztosításáért.
- b) *Egy naplózás- és helyreállítási-kezelő* felelős a tranzakciók tartósságáért.

Ezeket az alkotórészeket az 1.2.4. részben fogjuk tovább vizsgálni.

1.2.3. A központi memória pufferei és a pufferkezelő

Egy adatbázis adatai normális esetben másodlagos adattárolón helyezkednek el. A mai számítógépes rendszerek esetében a másodlagos adattárolón általában mágneses lemezt kell érteni. Igen ám, de ahhoz, hogy bármilyen hasznos műveletet végezzünk az adatokon, az adatoknak a központi memóriában kell lenniük. Így aztán az adatbázis-kezelő rendszer egyik komponensének, nevezetesen a *pufferkezelőnek* a feladata, hogy a rendelkezésre álló memóriát pufferekre ossza fel. A *pufferek* azok a lap méretű területek, ahová a lemezblokkokat lehet betölteni. Ebből következik, hogy ha az adatbázis-kezelő valamelyik komponensének lemezen lévő információra van szüksége, akkor vagy közvetlenül, vagy a végrehajtomotoron keresztül kapcsolaba kell lépnie a pufferekkel és a pufferkezelővel. A különböző komponenseknek a következő információkra lehet szükségük:

1. *Adatok:* magának az adatbázisnak a tartalma.
2. *Metaadatok:* az adatbázisséma, mely leírja az adatbázis struktúráját és megszórításait.

3. *Statistikák:* az adatbázis-kezelő által az adatok tulajdonságairól (például az adatbázis különböző relációinak vagy más komponenseinek méreteiről, a bennük szereplő értékekről) összegyűjtött és tárolt információ.

4. *Indexek:* olyan adatszerkezetek, melyek támogatják az adatok hatékony elérését.

A pufferkezelő leírásának és szerepének sokkal teljesebb taglalását találhatjuk meg a 6.8. részben.

1.2.4. A tranzakció feldolgozása

Ahogy már említettük, természetes dolog, hogy egy vagy több adatbázis-műveletet egy *tranzakcióba* csoportosítsunk, mely egy olyan munkakegység, amit atomosan és más tranzakcióktól látszólag elkülönítve kell végrehajtani. Ezenfelül az adatbázis-kezelő rendszer a tartósságot is garantálja: azaz egy befejezett tranzakció munkája sosem veszhet el. Így a *tranzakciókezelő* fogadja az alkalmazás *tranzakciós parancsait*. Az alkalmazás azt is megmondja a tranzakciókezelőnek, hogy mikor kezdődnek és végződnek a tranzakciók, és még egyéb információ is ad az alkalmazás elvárásairól (például lehet, hogy nem akarja megkövetelni az atomosságot). A tranzakciófeldolgozó a következő feladatokat hajítja végre:

1. *Naplózás:* annak érdekében, hogy a tartósságot biztosítani lehessen, az adatbázis minden változását külön feljegyezzük (naplózzuk) lemezen. A *naplókezelő* (log manager) többféle eljárásmod közül választja ki azt, amelyiket követni fog. Ezek az eljárásmodok biztosítják azt, hogy teljesen mindegy, mikor történik a rendszerhiba vagy a rendszer összeomlása, a *helyreállítási-kezelő* (recovery manager) meg fogja tudni vizsgálni a változások naplóját, és ez alapján vissza tudja állítani az adatbázist valamilyen konzisztens állapotába. A naplókezelő először a pufferekbe írja a naplót, és egyeztet a pufferkezelővel, hogy a pufferek alkalmas időpillanatokban garantáltan íródjanak ki lemeze, ahol már az adatok túlélhetik a rendszer összeomlását.

2. *Konkurenciavezérlés:* a tranzakcióknak úgy kell látszódnuk, mintha egymástól függetlenül, elkülönítve végeznénk el őket. A legtöbb rendszerben igazából sok tranzakciók kell egyszerre végrehajtani. Így aztán az ütemező (konkurenciavezérlés-kezelő) feladata, hogy meghatározza az összetett tranzakciók résztevékenységeinek egy olyan sorrendjét, amely biztosítja azt, hogy ha ebben a sorrendben hajtjuk végre a tranzakciók elemi tevékenységeit, akkor az összehatás megegyezik azzal, mintha a tranzakciókat tulajdonképpen egyenként és egyszeres egészként hajtottuk volna végre. A tipikus ütemező ezt a munkát azáltal látja el, hogy az adatbázis bizonyos részére elhelyezett *zárrakat* (lock) karbantartja. Ezek a zárrak megakadályoznak két tranzakciót abban, hogy rossz kölcsönhatással használják ugyanazt az adatrészt. A zárrakat rendszerint a központi memória *lock-táblájában* (lock table) tárolja a rendszer, ahogy ez az 1.1. ábrán is látható. Az ütemező azzal befolyásolja a lekérdezések és más adatbázis-műveletek végrehajtását, hogy megtiltja a végrehajtomotorok, hogy hozzányúljon az adatbázis zár alá helyezett részéhez.

A tranzakciók ACID-tulajdonságai

A helyesen implementált tranzakciókról rendszerint azt szokás mondani, hogy eleget tesznek az ACID-tulajdonságoknak, ahol a betűk a tulajdonságok angol megfelelőinek kezdőbetűit jelölik:

- Az „A” jelöli az atomosságot (atomicity), azaz a tranzakciók „mindent vagy semmit” jellegű végrehajtását.
- Az „I” jelenti az elkülönítést (isolation), vagyis azt a tényt, hogy minden tranzakciónak látszólag úgy kell lefuthatnia, mintha ez alatt az idő alatt semmilyen másik tranzakció sem hajtanáak végre.
- A „D” a tartósságot (durability) jelöli, azaz azt a feltételt, hogyha egyszer egy tranzakció befejeződött, akkor már soha többé nem veszhet el a tranzakciónak az adatbázison kifejtett hatása.

A hiányzó „C” betű a konzisztenciát helyettesíti (consistency). Minden adatbázisnak vannak konzisztenciamegőrzési, vagy másképpen az adatellenek közti kapcsolatokra vonatkozó elvárásai. Ilyen például, hogy egy bizonyos attribútumkulcs vagy a hallgató nem vehet fel nyolcnál több tárgyat egyszerre stb. A tranzakcióktól elvárjuk, hogy megőrizzzék az adatbázis konzisztenciáját. Ezt a témát a 9.1. részben fogjuk részletesebben tanulmányozni.

3. *Holpont feloldása*: a tranzakciók az ütemező által engedélyezett záruk alapján versenyeznek az erőforrásokért. Így előfordulhat, hogy olyan helyzetbe kerülnek, amelyben egyiküket sem lehet folytatni, mert mindegyiknek szüksége lenne valamire, amit egy másik tranzakció birtokol. A tranzakciókezelő feladata, hogy ilyenkor közbeavatkozzon, és töröljön, abortáljon egy vagy több tranzakciót úgy, hogy a többi már folytathatni lehessen.

1.2.5. A lekérdezésfeldolgozó

Az adatbázis-kezelő rendszer *lekérdezésfeldolgozó* részének van a legnagyobb hatása arra, amit a felhasználó is lát a működés hatékonyságából. Az 1.1. ábrán a lekérdezésfeldolgozó két részről áll:

1. A *lekérdezésfordító* a kérdést belső formátumra fordítja le. Ezt *lekérdezőtérnek* hívjuk. Ez valószínűleg nem más, mint azoknak a műveleteknek a sorozata, amelyeket az adatokon el kell végezni. A lekérdezőtér műveletei gyakran a jól ismert relációs algebra műveleteinek implementálásai, ahogy ezt majd a 6.1. részben látni fogjuk. A lekérdezőfordító három fő részből áll:

- a) *Lekérdezéselenző*, mely a kérdés szöveges formájából egy fástruktúrát hoz létre.
- b) *Lekérdező-előfeldolgozó*, mely egyrészt tartalmilag, azaz szemantikusan ellenőrzi a kérdést (például megbizonyosodik afelől, hogy a kérdésben szereplő összes reláció létezik), másrészt a nyelvi elemző fát (parse tree) átalakítja a kiindulási lekérdezőtérrel reprezentálható, algebrai műveletekből álló fává.
- c) *Lekérdezősoprtimalizáló*, mely a kiindulási lekérdezőtérrel átalakítja az aktuális adatok alapján a lehető legjobb műveletsorozatát.

A lekérdezésfordító a metaadatok és az adatokra vonatkozó statisztikák alapján tudja nagy valószínűséggel eldönteni, hogy melyik műveletsorozat lesz a leggyorsabb. Például egy index létezése az egyik tervet sokkal gyorsabbá teheti egy másik tervhez képest.

2. A *végrehajítómotor* felelős a választott lekérdezőtér minden egyes lépésének a végrehajtásáért. A lekérdezőmotor az adatbázis-kezelő rendszer egyéb komponenseinek legfőbbjével is kapcsolatba lép. Ez történhet közvetlenül vagy a puffereken keresztül. Ahhoz, hogy kezelhessük az adatokat, előbb az adatbázisból be kell hozni őket a pufferekbe. Ehhez kapcsolatba kell lépni az ütemezővel is, nehogy zár alá helyezett adathoz akarjunk hozzáférni. Ezenkívül a naplókezelővel is kapcsolatba kell lépni, hogy az adatbázis-változások biztosan megfelelő módon legyenek naplózva.

1.3. A könyv vázlatos felépítése

Az adatbázisrendszerek implementálásának témaköre nagy vonalakban három részre osztható fel:

1. *Tárkezelés*: hogyan használjuk hatékonyan a másodlagos tárat az adatok tárolására és a gyors elérésre.
2. *Lekérdezésfeldolgozás*: hogyan lehet hatékonyan végrehajtani a nagyon magas szintű nyelven (például SQL-ben) megfogalmazott lekérdezéseket.
3. *Tranzakciókezelés*: hogyan támogathatók az 1.2.4. részben leírt, ACID-tulajdonságú tranzakciók.

A fenti témák mindenyikét a könyv több fejezetben kereszttül tárgyalja.

1.3.1. Előismeretek

A könyv ugyan nem tételez fel semmilyen előismeretet az adatbázis-kezelők implementálásáról, mégis olyan tankönyvnek szántuk, amely egy adatbázissal foglalkozó tantárgysorozat második tantárgyához vagy egy egyféléves, de átfogóbb tantárgyhoz nyújt segítséget. A könyv a Jeff Ullman és Jennifer Widom: *Adatbázisrendszerek*.

Alapvetés című könyv folytatásának is tekinthető. Ez a korábbi könyv a következőkkel foglalkozik:

1. *Adatbázis-tervezés*: az adatbázisséma közvetlen, magas szintű specifikálása, valaminten jelölésrendszert alkalmazva, ami lehet az egyed/kapcsolat modell vagy akár az ODL (Object Description Language – Objektumleíró nyelv), valamint a tervek implementálása az SQL-nyelv adatdefíniós részének segítségével.
2. *Adatbázis-programozás*: lekérdezések és adatbázis-módosító parancsok írása valaminten megfelelő nyelven, például SQL-ben.

Az adatbázis-tervezés technológiájának nem sok köze van az adatbázis-kezelő rendszerek implementálásához, de szükséges, hogy az olvasó ismerje a relációs modellt és azt, hogy az adatokat hogyan reprezentáljuk a relációk segítségével. Erre azért van szükség, mivel a könyvben elmondottak jó része azzal foglalkozik, hogyan tároljuk a relációkat, hogyan optimalizáljuk a relációkra vonatkozó lekérdezéseket, és hogyan vezéreljük a relációkhoz történő hozzáféréseket valamilyen (például zárolási) módzerekkel. Emellett csak akkor látjuk tisztán a lekérdezések feldolgozása mögött rejlő technológiát, ha az SQL-programozást is ismerjük. Ezeknek a témáknak egy gyors áttekintését adja a 1.4. rész.

Továbbá ismertek tételezzük fel a *fájlok* (azaz adatok tárolására használható neve-sített tárterületek) fogalmát. Azt is elvárjuk, hogy az olvasónak legyenek előismeretei egy hagyományos fájlrendszer felépítéséről. A fájlrendszer az operációs rendszernek az a része, amely kezeli az operációs rendszer fájljait. Egy adatbázis-kezelő rendszer egészen másképp kezeli a fájlokat, de ennek a fontos témakörnek is átvesszük majd az alapjait.

1.3.2. A tárkezelés áttekintése

A könyv tárkezelésről szóló fejezeteivel kezdődik. A 2. fejezet bevezeti a memóriahierarchiát, de a legnagyobb részletességgel majd azt fogjuk inkább vizsgálni, hogy milyen módon kell tárolni és elérni a lemezen az adatokat, ugyanis a másodlagos adatokra, különös tekintettel a lemezekre, központi jelentőségűek abban, ahogy egy adatbázis-kezelő rendszer kezeli az adatokat. A lemezalapú adatok vonatkozásában bevezetjük a blokkmodellt, amely az adatbázisrendszerben történő majdnem minden tevékenységre kihat.

A 3. fejezet kapcsolatot teremt az adatelelemek (relációk, sorok, attribútumértékek, illetve más adatmodell esetén az előbbieknak megfelelő fogalmak) tárolása és az adatok blokkmodelljére vonatkozó követelmények között. Ezután megvizsgáljuk azokat a fontos adatszerkezeteket, amelyeket indexek készítésére használunk. Már említettük, hogy az index egy olyan adatstruktúra, amely az adatok hatékony elérését támogatja.

A 4. fejezet áttekinti a fontosabb egydimenziós indexstruktúrákat, indexszekvenciális fájlokat, B-fákat és tördelőtáblákat. Ezeket az indexeket használják hagyományosan az adatbázis-kezelő rendszerek az olyan lekérdezések támogatására, ahol egy attribú-

tumértéket adunk meg, és azokat a sorokat keressük, amelyekben ez az érték szerepel. Az 5. fejezet a többdimenziós indexeket tárgyalja, amelyek olyan adatszerkezetek, amik speciális alkalmazásokhoz, például földrajzi adatbázisokhoz használhatók. Ezeknél a tipikus kérdések valamilyen tartomány, terület tartalmát kérdezik le.

A fenti indexstruktúrák az olyan összetett SQL-lekérdezéseket is támogatni tudják, amikor kettő vagy több attribútumértékre szól a korlátozás. Így nem esoda, hogy némelyikük felhívjuk a kereskedelmi forgalomban kapható adatbázis-kezelő rendszerekben is.

1.3.3. A lekérdezésfeldolgozás áttekintése

A 6. fejezet bevezeti a relációs algebrát, amivel a lekérdezések végrehajtása is leírható. Ebben a fejezetben találjuk meg a lekérdezések végrehajtásával foglalkozó alapokat, beleértve számos olyan algoritmust, amelyek a kulcsfontosságú relációs műveletek (például relációk összekapcsolása) hatékony implementálására szolgálnak.

A 7. fejezetben áttekinjük a lekérdezésfordítónak és az optimalizálónak a felépítését. A vizsgálatot a lekérdezések elemzésével és a tartalmi szemantikus ellenőrzésükkel kezdjük. A következő részben megvizsgáljuk, hogyan alakítható át egy lekérdezés SQL-ből relációs algebrába. Ezután a *logikai lekérdezőterv* kiválasztásával foglalkozunk. A logikai lekérdezőterv egy algebrai kifejezés, amely az adatokon elvégzendő speciális műveleteket reprezentálja a műveletek sorrendjére vonatkozó szükséges megszorításokkal együtt. Végül a *fizikai lekérdezőterv* kiválasztását nézzük meg. Ez utóbbi magában foglalja a műveletek speciális sorrendjének megadását és az egyes műveleteket implementáló algoritmust is.

1.3.4. A tranzakciófeldolgozó áttekintése

A 8. fejezetben látni fogjuk, miképp támogatja egy adatbázis-kezelő rendszer a tranzakciók tartósságát. Az alapötlet, hogy készítsünk egy olyan naplót, amibe az adatbázis minden változtatása bekekerül. Tudjuk, hogy bármi elveszhet, ami a központi memóriában van és nem lemezen egy összeomlás (például áramkimaradás) esetén, ezért körültekintőnek kell lennünk, hogy az adatbázis változását és a változást rögzítő naplót alkalmass sorrendben vigyük át a puffertől lemeze. Sok különböző naplózási módszer létezik, de mindegyik valamennyire korlátozza a további tevékenységünk szabadságát.

Ezután a 9. fejezetben rátérünk a konkurenciakézelésre, amely biztosítja az atomosságot és az elkülönítést. A tranzakciókat adatbáziselemek olvasási és írási műveleteiből álló sorozatnak tekintjük. Ebben a fejezetben az a fő kérdés, hogyan kezeljük az adatbáziselemekre elhelyezett záratokat, milyen különböző típusú zárok használhatók, hogyan lehet a tranzakcióknak megengedni, hogy záratokat töltsenek az elemekre, vagy elengedjék, feloldják a záratokat. Emellett azt is tanulmányozni fogjuk, hogyan biztosítható az atomosság és elkülönítés akkor, ha nem használunk záratokat.

A 10. fejezet összefoglalja a tranzakciófeldolgozásról tanultakat. Áttekinjtjük azt,

hogy hogyan egyeztethetők össze a naplózásra és a konkurenciára vonatkozó elvárásaink. Az ezekre vonatkozó követelményeket a 8. és 9. fejezetben fektettük le. A tranzakciókezelő másik fontos feladatát, a holtpontkezelést is itt nézzük meg alapsabban. A 10. fejezetben találjuk meg a konkurenciavezérlés kiterjesztését osztott környezetre. Végül bevezetjük annak a lehetőségét, hogy a tranzakciók hosszúak is lehetnek, azaz nem a másodperc ezredéséig tartanak, hanem órákig vagy akár napokig is. Ha egy hosszú tranzakció zárolja az adatokat, akkor ezzel nagy káoszt idézhet elő azok között a lehetséges felhasználók között, akik éppen ezeket az adatokat használják. Ez arra késztet bennünket, hogy újragondoljuk a konkurenciavezérlési azokra az alkalmazásokra, amelyekben hosszú tranzakciók is szerepelhetnek.

1.3.5. Az információintegráció áttekintése

Az adatbázisrendszerek fejlődését az utóbbi években nagyban meghatározta az a törekvés, hogy megereimdjön a lehetőség arra, hogy különböző *adatforrások* (adatbázisok és/vagy nem adatbázisrendszerrel kezelt információforrások) úgy működjenek együtt, mintha egy nagy egésznek a részei lennének. A 11. fejezetet annak szenteljük, hogy áttekintjük ennek az *információintegrációnak* nevezett új technológiának a fontosabb szempontjait. Tártyaljuk az integráció alapvető módszereit, amelyek sorában foglalkozunk az adatforrások lefordított és integrált másolataival, amit *adatraktárnak* vagy *adatirrháznak* (data warehouse) hívunk, és az adatforrások halmozásának virtuális nézetével, amit *közvetítőnek* (mediator) nevezünk.

1.4. Az adatmodellek és nyelvek áttekintése

Ebben a részben röviden áttekintjük az SQL-t és a relációs modellt. Ezenkívül ismeretjük az objektumok fogalmát, ahogy azt az objektumorientált adatbázisokban használjuk. A példákat az Ullman–Widom: *Adatbázisrendszerek. Alapvetés* című könyvből vesszük át.

1.4.1. A relációs modell áttekintése

A *reláció soroknak* a halmaza, a sorok pedig értékekből álló listák. Egy reláció minden sorának ugyanannyi számú komponense van, és a különböző sorokból vett megfelelő komponensek ugyanolyan típusúak. Egy relációt úgy jeleníthetünk meg, hogy minden egyes sorát felsoroljuk egy táblázat soraként. Az oszlopok fejlécét *attribútumnak* hívjuk. Az attribútumok a sorok komponenseinek értelmét jelentik. A reláció neve, az attribútumok nevei és az attribútumokhoz tartozó típusok együtt alkotják a reláció *sémáját*.

1.3. példa: A példákban gyakran fogjuk idézni a Film relációt, mely tartalmazhatná a következő sorokat:

Filmcím	Év	Hossz
Csillagok háborúja	1977	124
Erős kacsák	1991	104
Wayne világga	1992	95

Ennek a relációnak a sémája az alábbi:

Film(filmcím, év, hossz)

Az attribútumai a filmcím, az év és a hossz, melyekről feltehetjük, hogy rendre karakterlánc, egész, egész típusúak. A vonal alatti három sor mindegyike egy-egy relációsor. Az első sor például azt mondja, hogy a „Csillagok háborúját” 1977-ben készítették és 124 perc hosszú. □

Az *adatbáziséma* relációsémáknak a halmaza. A filmekkel kapcsolatos állandó példánkban gyakran fogjuk használni a következő relációkat:

Film(filmcím, év, hossz, stúdiónév⁸)

Filmszínész(név, cím, neme, születési_idő)

Szerepel(filmcím, év, színésznev)

Stúdió(név, cím)

Az első reláció majdnem ugyanaz, mint az 1.3. példában szereplő Film reláció, azaz a különbséggel, hogy a sémához még hozzávetettük a filmet gyártó stúdió nevét is, azért hogy szükség esetén további kapcsolatokat tudjunk majd a példánkban létrehozni. A második reláció a filmszínészektől nyújt információt, míg a harmadik összekapcsolja a filmeket a szereplőikkel. A negyedik reláció pedig a stúdióktól ad információt. A különböző attribútumok jelentése az attribútumok nevéből kézenfekvően következik.

1.4.2. Az SQL áttekintése

Az SQL nevű adatbázisnyelv nagyon sok lehetőséggel rendelkezik. Ezek között megtalálhatók azok az utasítások, amelyek az adatbázist lekérdezik, illetve módosítják. Az

⁸ Az attribútum neve általában egy szó, ezért angolul a stúdioName jelölést használják, ami utal arra, hogy angolban a stúdiónév eredetileg két szó (*studio name*) lenne. A fordításban is ezt a konvenciót fogjuk használni, azaz két szóból szükség esetén úgy alkotunk egyet, hogy nagybetűvel kezdjük a második tagot. Mivel a magyar helyesírás szerint a stúdiónév eleve egybetrandó, ezért itt nem kell használnunk a megkülönböztető nagybetűt. A *fordító megjegyzése*.

adatbázis módosítása három paranccsal (INSERT, DELETE és UPDATE) keresszük törté-
nik, melyek formális megadását, szintaxisát most itt nem adjuk meg. A lekérdezéseket
általában „select-from-where” utasításokkal fejezzük ki, melyek általános formáját
ténylegesen az 1.2. ábra mutatja. Az utasításnak csak a SELECT és FROM szavakkal
kezdődő első két sorát, vagy másképpen *záradékait*⁹ (clause) kötelező megadni.

```
SELECT <attribútumok listája>  
FROM <relációk listája>  
WHERE <feltétel>  
GROUP BY <attribútumok listája>  
HAVING <feltétel>  
ORDER BY <attribútumok listája>
```

1.2. ábra. Egy SQL-lekérdezés általános formája

Egy ilyen lekérdezés eredményét, még ha nem is a lehető legjobb módon, de ki-
számolhatjuk az alábbiak szerint:

1. Vegyük a FROM záradékban szereplő relációkból a sorok összes lehetséges kombi-
nációját.
2. Dobjuk el azokat a kombinációkat, amelyek nem elégítik ki a WHERE záradék fel-
tételét.
3. A megmaradt sorkombinációkat csoportosítjuk a GROUP BY záradékban felsorolt
attribútumokhoz (ha van ilyen egyáltalán) tartozó értékeik alapján.
4. Ellenőrizzük az összes csoportra a HAVING záradékban szereplő feltételt (ha egy-
általán megadtunk ilyen), és hagyjuk el az összes olyan csoportot, amely nem felel
meg ennek a feltételnek.
5. Számítsuk ki a sorokat a SELECT záradékban megadott attribútumokból és attribú-
tumok összesítéséből¹⁰ (aggregation) (például a csoportokon belüli összegek
képzéséből).
6. Rendezzük az eredményül kapott sorokat az ORDER BY záradékban megadott
attribútumlistának megfelelő értékeik alapján.

1.4. példa: Az 1.3. ábra egy olyan egyszerű SQL-lekérdezést mutat be, amelynek csak
az első három záradéka van megadva. Ez a lekérdezés a Paramount stúdió által gyár-
tott filmek címeit és a bennük szereplő színészek nevét adja meg. Megjegyezzük,
hogy a filmcím és év együtt a Film reláció kulcsa, mivel két filmnek lehetne ugyan
meg egyező címe, de reményeink szerint ekkor viszont nem ugyanabban az évben ke-
szültek.

⁹ Szokás még mondatrészeknek, klauzuláknak vagy klónnak is hívni. A *fordító megjegyzése*.

¹⁰ Ezeket a statisztikai, összesítő függvényeket aggregátoroknak is szokták fordítani. A *for-
dító megjegyzése*.

```
SELECT színészNév, Film.filmcím  
FROM Film, Szerepel  
WHERE Film.filmcím = Szerepel.filmcím AND  
Film.év = Szerepel.év AND  
stúdióNév = 'Paramount';
```

1.3. ábra. A Paramount színészeinek megkeresése

1.5. példa: Az 1.4. ábra egy bonyolultabb lekérdezést mutat be. Először is azt kell
megkeresnünk, hogy kik azok a színészek, akik legalább három filmben szerepeltek.
A lekérdezésnek ezt a részét úgy tudjuk megvalósítani, hogy a Szerepel sorait a
GROUP BY záradékkal a színészek neve szerint csoportosítjuk, és aztán a HAVING zá-
radékkal kiszűrjük azokat a csoportokat, amelyeknek kettő vagy kevesebb sora van.

```
SELECT színészNév, MIN(év) AS először  
FROM Szerepel  
GROUP BY színészNév  
HAVING COUNT(*) >= 3  
ORDER BY először;
```

1.4. ábra. Azoknak a legkorábbi éveknak a megkeresése, amikor a legalább három filmben játszó színészek először szerepeltek filmben

Ezután a SELECT záradék azt mondja, hogy a megmaradó csoportokból elő kell
állítani a színészek nevét és a legkorábbi évet, amikor a színész először szerepelt egy
filmben. A select-lista második tagját, vagyis a MIN(év)-et először-re nevezzük át.
Végül az ORDER BY azt mondja, hogy az eredményül kapott sorokat az először ér-
tékei szerint növekvő sorrendben kell listázni, vagyis a színészek az első filmjük évé-
nek sorrendjében fognak következni.

Alkérdezések

Az SQL egyik legerőteljesebb sajátossága az, hogy a WHERE, FROM vagy HAVING zá-
radékokon belül lehetőségünk van alkérdezések (subquery) használatára. Az alkérde-
s egy olyan „select-from-where” utasítás, amelynek az értékét a fent említett záradékok
ellenőrzik.

1.6. példa: Az 1.5. ábra egy olyan SQL-lekérdezést mutat, amely egy alkérde-
ssel is rendelkezik.

```
SELECT filmcím, év  
FROM Film  
WHERE stúdióNév IN (  
SELECT név  
FROM Stúdió  
WHERE cím NOT LIKE '%Hollywood%'  
);
```

1.5. ábra. A nem Hollywoodban készült filmek megkeresése

A teljes lekérdezés a nem Hollywoodban gyártott filmek címét és gyártási évét keresi meg, míg a

```
SELECT név
FROM Stúdió
WHERE cím NOT LIKE '%Hollywood%'
```

alkérdés azt az egyoszlopos relációt adja vissza, amely azon stúdiók nevéből áll, melyek címben nem fordul elő a „Hollywood” szó. Ezután ezt az alkérdést használja a külső lekérdezés WHERE záradéka arra, hogy beazonosítsa azokat a filmeket, amelyek stúdiója nem szerepel az alkérdés által meghatározott stúdiónevek halmazában. □

Nézetablák

Az SQL egy másik fontos lehetősége, hogy *nézetablákat* vagy röviden *nézeteket* (view) lehet definiálni. A nézetablák valójában relációk leírásai. Ezeket a relációkat nem tároljuk, de szükség esetén elő tudjuk állítani a tárolt relációkból.

1.7. példa. Az 1.6. ábra egy nézetábla definícióját mutatja. Ez a nézetábla a Paramount stúdiók által készített filmek címét és gyártási évét határozza meg. A ParamountFilm nézetábla definícióját az adatbázisséma részeként tárolja a rendszer, de a hozzá tartozó sorokat most még nem számolja ki. Csak akkor állítja elő a sorait, amikor egy lekérdezés a ParamountFilm relációt használja. Ha ennek a lekérdezésnek nincs szüksége a teljes relációra, akkor csak a sorának egy szükséges részahalmazát állítja elő megfelelőképpen. Mindezt azzal éri el, hogy a nézetábla definícióját beépíti a lekérdezésbe. A nézetábla sorait így valójában sohasem tároljuk az adatbázisban. □

```
CREATE VIEW ParamountFilm AS
SELECT filmcím, év
FROM Film
WHERE stúdiónev = 'Paramount';
```

1.6. ábra. A Paramount filmjeit meghatározó nézetábla

1.4.3. A relációs és objektumorientált adatok

A könyvben tárgyaltak többsége azt tételezi fel, hogy az adatbázis relációs adatbázis: vagyis az adatokat táblákkal modellezzük, az adatelemeket relációsorokként vagy a tábla soraként. A soroknak rögzített számú komponense van, és ezek mindegyike a relációsémában rögzített típusúval rendelkezik. Az adatoknak ezt a szemléltetését sugallta az 1.3. példa. Egy ettől különböző szinten úgy is gondolhatunk egy sorra, mint egy struktúrára (C nyelven „struct”-ra) vagy egy olyan rekordra, amelynek minden mezője egy attribútumértéknek felel meg.

Egyes adatbázisrendszerekben másmilyen adarmodelli használnak, nevezetesen az adatokat objektumoknak is lehet tekinteni. Ebben a modellben az elemi adatelem az *objektum*. Az objektumokat *osztályokba* csoportosítjuk, és minden osztálynak van egy sémája, ami az osztály *jellemzőinek, tulajdonságainak* listája.

1. A jellemzők között lehetnek attribútumok, melyeket a relációsorok attribútumaihoz hasonlóan lehet feltüntetni.
2. A *kapcsolatok* (relationship) is jellemzők. Ezek kötnék össze egy objektumot egy vagy több másik objektummal. Az implementációs szinten úgy gondolhatunk egy kapcsolatra, mint az ezekre az objektumokra mutató mutatók (pointerek) listájára.
3. Vannak még *metódusoknak* (módszerek, eljárások) hívott jellemzők is, melyek tulajdonképpen olyan függvények, amiket az osztály objektumaira lehet alkalmazni.

Aztól eltekintve, hogy a metódusok kódja tipikusan az objektumokon kívül lesz tárolva, az adatok objektumorientált formalizálásának többi része jól illeszkedik az általános keretünkbe. Ugyanis általában úgy gondolunk a fájlokra, mint a legnagyobb adategységekre. A fájlokat tekinthetjük egyszerűen névvel ellátott adatgyűjteményeknek. A fájlok rendszerint kisebb egységekből állnak, amelyekre a következő elnevezéseket használjuk:

- a) A legelső adatbázisokban a fájlok *rekordokból* álltak, a rekordok pedig *mezőkből*. A rekord a C nyelv és a lezármazott programozási nyelvek (C++, Java) „struct”-jával rokon fogalom.
- b) A relációs adatbázisban a fájlok *relációk*, melyek relációsorokból állnak. A relációsorokat pedig *attribútumok* alkotják.
- c) Az objektumorientált adatbázisokban a fájlok az *osztályok előfordulásai* (extent). Egy ilyen osztály-előfordulás az osztályhoz adott pillanatban tartozó objektumok halmazát jelenti. Az osztály-előfordulások *objektumokból* állnak, az objektumoknak pedig *mezői* vannak. A mezőket másképpen példányváltozóknak (instance variable) is nevezik, melyek értékei jelenthetik az objektum attribútumait vagy jelenthetik a kapcsolódó objektumok halmazát is valamilyen kapcsolaton keresztül.

Hasznos lehet, ha összegezzük a következő hasonlóságokat:

1. A fájl, a reláció és az osztály-előfordulás hasonló fogalmak. Mind olyan értéket jelentenek, melyek valamilyen kisebb, de közös sémával rendelkező elemekből (rekordok, relációsorok vagy objektumok) állnak.
2. Egy fájl vagy reláció sémája és egy osztály definíciója szintén hasonló fogalmat takar. Mind azt írja le, hogy milyenek egy fájlnak, relációnak vagy egy osztály-előfordulásnak az elemei.
3. A rekordok, relációsorok és objektumok is hasonló fogalmak. Mindegyiküket gyakran implementálhatjuk úgy, mint ha „struct”-ok lennének a C nyelvben.

1.5. Összefoglalás

- *Adatbázis-kezelő rendszerek:* Ezek a rendszerek azzal a képességgel jellemezhetők, hogy támogatják a nagyon nagy mennyiségű adatok hatékony elérését, és ezek az adatok hosszú ideig megőrződnek. További jellemvonásuk még, hogy támogatják a nagy kifejezőerővel rendelkező lekérdezőnyelveket. Támogatják a tartós tranzakciók konkurens végrehajtását is oly módon, hogy a tranzakciók atomosnak és más tranzakcióktól függetlennek lássanak.
- *Összehasonlítás a fájlrendszerekkel:* A hagyományos fájlrendszer alkalmatlan adatbázisrendszernek, mert nem támogatja a hatékony keresést, a kisebb adatok hatékony módosítását, az összetett lekérdezéseket, a hasznos adatok vezérelt puffelését a központi memóriában, és nem támogatja a tranzakciók atomos és független végrehajtását sem.
- *Az adatbázis-kezelő részeli:* Az adatbázis-kezelő rendszer fő részei a tárkezelő, a lekérdezőfeldolgozó és a tranzakciókezelő.
- *A tárkezelő:* Ez a komponens felelős az adatok, metaadatok (adatszerkezeteket, sémát leíró információk), indexek (adatok gyors elérését biztosító adatszerkezetek) és naplók (az adatbázis változásairól szóló feljegyzések) lemezen történő tárolásáért. A tárkezelő fontos eleme a puffkezelő, mely a lemez tartalmának egy részét a központi memóriában tartja.
- *A lekérdezőfeldolgozó:* Ez a komponens elemzi a lekérdezéseket, egy lekérdezőterv kiválasztásával optimalizálja őket, és végrehajtja a tervet a tárolt adatokon.
- *A tranzakciókezelő:* Ez a komponens felelős az adatbázis változásainak naplózásáért. Ezáltal támogatja, hogy egy rendszerösszeomlás után helyre lehessen állítani a rendszert. Emellett felelős a tranzakciók konkurens végrehajtásáért oly módon, amely biztosítja az atomosságot (egy tranzakciót vagy teljesen végrehajtunk, vagy egyáltalán nem hajtunk végre) és az elklüönítést (a tranzakciók végrehajtása olyan, mintha más konkurens tranzakciót nem hajtanánk végre).
- *SQL:* Ez egy fontos, szabványos, relációs modellen alapuló lekérdezőnyelv. Mind a nyelv, mind a relációs modell központi jelentőségű a könyvünk nagy részében.
- *Adatfogalmak:* A fájlrendszereknek, a C-hez hasonló hagyományos programozási nyelveknek, a relációs modellnek és az objektumorientált adatmodellnek sok közös fogalma van, bár ezekre gyakran különböző szakkifejezéseket használnak. Párhuzamot lehet vonni a „struct”-ok, relációsorok és objektumok vagy a fájlok, relációk és oszlopok között.

1.6. Irodalomjegyzék

A ma már közvetlenül (on-line módon) kereshető bibliográfiákban megtalálható majd-nem minden aktuális adatbázisrendszerrel foglalkozó cikk. Emiatt könyvünkben nem szándékozunk teljes irodalomjegyzéket adni, inkább csak a történelmi fontosságú cikkeket, a hasznos áttekinthető tanulmányokat, és azokat a másodlagos leltőhelyeket ad-

juk meg, ahol további cikkeket lehet találni. Michael Ley [6] elkészítette az adatbázis-kutatással foglalkozó cikkeknek egy keresésre alkalmas tárgymutatóját. Lehet keresni Alf-Christian Achilles könyvtárában is, aki folyamatosan karbantartja az adatbázis témával kapcsolatos lényeges indexeket [1].

Az ehhez a könyvhöz szükséges háttérismereket [8]-ból lehet elsajátítani. Az SQL2 és SQL3 szabványok letölthetők az [5] anonim FTP-szerverről. Akik egy SQL2 kézikönyvet szeretnének, azoknak [4]-et ajánljuk.

A témakör technológiájához sok prototípus adatbázisrendszer implementálása járult hozzá, ezek közül a két legismertebb az IBM Almaden Kutatási Központjának System R rendszere [2], és az INGRES projekt, amit Berkeleyn fejlesztettek [7]. Mindkettő olyan korai relációs rendszer, melynek köszönhetően a relációs rendszer lett a meghatározó adatbázis-technológia.

Az 1998-as „Asilomar report” [3] a legfrissebb az adatbázisrendszerek kutatásáról és irányvonalairól szóló jelentések sorában. Ebben is találunk hivatkozásokat korábbi hozzá hasonló jelentésekre.

1. <http://www.ira.uka.de/db1biography/Database>.
2. M. M. Astrahan et al., „System R: a relational approach to database management”, *ACM Trans. on Database Systems* 1:2 (1976), pp. 97–137.
3. P. A. Bernstein et al., „The Asilomar report on database research”, http://s2k-ftp.cs.berkeley.edu:8000/postgres/papers/Asilomar_Final.htm.
4. Date, C. J. and H. Darwen, *A Guide to the SQL Standard*, Fourth Edition, Addison-Wesley, Reading, MA, 1997.
5. <ftp://jerry.ece.umassd.edu/isow93>.
6. <http://www.informatik.uni-trier.de/~ley/db/index.html>. Egy másolat található a következő helyen: <http://www.acm.org/stgmod/db1p/db/index.html>.
7. M. Stonebraker, E. Wong, P. Kreps, and G. Held, „The design and implementation of INGRES”, *ACM Trans. on Database Systems* 1:3 (1976), pp. 189–222.
8. J. D. Ullman and J. Widom, *A First Course in Database Systems*, Prentice-Hall, Englewood Cliffs NJ, 1997. (Magyarul Ullman-Widom: *Adatbázisrendszerek Alapvetés*, Panem Könyvkiadó Kft., Budapest, 1998.)

Adattárolás

Az egyik legfontosabb különbség az adatbázis-kezelő rendszerek és más rendszerek között az, hogy az adatbázis-kezelő rendszerek nagyon sok adatot is hatékonyan tudnak kezelni. Ebben és a következő fejezetben megismerjük azokat az alapvető technikákat, amelyek arra vonatkoznak, hogyan kell kezelni az adatokat a számítógépből. A tanulmányozásunk két részre osztható:

1. Hogyan tárolja és kezeli egy számítógépes rendszer a nagyon nagy mennyiségű adatokat?
2. Milyen reprezentációk és adatszerkezetek támogatják a legjobban ezeknek az adatoknak kezelését?

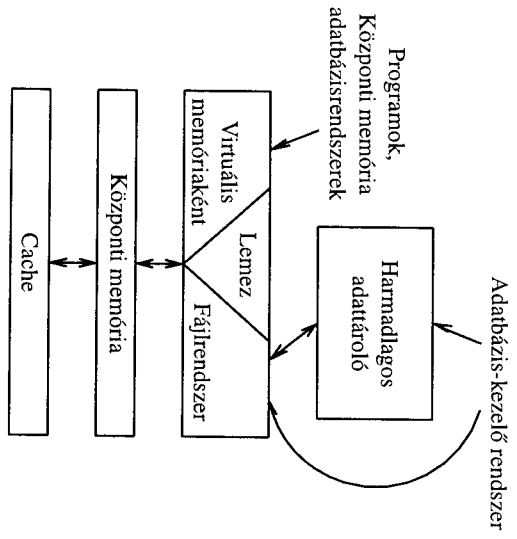
Az első kérdést ez a fejezet fedi le, míg a második a következő három fejezet témája lesz.

Ez a fejezet azzal kezdődik, hogy milyen technika használható a nagy mennyiségű adatok fizikai szintű tárolására. Megvizsgáljuk az információ tárolására használható eszközöket, elsősorban a forgó mozgást végző lemezeket. Bevezetjük a „memóriahierarchiát”, és megnézzük, hogyan függ a nagyon nagy tömegű adatokra vonatkozó algoritmusok hatékonysága a különböző adatmozgatási sémáktól, ahol az adatok mozgása a központi és a másodlagos adattároló eszközök (tipikusan lemezek) vagy esetleg „harmadlagos adattároló eszközök” között történik. Ez utóbbiak olyan robot-eszközök, melyek nagyszámú optikai lemezt vagy mágneses szalagot tároló karretta (tape cartridge) tárolására és elérésére szolgálnak. A kétfázisú, többszörös, összefüggő rendezés nevű speciális algoritmust használjuk arra, hogy például mutassunk olyan algoritmusra, amely hatékonyan használja a memóriahierarchiát.

A 2.4. részben számos technikát fogunk nézni arra is, hogyan lehet csökkenteni a lemeztől történő adatolvasáshoz, illetve lemezeztől történő adatíráshoz szükséges időt. Az utolsó két rész olyan módszereket vizsgál meg, amelyekkel javítani lehet a lemezek megbízhatóságán. A felvetett problémák tartalmazzzák az időszakos írási, olvasási hibákat és a lemez összeomlását is, vagyis amikor az adatok tartósan olvashatlanná válnak.

2.1. A memóriahierarchia

Egy tipikus számítógépes rendszernek több különböző olyan része is van, amelyekben adatokat lehet tárolni. Ezeknek a részeknek az adatkapacitása akár hét nagyságrenddel is elérhet egymástól, és a bennük tárolt adatok elérési sebessége is hét vagy még több nagyságrendben különbözhet. Az egy bájtra vonatkoztatott költség is különbözik ezeknél a komponenseknél, bár itt már sokkal kisebb az elérés, talán három nagyságrend van a legolcsóbb és a legdrágább táruk ára között. Nem meglepő, hogy a legkisebb kapacitású eszközök kínálják a leggyorsabb elérési sebességet a legdrágább bájtra vetített költségért. A memóriahierarchia egy lehetséges sémája a 2.1. ábrán látható.



2.1. ábra. A memóriahierarchia

2.1.1. Cache¹

A hierarchia legalsó szintjén találjuk a *cache*-t. A *cache* egy integrált áramkör („chip”) vagy egy processzor chipjének a része. Alkalmas arra, hogy adatokat vagy gépi utasításokat tároljon. A *cache*-ben tárolt adatok, utasítások a központi memória (ami a memóriahierarchia következő szintje) egy részének, illetve bizonyos helyeinek a másolata. Időnként a *cache*-ben lévő értékek is megváltoznak, de a központi memória megfelelő módosítása általában későbbre halasztódik. Ennek ellenére a *cache* minden értéke bármelyik időpillanatban megfelel egy helynek a központi memóriában. A központi memória és a *cache* közti adatátvitel egysége általában csak néhány bájti. Ebből kifolyólag azt gondolhatjuk a *cache*-ről, hogy egyedi gépi utasításokat, egész vagy lebegő pontos számokat, vagy rövid karakterstringokat tárol.

¹ A *cache* szó rejtett memóriát jelent, de a magyar szakirodalomban az angol szó terjedt el, bár a gyorsítótár kifejezés is használatos rá. A fordító megjegyzése.

A gép cache-je gyakran két szintre osztható. A *beépített cache* (on-board cache) ugyanazon a chipen található, mint maga a mikroprocesszor, míg a *második szintű cache* (level-2 cache) egy másik chipen helyezkedik el.

Amikor a gép utasításokat hajt végre, akkor mind az utasításokat, mind az utasítások által használt adatokat megpróbálja megkeresni a cache-ben. Ha nem találja ott őket, akkor a központi memóriából a cache-be másolja az utasításokat vagy az adatokat. Mivel a cache csak kevés adatot tud tárolni, ezért a cache-ből rendszert el kell távolítanunk valamit, azért, hogy helyet biztosítsunk benne az új adatoknak. Ha az, amit el akarunk távolítani a cache-ből, nem változott meg, mióta a cache-be másoltuk, akkor semmi más teendők nincs, mint hogy kidobjuk a cache-ből. Ezzel szemben, ha a cache-ből kidobásra szánt értékek módosultak, akkor az új értéket a központi memóriára megfelelő helyére kell bemásolni.

Az egyszerű, egyprocesszoros számítógépnek nem szükséges aktualizálni a központi memória megfelelő helyét még abban az esetben sem, ha a cache-ben módosultak az adatok. Ezzel szemben az olyan multiprocesszoros rendszerben, ahol több processzor is hozzáférhet ugyanahhoz a központi memóriához és mindegyik processzor saját cache-t tart fenn, gyakran szükséges, hogy a cache módosulását azonnal *áthetessék*, azaz azonnal módosítsák a központi memória megfelelő helyét.

Az ezredforduló idején használatos cache-ek legfeljebb egy megabájt kapacitásúak. A cache és a processzor közt az adatokat a processzor művelési sebességével lehet olvasni vagy írni, ami rendszert 10 nanoszekundumot (10^{-8} másodperc) vagy kevesebbet jelent. Másrészt a cache és a központi memória közti adatellenek vagy utasítások mozgata sokkal tovább tart, körülbelül 100 nanoszekundumot (10^{-7} másodperc) vesz igénybe.

2.1.2. A központi memória

A tevékenységek középpontjában a számítógép *központi memóriája* áll. Mindenne, ami a számítógépben történik – utasítások végrehajtása, adatok kezelése – úgy gondolhatunk, mintha a központi memóriában jelen levő információval dolgoznánk (bár a gyakorlatban az is normális, ahogy a 2.1.1. részben tárgyalt módon a cache-be töljtjük át az adatok, utasítások egy részét).

1999-ben egy tipikus számítógép körülbelül 100 megabájt (10^8 bájt) központi memóriával rendelkezett, bár lehetett kapni sokkal nagyobb, 10 vagy még több gigabájt (10^{10} bájt) központi memóriájú gépeket is.

A központi memóriák *véletlen hozzáférések* (random access), ami azt jelenti, hogy bármelyik bájtot ugyanannyi idő alatt lehet megkapni.² A központi memória adatainak tipikus elérési ideje 10 és 100 nanoszekundum között változik (azaz 10^{-8} és 10^{-7} másodperc között).

² Ezzel szemben egyes modern, párhuzamos működésű számítógépek központi memóriáján több processzor is osztozik. Ekkor a memória bizonyos részeinek elérési ideje eltérő lehet a különböző processzorokra vonatkozóan. Az egyik processzor akár háromszor olyan gyorsan is el tudja érni a memória egy részét, mint a másik processzor.

A számítógépes mennyiségek 2 hatványai

Általában úgy beszélünk a számítógép komponenseinek méretéről, kapacitásáról, mintha 10-nek lenne valamilyen hatványa, azaz megabájtot, gigabájtot és hasonlíkat mondunk. A valóságban ezek a számok tulajdonképpen a legközelebbi 2-hatvány rövidítései. Elmögött az húzódik meg, hogy úgy a legcélszerűbb megtervezni a komponenseket, például memóriachipeket, hogy 2-hatvány számú bitet tároljon. Mivel $2^{10} = 1024$ nagyon közel van ezerhez, ezért gyakran úgy tesszünk, mintha 2^{10} egyenlő lenne ezerrel, és emiatt 2^{10} esetében a „kilo”, 2^{20} esetében a „mega”, 2^{30} esetében a „giga”, 2^{40} esetében a „tera” és 2^{50} esetében a „peta” szócskákat használjuk előtagként, még akkor is, ha ezek az előtagok a tudományos beszédmódban valójában a 10^3 , 10^6 , 10^9 , 10^{12} , 10^{15} számokra utalnak. Az eltérés annál nagyobb, minél nagyobb számokról beszélünk. Egy „gigabájt” valójában $1,074 \times 10^9$ bájt.

Ezekre a számokra a következő szabványos rövidítéseket használjuk: K felel meg a kilónak, M a megának, G a gigának, T a terának és végül P a petának. Így tehát 16 Gbájt 16 gigabájtot jelent, ami szigorúbb értelemben 2^{34} bájt. Mivel időnként olyan számokról is akarunk beszélni, amelyek 10-nek hagyományos értelemben veti hatványai, például ezer, millió stb., ezért ezeket az elnevezéseket továbbra is megtartjuk ezekre a hagyományos számokra, azaz ilyen esetben nem használjuk a „kilo”, „mega” stb. előtagokat. Például „egymillió bájt” 1 000 000 bájtot, míg „egy megabájt” 1 048 576 bájtot jelent.

2.1.3. Virtuális memória

Program írásakor az általunk használt adatok – a program változói, a beolvasott fájlok stb. – egy *virtuális memória címtartületet* foglalnak le. A program utasításai szintén a nekik megfelelő címtartületet foglalják le. Sok gép használ 32 bites címeket, vagyis összesen 2^{32} , azaz körülbelül 4 milliárd különböző címet lehet megadni. Mivel minden bájtának szüksége van saját címre, ezért a tipikus virtuális memóriát 4 gigabájtnak tekinthetjük.

Abból kifolyólag, hogy a virtuális memória területe sokkal nagyobb, mint a szokásos központi memóriáé, az következik, hogy a teljesen kitöltött virtuális memória tartalmának legnagyobb részét valójában a lemezen tároljuk. A lemezműveletek közül a legtipikusabbakat a 2.2. részben foglalkoztatjuk. Pflananyagilag elég annyit tudnunk, hogy a lemez logikailag *blokkokra* van felosztva. A szokásos lemezek blokkmérete 4 K és 56 K, azaz 4 és 56 kilobájt között van. A virtuális memóriát a lemez és a központi memória között teljes blokkokban mozgathatjuk. A központi memóriában a blokkokat *lapoknak* (page) hívjuk. A gép hardvere és az operációs rendszer megengedi, hogy a virtuális memória lapjait a központi memória tetszőleges részére lehessen behozni, miközben a blokkok minden bájtjára szabályosan lehet hivatkozni a virtuális memória címük alapján. A könyvünkben nem foglalkozunk azzal, hogy ezt milyen mechanizmussal lehet elérni.

Moore törvénye

Gordon Moore sok évvel ezelőt megfigyelte, hogy az integrált áramkörök sok jellemzőjének fejlődése exponenciális görbét követ, meghozzá olyat, amely az értéket 18 hónaponta megduplázza. Az alábbiakban megadunk néhány paramétert, melyek Moore törvényének engedelmeskednek.

1. A processzorok sebessége, vagyis a másodpercenként végrehajtott utasítások száma és a processzor sebességének és árának aránya.
2. A központi memória egy bite jutó ára és az egy chipbe tehető bitek száma.
3. A lemez egy bite eső ára és a lemezen tárolható bájtok száma.

Másrészt vannak olyan fontos paraméterek is, melyek nem követik Moore törvényét, mivel lassabban vagy egyáltalán nem nőnek. Ezek között a lassan növekedő paraméterek között szerepel az a sebesség, hogy a központi memóriában milyen gyorsan lehet az adatokat elérni, vagy az a sebesség, amilyen gyorsan a lemez forog. Mivel ezek lassan nőnek, ezért a lemaradásuk egyre nagyobb. Emiatt az az idő, ami alatt az adatokat a memóriahierarchia szintjei között mozgatjuk, egyre tovább növekszik a számítási időhöz viszonyítva. Ennélfogva az elkövetkezendő években az várható, hogy a központi memória a cache-hez képest sokkal távolabb kerül a processzortól, és ugyanakkor a lemez adatai is még távolabb kerülnek a processzortól a fenti paraméterek tekintetében. Ennek az érzékelhető távolviságnak 1999-ben már komoly hatásaival kellett számolni.

A virtuális memóriát is magában foglaló 2.1. ábrán az útvonal a hagyományos programok és alkalmazások kezelését reprezentálja. Ez nem egyezik meg egy adatbázis adatlainak tipikus kezelésével. Ennek ellenére egyre nő az érdeklődés a *központi memória adatbázisrendszerek* iránt, melyek az adatokat ténylegesen a virtuális memórián keresztül kezelik, és ehhez az operációs rendszerre támaszkodnak, mivel ennek a lapozási mechanizmusával hozzák be a szükséges adatokat a központi memóriába. A központi memória adatbázisrendszerek a legújabb alkalmazáshoz hasonlóan akkor a leghasznosabbak, ha az adatok mérete eléggé kicsi ahhoz, hogy a központi memóriában maradjanak anélkül, hogy az operációs rendszernek ki kelljen vinnie őket. Ha a gép 32 bites címtáblával rendelkezik, akkor a központi memória adatbázisrendszerek olyan alkalmazásokhoz jók, amelyeknél nem szükséges 4 gigabájtól több adatot egyszerre a memóriában tartani. (A 4 gigabájtól kisebb értéket kell venni, ha a gép tényleges központi memóriája 2³² bájtól kisebb.) Ez a támmemória igen sok alkalmazás számára elegendő, de már nem elég az adatbázisrendszerek nagy és igényes alkalmazásához.

A fentiek miatt a nagyméretű adatbázisrendszerek az adatokat közvetlenül a lemezen kezelik. Ezeknek a rendszereknek a méretét csak az korlátozza, hogy összesen mennyi adatot lehet tárolni az összes lemezen és a számítógépes rendszer által elérhető egyéb adattároló eszközön. Ezzel a működési móddal a következőkben fogunk foglalkozni.

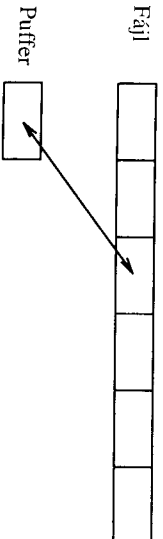
2.1.4. Másodlagos tárolás

Lényegében minden számítógépnek van valamilyen *másodlagos tárolója*, amely a tárolónak olyan formája, mely lényegesen lassabb, de lényegesen nagyobb kapacitással, mint a központi memória, ugyanakkor lényegében véletlen hozzáféréstű, azaz különböző adatelemek eléréséhez szükséges időtartamok között viszonylag kicsi az eltérés. (Az eltérés okairól a 2.2. részben lesz majd szó.) A modern számítógépes rendszerek másodlagos memóriának valamilyen lemezt használnak. Ez a lemez rendszerint mágneses lemez, bár néha optikai vagy mágneses-optikai lemezeket is használnak. Ez utóbbi típusok olcsóbbak, de lehet, hogy nem támmogatják azt, hogy könnyen lehessen adatokat írni rájuk, ha egyáltalában ez lehetséges, így ezeket általában csak olyan adatok archiválására szokták használni, melyek nem változnak.

A 2.1. ábrán észrevehetjük, hogy a lemezt úgy tekintjük, mint amely mind a virtuális memóriát, mind a fájlrendszert támmogatja. Tehát, míg bizonyos lemezblokkok arra a célra szolgálnak, hogy egy alkalmazói program virtuális memóriájának lapjait tárolják, addig más lemezblokkok fájlkat vagy fájlkat részeit támmogazzák. A fájlkat a lemez és a központi memória között blokkokban mozgatjuk, és ezt a folyamatot az operációs rendszer vagy az adatbázisrendszer felügyeli. Egy blokk mozgását a lemeztől a központi memóriába *lemezolvasásnak* hívjuk, míg a központi memóriából a lemeze mozgást *lemezírásnak* nevezzük. Mindkét műveletre *lemez I/O³-ként* fogunk hivatkozni. A központi memória bizonyos részei arra használhatók, hogy a fájlkat *pufferreljék*, vagyis ezeknek a fájlloknak blokkméretű darabjait tárolják.

Például, mikor olvasásra nyitunk meg egy fájl, akkor az operációs rendszer valószínűleg lefoglal egy 4 KB-os blokkot a központi memóriában, ami egy puffer lesz ehhez a fájlhoz, feltéve, hogy a lemezblokkok 4 Kbájt méretűek. Kezdetben a fájl első blokkja másolódik a pufferbe. Mikor az alkalmazói program feldolgozza a fájlkat ezt a 4 Kbájtát, akkor a fájl következő blokkja másolódik a pufferbe, lecsereelve annak régi tartalmát. Ez a 2.2. ábrán látható folyamat addig folytatódik, míg a teljes fájl beolvasásra nem kerül, vagy amíg a fájl le nem zárjuk.

Az adatbázis-kezelő rendszer a lemezblokkokat saját maga kezeli, és nem hagyatkozik az operációs rendszer fájlkezelőjére, mikor blokkokat kell mozgatni a központi és a másodlagos memória között. Az is igaz viszont, hogy a kezelés alapvető pontjai lényegében megegyeznek, függetlenül attól, hogy egy fájlrendszert vagy egy adatbázis-



2.2. ábra. Egy fájl és a központi memóriában hozzá tartozó puffer

³ Az I az input (bemenet), az O az output (kimenet) szavak kezdőbetűjét jelölik. A fordító megjegyzése.

zrendszerét nézzük. Durván 10–30 milliszekundum (0,01–0,03 másodperc) ideig tart, hogy egy lemezblokkot olvassunk vagy írjunk. Ez alatt az idő alatt egy tipikus gép egymillió műveletet is végrehajthat. Ennek az a következménye, hogy általában egy lemezblokk olvasására vagy írására fordított idő a domináns feldolgozási idő, függetlenül attól, hogy a blokk tartalmával mit csinálunk. Emiatt lényeges fontosságú, hogy ha lehetséges, akkor az a lemezblokk, amiben a szükséges adatok szerepelnek, már a központi memóriának egy pufferebben legyen jelen, mert ezután nem kell a lemez I/O-költségével számolnunk. Erre a problémára a 2.3. és a 2.4. részekben fogunk vissza térni, ahol példákat fogunk látni arra, hogyan foglalkoztunk azzal a nagy költséggel, amivel a memóriahierarchia szintjei közti adatmozgatás jár.

1999-ben a lemezeységek mérete általában 1-től 10-ig, vagy még ennél is több gigabájtig terjed. Ezenfelül a gépek több lemezeységet is használhatnak, így egy önálló gép reálisan akár 100 gigabájt kapacitású másodlagos memóriával is rendelkezhet. Végül is a másodlagos memória 10^5 nagyságrendben lassabb, de legalább 100-szor nagyobb kapacitási, mint a tipikus központi memória. A másodlagos memória jelentősen olcsóbb, mint a központi memória. 1999-ben a mágneses lemezeységek egy megabájtra jutó ára 5 és 10 cent közötti, míg a központi memóriánál ugyanez 1 és 2 dollár között mozog.

2.1.5. Harmadlagos tárolás

Bármilyen nagy is lehet több lemezeységségs együttes kapacitása, vannak olyan adatbázisok, amelyek sokkal nagyobbak annál, hogy egy vagy akár sok gép lemezein lehessen tárolni őket. Például áruházláncok terabájnyi adatot őriznek a forgalmukkal kapcsolatban. A műholdas felvételek alapján össze gyűjtött adatok is gyakran terabájtokban mérhetők, sőt a műholdak a közeljövőben petabájt (10^{15} bájt) információt fognak visszaadni évente.

Ezeknek az igényeknek a kielégítésére fejlesztették ki a *harmadlagos tárolókat*, melyek terabájtokban mérhető adatmennyiséget tudnak tárolni. A harmadlagos tárolókat azzal lehet jellemezni, hogy a másodlagos tárolókhöz képest sokkal lassabban olvassák vagy írják az adatokat, de ugyanakkor sokkal nagyobb a kapacitásuk, és ráadásul a mágneses lemezekhez viszonyítva kisebb az egy bájtra jutó költségük. Míg a központi memória egyforma idő alatt ér el tetszőleges adatot, a lemez esetében a különböző adatok elérési ideje csak kis ténylegében különbözik, addig a harmadlagos tárolóeszközök esetében az elérési időik széles sávban változhatnak attól függően, hogy milyen közel van az adat az olvasási/írási ponthoz. Az alapvető harmadlagos tárolóeszközök a következők:

1. *Ad hoc szalagos tárolás*: A legegyszerűbb – és az elmúlt években sokáig az egyetlen – megvalósítása a harmadlagos tárolónak az, hogy az adatokat orsós vagy kazettás szalagokra menjük ki, és ezeket utána tárolókba helyezzük. Amikor valamilyen információra van szükség a harmadlagos tárolóról, akkor a kiszolgáló személyzetből egy operátor megkeresi, és felteszi a szükséges szalagot a szalagolvasóra. Az információ megkeresése úgy történik, hogy a szalagot a megfelelő pozícióra

tekerjük, és az információt a szalagról bemásoljuk a másodlagos vagy a központi memóriába. A harmadlagos tárolóra írás ennek a megfordítottja, azaz a megfelelő szalagon a megfelelő helyet keressük meg, és a lemeztől kimásoljuk az információt a szalagra.

2. *Többlemez optikai lemeztár* (juke box⁴). A lemeztár CD-ROM-tárolókból áll. (CD = kompaktlemez [compact disk]; ROM = csak olvasható memória [read-only memory]). Ezek azok az optikai lemezek, melyeken általában a szoftvereket forgalmazzák.) Az optikai lemezeken a biteket fekete vagy fehér kis területek reprezentálják, így a biteket úgy lehet elolvasni, hogy lézerral megvilágítjuk az adott helyet, és megmérzük, hogy a fény visszaverődik-e. Egy robotkar is része a lemeztárnak, amely gyorsan ki tudja emelni bármelyik CD-ROM-ot, és be tudja helyezni egy CD-olvasóba. A CD tartalmát vagy annak egy részét így be lehet olvasni a másodlagos memóriába. Speciális eszközök nélkül általában nem lehet írni a CD-kre. Már kaphatók viszonylag olcsó CD-írók, és valószínűleg hamarosan gazdaságos lesz olyan harmadlagos tárolókat készíteni, amelyek írni és olvasni is tudják az optikai lemezeket.

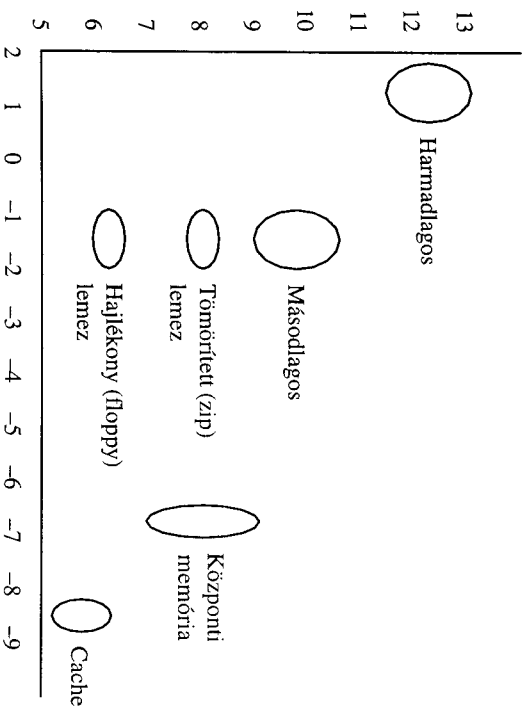
3. *Szalagsítlók*. A sítló egy szoba nagyságú eszköz, mely szalagátrolókat tartalmaz. A szalagokat robotkarokkal lehet elérni, melyek aztán a kiválasztott szalagot a szalagolvasók valamelyikéhez tudják vinni. Ekképpen a sítló a korábbi ad hoc szalag tárolásnak egy automatizált formája. Mivel a berendezés számítógéppel vezérelt és a szalag megkeresése is automatikus, ezért a működés legalább egy nagyságrenddel gyorsabb, mint az emberi szemlélyzettel működtetett ad hoc rendszeré.

A mágneses szagot tartalmazó kazetta kapacitása 1999-ben körülbelül 50 gigabájt. Ennél fogva a szalagsítlók sok terabájtot tudnak tárolni. A CD-k szabványosan 2/3 gigabájtot képesek tárolni, de a következő generációs szabvány szerint ez a kapacitás körülbelül 2,5 gigabájtra fog emelkedni. Már most is kaphatók olyan CD-ROM-lemeztárak, melyek sok terabájt adat tárolására képesek.

A harmadlagos tárolóeszközök esetében az adat-hozzáférési idő néhány másodpercetől néhány percreig is terjedhet. Egy lemeztár vagy szalagsítló robotkarja a kívánt CD-ROM-ot vagy kazettát másodpercek alatt megtalálja, ezzel szemben egy embernek valószínűleg percekkel vesz igénybe, hogy megkeresse a szalagokat. Ha már egyszer betöltötték a CD-t az olvasóba, akkor a CD bármely része a másodperc törtérsze alatt elérhető, viszont ehhez képest sokkal több másodpercig is tarthat, míg egy szalag megfelelő részét a szalagolvasó eszközje alá tudjuk esévélni.

Végéredményben harmadlagos tárolás esetén az adatelérés körülbelül 1000-szer is lassabb lehet, mint a másodlagos memória esetében (az első esetében másodpercekben, a második esetben a másodperc ezredrészében mérhető ez az idő). Ezzel szemben a harmadlagos tárolóeszközök 1000-szer nagyobb kapacitásúak, mint a másodlagos eszközök (itt terabájtok állnak szemben gigabájtokkal). A 2.3. ábra egy log-log skálán mutatja az elérési időket és a kapacitások közti kapcsolatot a memóriahierarchia mind a négy általunk tárgyalt szintjén. Az ábrán szerepeltejük a tömörített (zip) és a hajlé-

⁴ A juke box szó szerint több lemezzel működő zenegépet jelent. A fordító megjegyzése.



2.3. ábra. Az elérési idő és a kapacitás összevetése a memóriahierarchia különböző szintjein

kony (floppy) lemezt is, mivel ezek szintén általában használt tárolóeszközök, habár adatabázisok esetén másodlagos tárolásra nem túl tipikus a használatuk. A vízszintes tengely beosztása tízhavvány értékű másodperceket jelöl, azaz például a -3 valószínűleg 10⁻³ másodpercet, vagyis egy milliszekundumot jelöl. A függőleges tengely a bájtokat szintén 10 havványként jelöli, azaz például a 8-as 100 megabájtot jelent.

2.1.6. Felejtő és nem felejtő tárolás

Egy további megkülönböztetés az adattároló eszközök terén az, hogy vajon *felejtelenek* vagy *nem felejtelenek*. A felejtő eszköz, mint ahogy a neve is mondja, elfelejt minden benne tárolt adatot, ha áramszünet történik. Ezzel szemben a nem felejtő eszköz hosszú ideig épségben megőrzi a tartalmát még akkor is, ha kikapcsolják vagy ha áramkimaradás történik. A felejtés kérdése valóban fontos, mivel az adatabázis-kezelő rendszerek egyik jellemző vonása éppen az, hogy képesek megőrizni az adatokat áramszünet esetén is.

A mágneses anyagok megtartják a mágnességüket áramhiány esetében is, így a mágneses lemezek és szalagok nem felejtő tárolók. Ugyanígy a CD-hez hasonló eszközök is megtartják a beléjük égetett feketé és fehér pöttyöket áram jelenléte nélkül is. Sok ilyen eszköz esetében valójában egyetlen megváltoztatható, amit a felhírtükre írtak. Ebből kifolyólag az összes másodlagos és harmadlagos tárolóeszköz nem felejtő típusú.

A központi memória viszont általában felejtő típusú. Kiderült, hogy egy memóriachipet egyszerűbb áramkörökből lehet megtervezni, ha az is megengedett, hogy egy bit értéke bizonyos idő, például egy perc múlva megváltozzon. Ez az egyszerűsítés csökkenti a chip egy bitjére eső költségét. Tulajdonképpen az történik, hogy a bitet reprezentáló elektromos töltés lassan elfolyik abból a tartományból, mely ennek a bit-

nek volt kijelölve. Emiatt egy úgynevezett *dinamikus véletlen hozzáférési memóriachipre*, DRAM-ra (dynamic random-access memory) van szükség, amely periodikusan olvassa és újratírja a teljes tartalmát. Ha kimegy az áram, akkor ez a frissítés nem működik, és a chip hamarosan elveszi, amit tárolt.

Ha egy adatabázisrendszer felejtő központi memóriájú gépen fut, akkor minden változást ki kell menteni a lemezre, mielőtt a változást az adatabázis részének tekintnénk, mert különben azt kockáztatjuk, hogy áramkimaradás esetén elveszítjük az információt. Ennek következményeként a lekérdezések és az adatabázis-módosítások nagyszámú lemezírást fognak eredményezni, melyek bizonyos részére nem is lenne szükség, ha nem lennénk kénytelenek minden információt minden időben megőrizni. Egy másik lehetőség az, hogy olyan központi memóriát használjunk, amely nem felejtő. A nem felejtő memóriachipek egyre olcsóbbá válnak, új típusa a *flash memória*. Egy további lehetőség, hogy a hagyományos memóriachipekből egy úgynevezett RAM-lemezt készítsünk, melynek a tápegységét egy elemmel is kiegészítjük.

2.1.7. Feladatok

2.1.1. feladat: Tegyük fel, hogy 1999-ben egy tipikus számítógép 500 megahertz gyorsaságú processzorral rendelkezik, 10 gigabájt lemeze és a központi memóriája 100 megabájtos. Tegyük fel továbbá, hogy Moore törvénye, miszerint ezek a tényezők 18 hónaponta megduplázódnak, az idők végtelenségéig igaz marad.

- * a) Mikor lesznek tipikusak a terabájt kapacitású lemezek?
- b) Mikor lesz tipikus a gigabájt központi memória?
- c) Mikor lesz tipikus a terahertz gyorsaságú processzor?
- d) Milyen lesz egy tipikus konfiguráció (processzor, lemez, memória) 2008-ban?

! 2.1.2. feladat: Data paramcsnok, aki a *Star Trek: The Next Generation* közismert amerikai sci-fi szappanopera sorozat egyik android szereplője a 24. századból, egyszerűen kijelentette, hogy az ő processzora 12 teraop gyorsaságú, azaz 12 tera műveletet hajt végre másodpercenként. Igaz, hogy a műveletek száma és a processzor órajelének frekvenciája általában nem egyezik meg, de most tegyük fel, hogy azonosak, azaz Data processzora 12 terahertz gyorsaságú. Tegyük fel, hogy Moore törvénye még 400 évig fennáll. Ekkor valójában milyen gyorsaságú lenne Data paramcsnok processzora?

2.2. Lemezek

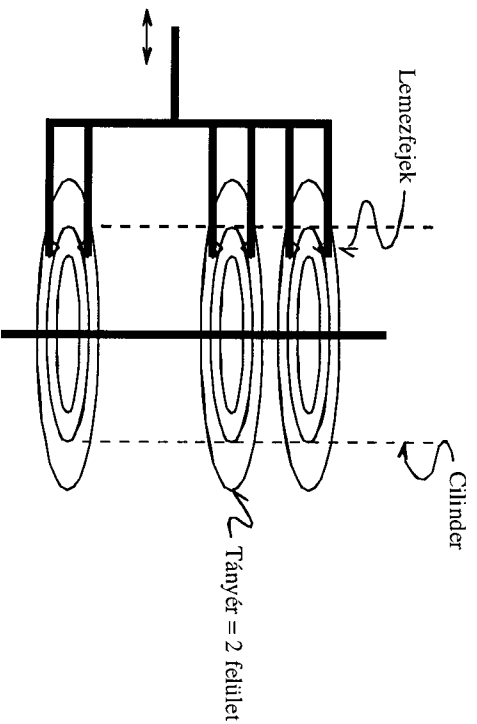
Az adatabázis-kezelő rendszerek egyik fontos jellemzője a másodlagos tárolók használata. A másodlagos tárolás szinte kizárólag mágneses lemezeken alapul. Így tehát ahhoz, hogy az adatabázis-kezelő rendszerek implementálásában használhatjuk a megindokolhassuk, először meg kell vizsgálnunk részletesen a lemezek működését.

2.2.1. A lemezek mechanikája

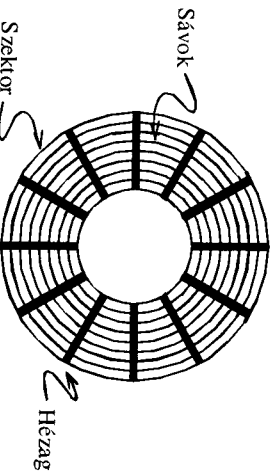
A 2.4. ábrán egy lemezmeghajtó két alapvető fontosságú mozgó részét láthatjuk, ezek a *lemezgyűjtemény* és a *fejserelvény*. A lemezgyűjtemény egy vagy több kör alakú *tányértől* áll, melyek egy központi tengely körüli forognak. A tányérok felső és alsó felülete is vékonyan be van vonva egy mágneses anyaggal, amelyek a biteket tárolják. A 0-t egy kis terület mágnességének egyik iránya reprezentálja, míg az 1-et pont az ellenkező irányú mágnesség jelenti. A lemeztányérok manapság szokásos átmérete 3,5 hüvelyk, bár készülték lemezek egy hüvelyktől több láb hosszú terjedő átmérettel is.

A bitek tárolására szolgáló helyet *sávokba* (track) szerveztük. A sávok koncentrikus köröket jelentenek egyetlen tányéron. A sávok a tányér felületét szinte teljesen betöltik, kivéve a tengelyhez legközelebbi területet, ahogy az a 2.5. ábrán felülnezetben látható. Egy sáv sok pontból áll, melyek mindegyike egyetlen bitek reprezentálással, hogy milyen irányú abban a pontban a mágnesség iránya.

A sávok *szektorokba* vannak szervezve. A szektorok a kör olyan szeletei, melyek *hézaggal* (gap) vannak elválasztva. A hézag attól hézag, hogy semelyik irányba sincs



2.4. ábra. Egy tipikus lemez



2.5. ábra. Egy lemez felszínének felülnezeze

Szektorok kontra blokkok

Ne feledjük, hogy a szektor a lemeznek egy fizikai egysége, míg a blokk egy logikai egység. A blokk a lemezt használó szoftverrendszer – operációs rendszer vagy például adatbázis-kezelő rendszer – alkotása. Ahogy már említettük, manapság a blokkok tipikusan legalább akkorák, mint a szektorok, és a blokkok egy vagy több szektorból állanak. Ennek ellenére nem világos, hogy a blokkok miért ne lehetnének egy szektoronak a töredékei, miáltal több blokkot lehetne egy szektorba betenni. Tulajdonképpen léteztek olyan korábbi rendszerek, amelyek pont ezt a stratégiát követték.

mágneseszev⁵ A lemez olvasásának és írásának tekintetében a szektor képezi a felbonthatatlan egységet. Ugyanígy a hibák tekintetében is oszthatatlan egységet képez. Ha történetesen a mágneses felület egy kis részen valahogy elromlik úgy, hogy ez a rész nem tud információt tárolni, akkor az ezt a részt tartalmazó szektor teljes egészében használhatatlanná válik. A hézagok, melyek gyakran a teljes sáv 10%-át is kitevő, arra használhatók, hogy segítséget nyújtsanak a szektorok elejének megtalálásához. A 2.1.3. részben említettük, hogy a blokkok olyan logikai adategységek, melyeket a lemez és a központi memória között mozgatunk. Egy ilyen blokk egy vagy több szektorból állhat.

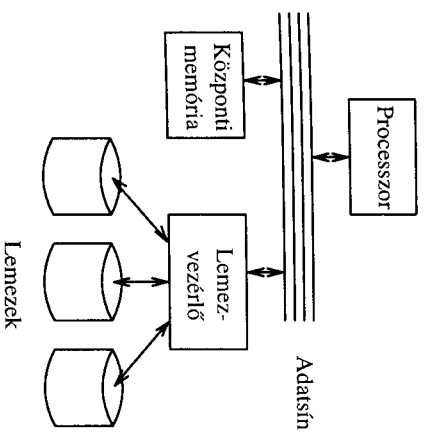
A 2.4. ábrán látható második mozgatható darab a fejserelvény, amely a *lemezfejeket* tartja, áll. Mindegyik felülethez egy fejtartozék, mely igen-igen közel kerül a felülethez, de sohasem érintkezik vele (mert ha belezuhana a fejtartozék, akkor a lemez tönkremegy, és minden elvesz, amit rajta tároltunk). A fejtartozék halad a lemez felületén, és arra is képes, hogy megváltoztassa ezt a mágnességet, miáltal információt ír a lemezre. A fejek mindegyike egy-egy karhoz csatlakozik. A különböző felületekhez tartozó karok együttesen mozognak ki vagy be. Így a karok is részei a merev fejserelvénynek.

2.2.2. A lemezvezérlő

Az egy- vagy többlemezes meghajtó vezérlését a lemezvezérlő végzi, mely a következő képességekkel jellemezhető kis processzor:

1. Vezérli azt a mechanikus szerkezetet, mely a fejserelvényt mozgatja úgy, hogy a fejeket pozicionálja egy adott sugár mentén. Ezen a sugáron minden felülethez tartozó fejtartozék alatt a felület egy sávja fog elhelyezkedni, mely így írható és olvasható.

⁵ A 2.5. ábrán minden sáv ugyanannyi szektorra van felosztva, bár a 2.1. példában látni fogunk olyan esetet, mikor a szektorok száma sávonként különböző lehet, nevezetesen a külső sávon több szektor van, mint a belsőkn.



2.6. ábra. Egy egyszerű számítógépes rendszer vázlatos felépítése

Az összes fejt alatti egy időben elhelyezkedő összes sávot együttesen *cilindernek* nevezzük.

2. Kiválasztja azt a felületet, amelyről olvasni kell, vagy amelyre írní akarunk. Kiválasztja az ezen a felületen a fejt alatti sáv egy szektorát. A vezérlő felelős azért is, hogy észrevegye, hogy a forgó tengely mikor jutott el abba a helyzetbe, hogy a kívánt szektor éppen a fejt alatt kezdődik.

3. Átadja a kívánt szektorból kiolvasott biteket a számítógép központi memóriájába vagy fordítva, a központi memóriából vett biteket kiírja a kívánt szektorba.

A 2.6. ábra egy egyszerű egyprocesszoros számítógépet mutat. A processzor az adatsínen (data bus) keresztül tartja a kapcsolatot a központi memóriával és a lemezvezérlővel. A lemezvezérlő sok lemezt is irányíthat; az ábrán ennek a számítógépnek három lemezt tüntettük fel.

2.2.3. A lemeztárolók jellemzői

A lemeztechnológiák állandóan változnak, ahogy egyre csökken az egy bit tárolásához szükséges hely. 1999-ben a lemezekkel kapcsolatos tipikus mértékszámok a következők:

- A lemezgyűjtemény *forgási sebessége*. Általános az 5400 fordulat/perc (rpm), azaz egy körtüfordulás 11 milliszekundum alatt történik, de találhatóak ennél magasabb és alacsonyabb értékek is.
- Az *egysegységhez tartozó tányérók száma*. A tipikus meghajtónak körülbelül 5 tányérja és emel fogva 10 felülete van, de 30 felülettel rendelkező lemez meghajtókat is lehet találni. A szokásos hajlékony- és tömörített lemezeknek egyetlen tányérjuk van két felülettel. Az egyoldalas hajlékonylemez már elavult, de azért még ilyen is lehet találni. Ennek egy tányérja van, de csak egyetlen felületét lehet használni.

- A *felületen található sávok száma*. Egy felületen akár 10 000 sáv is lehet, bár a hajlékonylemezeknél ez a szám sokkal kisebb. Lásd a 2.2. példát.
- A *sávokra jutó bájtok száma*. A szokásos lemezmeghajtóban 10^5 vagy több bájttal jut egy sávra, de a hajlékonylemezek sávjai természetesen kevesebb bájttal tartalmazzanak. Már említettük, hogy a sávokat szektorokra osztjuk. A 2.5. ábrán 12 szektor szerepel minden sávban, de a modern lemezek esetében 500 szektor is juthat egy sávba. A szektorok egyenként körülbelül 512 és 4096 közé eső számú bájttal tartalmazhatnak.

2.1. példa: A *Megarion 747* lemeznek a következő jellemző paraméterei vannak, melyek egyébként is tipikusak a középmeretű vintage-1999 lemezmeghajtó esetén.

- Négy tányérja van nyolc felülettel.
- 2^{13} , azaz 8192 sáv van mindegyik felületen.
- Átlagosan $2^8 = 256$ szektor van minden sávban.
- $2^9 = 512$ bájt van minden szektorban.

A lemez kapacitását úgy kapjuk, hogy összeszorozzuk ezeket a számokat, azaz 8 felület szer 8192 sáv szer 256 szektor szer 512 bájt, az annyi, mint 2^{33} bájt. Tehát a *Megarion 747* egy 8 gigabájtos lemez. Egy sáv 256-szor 512 bájt, azaz 128 Kbájttal tartalmaz. Ha a blokkok 2^{12} , azaz 4096 bájtosak, akkor egy blokk 8 szektor használat fel, így $256/8 = 32$ blokk jut egy sávra.

A *Megarion 747* felületének átmérője 3,5 hüvelyk. A sávok a felületek külső részére esnek, és a belső részen 0,75 hüvelyk nincs lefoglalva. A sugárirányban vett bitsűrűség így 8192 bit/hüvelyk, mert ennyi a sávok száma.

A sávokra számolt bitsűrűség sokkal nagyobb. Először tegyük fel, hogy minden sáv az átlagos számmal, 256-tal megegyező számú szektor tartalmaz. Tételizzük fel azt is, hogy a hézagok a sávok 10%-át foglalják el. Ekkor a sáv 128 Kbájta (ami 1 Mbit) a sáv 90%-át foglalja el. A legkülső sáv hossza 3,5 π , azaz körülbelül 11 hüvelyk. Ennek a hosszának a kilencven százalékka, azaz körülbelül 9,9 hüvelyk tartalmaz 1 Mbitet. Ennél fogva a sávnak a bielek tárolására lefoglalt részen a bitsűrűség körülbelül 100 000 bit/hüvelyk.

Másrészt a legbelső sáv átmérője csak 1,5 hüvelyk. Így $0,9 \times 1,5 \times \pi$, azaz 4,2 hüvelyk kellene 1 Mbitet tárolni. A bitsűrűség a belső sávokon tehát 250 000 bit/hüvelyk körül lesz.

Mivel a sűrűség a belső és külső sávokon nagyon eltér, ha a szektorok és bitek minden sávra egyformák lennének, ezért a *Megarion 747*, más modern meghajtókhöz hasonlóan a külső sávokon több szektor tárol, mint a belső sávokon. Például 256 szektor tárolhatunk sávonként a lemez középső harmadában, de csak 192 szektor a belső harmadában, ugyanakkor 320 szektor a külső harmadba eső sávokban. Ha így leszünk, akkor a sűrűség a legkülső és legbelső sávok bitsűrűsége között változna, azaz 114 000 bit/hüvelyk és 182 000 bit/hüvelyk közé esne.

2.2. példa: A lemezek összehasonlításának az egyik végen szerepel a szabványos 3,5 hüvelykes hajlékonylemez. Ennek két felülete van, mindegyiken 40 sáv található, azaz

összesen 80 sáv. A lemez kapacitása körülbelül 1,5 Mb/ajt adat, függetlenül attól, hogy MAC-en vagy PC-n formáztuk meg. Ez azt jelenti, hogy 150 000 bit (18 750 bájt) jut minden sávra. A rendelkezésre álló területnek körülbelül a negyedét a hézagok és más lemezadminisztrációs részek töltik ki, mindkét típusú formázás esetén. □

2.2.4. A lemezhozzáférés jellemzői

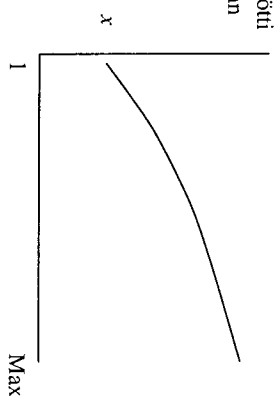
Az adatbázis-kezelő rendszerek tanulmányozása során nemcsak arra van szükségünk, hogy megértsük, miként tároljuk az adatokat a lemezeken, hanem arra is, hogy hogyan történik az adatok kezelése. Mivel minden számítás a központi vagy a cache memóriában történik, ezért az egyetlen lényeges kérdés a lemezekkel kapcsolatban az, hogyan mozgadjuk az adatblokkokat a lemez és a központi memória között. Már a 2.2.2. részben megemlítettük, hogy blokkokat (azaz azokat az egymás utáni szektorokat, melyek a blokkot tartalmazzák) akkor lehet írni vagy olvasni, mikor:

- a) a fejek arra a cillindere állnak, amelyik tartalmazza azt a sávot, melyen a blokk elhelyezkedik, és
- b) a blokkot tartalmazó szektorok a lemezfej alá kerülnek a teljes lemezgyűjtemény forgása által.

A blokkolvasási parancs kiadásának időpontja és a blokk tartalmának központi memóriába kerülésének időpontja közti eltelt időt a lemez késésének (latency) hívjuk. Ez a következő komponensekből áll össze:

1. A milliszekundum tört részével egyező idő telik el, míg a processzor és a lemezezerítő feldolgozza az igényt. Ezt az időt a továbbiakban elhanyagoljuk. Mivel más folyamatok is olvashatják és írhatják ugyanakkor a lemezt, ezért a kiszolgálásért versenyeznek a folyamatok a lemezezerítő és az adatsín esetében is, ami szintén késlekedést eredményezhet, de ezt az időt is elhanyagoljuk.
2. Időbe telik, míg a fejszerelvényt a megfelelő cillindere állítjuk. Ezt hívjuk *keresési időnek*, ami akár 0 is lehet, ha a fej véletlenül pont a megfelelő cillinderen áll. Ha

3x és 20x közötti tartományban

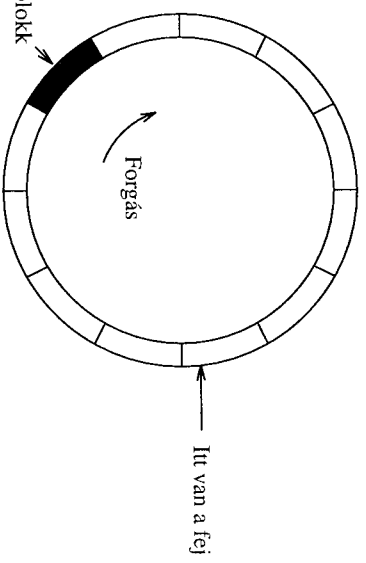


Az átléptett cillinderek száma

2.7. ábra. Keresési idő a meglett távolság függvényében

nem áll a cillinderen, akkor a fejszerelvénynek minimális időbe telik, míg elkezd mozogni, az is időbe telik, míg megáll, és ehhez adódik még a megtett távolsággal nagyjából arányos idő. Az elinduláshoz, a következő sávra lépéshez, a megálláshoz szükséges minimális idő tipikusan néhány milliszekundum, míg az összes sávon keresztülhaladás maximuma 10 és 40 milliszekundum között változik. A 2.7. ábra érkeketi, hogyan változik a keresési idő a távolság függvényében. Azt látni értől, hogy ha egycillindernyi távolság megtételéhez valamilyen x értékű keresési időre van szükség, akkor a maximális keresési idő $3x$ és $20x$ közé esik. Gyakran az átlagos keresési idővel jellemzik a lemez sebességét. A 2.3. példában látni fogjuk, hogy kell kiszámolni ezt az átlagot.

3. Ahhoz is idő kell, hogy a lemez úgy forduljon, hogy a blokkot tartalmazó szektorok közül az első kerüljön a fej alá. Ezt hívjuk *rotációs késésnek*. Egy tipikus lemez körülbelül egyszer fordul teljesen körbe 10 milliszekundum alatt. A kívánt szektor átlagosan felúton helyezkedik el a körön a fejekhez képest, így fél fordulatra van szükség, hogy elérjék a megfelelő cillindert. Az átlagos rotációs késés tehát körülbelül 5 milliszekundum. A 2.8. ábra mutatja be a rotációs késés problémáját.



2.8. ábra. A rotációs késés oka

4. *Árviteli időnek* nevezzük azt az időtartamot, ami alatt a blokk szektorai és a közöttük levő hézagok forgás közben elhaladnak a fej alatt. Mivel a tipikus lemez körülbelül 100 000 bájt tartalmaz sávonként, és nagyjából 10 milliszekundumonként fordul egyet, ezért ez az jelenti, hogy körülbelül 10 Mb/ajt lehet olvasni a lemeztől másodpercenként. Így az árviteli idő egy 4096 bájtos blokk esetén kevesebb, mint fél milliszekundum.

2.3. példa: Vizsgáljuk meg, hogy mennyi időbe telik egy 4096 bájtos blokkot beolvasni a *Megatron 747* lemeztől. Először is ismernünk kell a lemez egyes időparamétereit:

- A lemez forgási sebessége 3840 fordulat/perc (rpm); vagyis egy teljes körfordulatot 1/64 másodpercenként tesz meg.
- A fejszerelvény mozgásánál az elindulás és megállítás egy milliszekundumig tart. Minden 500 cillindertel történő elmozdulás további egy milliszekundumot jelent.

Tehát a fejek egy sávot 1,002 milliszekundum alatt tesznek meg. Így ahhoz, hogy a legkülső sávból a legbelső sávig terjedő 8191 sáv távolságot megtegyék összesen körülbelül 17,4 milliszekundum szükséges.

Számoljuk ki a 4096 bájtos blokk olvasásához szükséges minimális, maximális és átlagos időt. A minimális idő éppen az átviteli idő, mivel a vezérlőre vonatkozó sorban állási időtől és az egyéb adminisztrációs időtől eltekintünk. Ez azt jelenti, hogy ekkor már a blokkot tartalmazó sáv feletti tartózkodik a fej, és ráadásul éppen a blokk első szektorra fog elhaladni a fej alatt.

Mivel a *Megarion 747* lemezen egy szektorban 512 bájt van (a 2.1. példában adtuk meg a lemez fizikai jellemzőit), így a blokk nyolc szektorot foglal el. A fejnek ezért összesen nyolc szektor és a közöttük levő hét hézag fölött kell elhaladnia. Emlékezzünk vissza arra, hogy a hézagok a kör 10%-át foglalják el, és a szektoroké a maradék 90%. Egy kör mentén 256 szektor és 256 hézag található. Így a hézagok a 360 fokos szögéből összesen 36 fokot fednek le, míg a szektorok 324 fokot, azaz a 8 szektor és 7 hézag által lefedett szög összesen:

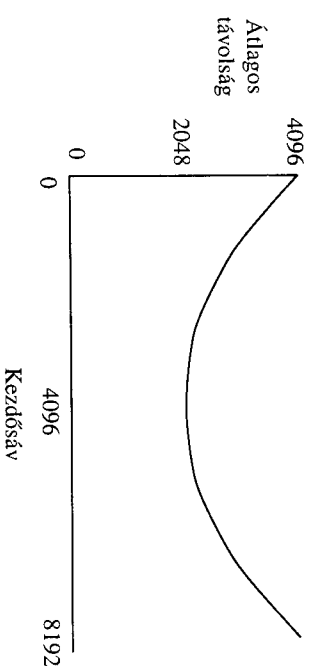
$$36 \times \frac{7}{256} + 324 \times \frac{8}{256} = 11,109.$$

Az átviteli idő emiatt (11,109/360)/64 másodperc: mivel osztani kell 360-nal, hogy megkapjuk, hogy a teljes körtőlfordulásnak milyen tört részére van szükség, és aztán osztani kell még 64-vel, mert a *Megarion 747* 64-szer fordul körbe egy másodperc alatt. Így az átviteli idő, ami egyben a minimális késés is körülbelül 0,5 milliszekundum.

Most nézzük meg, hogy egy blokk olvasásához maximálisan mennyi időre lehet szükség. Az a legrosszabb eset, ha a fejek a legbelső cilindren állnak, míg a kívánt blokk a legkülső cilindren helyezkedik el (vagy éppen fordítva). Először is a vezérlőnek el kell vinnie a fejeket a megfelelő helyre. Már az előbb észrevettük, hogy a *Megarion 747*-nek 17,4 milliszekundumra van szüksége ahhoz, hogy a fejeket az összes cilindren keresztül vigye. Ez a mennyiség az olvasáshoz szükséges keresési idő.

Irányzatok a lemezvezérlők felépítésében

Mivel a digitális hardverek ára szemmel láthatóan csökken, ezért a lemezvezérlők is kezdenek egyre inkább számítógépre hasonlítani, általános célu processzorral és tekintélyes véletlen hozzáféréstű memóriával is rendelkeznek. Sok mindenre lehet ezeket a kiegészítő hardvereket használni. A lemezvezérlők beolvashatják és tárolhatják egy lemeznek a teljes sávját még akkor is, ha csak a sáv egyetlen blokkjára van szükségük. Ezzel a lehetőséggel nagymértékben lehet csökkenteni az átlagos blokkhozáférési időt, ha egy sáv valamennyi vagy a legtöbb blokkjára szükségünk van. A 2.4.1. részben mutatunk néhány alkalmazást a teljes sáv vagy teljes cylinder olvasására, írására.



2.9. ábra. A megtett átlagos távolság a fej kezdeti helyének függvényében

A legrosszabb esetben az történhet, hogy mikor a fejek a megfelelő cilindert érnek, éppen akkor haladt el a kívánt blokk eleje a fej alatt. Ha azt tételezzük fel, hogy a blokkot az elejétől kezdve kell olvasni, akkor lényegében egy teljes fordulatot kell várni, azaz 15,6 milliszekundumra (a másodperc 1/64 részére) van szükség ahhoz, hogy a blokk eleje ismét elérje a fejet. Amint ez megtörténik, már csak az átviteli idő, 0,5 milliszekundumot kell megvárni, hogy teljesen beolvassuk a blokkot. Tehát a legrosszabb esetben $17,4 + 15,6 + 0,5 = 33,5$ milliszekundumra van szükség.

Végül számítsuk ki a blokkolvasáshoz szükséges átlagos időt. A késés két komponensét könnyű kiszámolni: az átviteli idő mindig 0,5 milliszekundum és az átlagos rotációs késés a lemez fél fordulatához szükséges idő, azaz 7,8 milliszekundum. Ismét feltehetjük, hogy az átlagos keresési idő megegyezik azzal, amennyi a sávok felén történő áthaladáshoz szükséges. Ez persze nem teljesen igaz, mivel az a tipikus, hogy a fejek kezdetben valahol középtájon helyezkednek el, ezért átlagosan a fél távolságnál kevesebbre van szükség, hogy a kívánt cilindert eljussanak.

Most egy részletesebb becslést adunk arra, hogy a fejnek átlagosan mennyi sávot kell elmozdítania. Tegyük fel, hogy a fej kezdetben a 8192 cylinder bármelyikén egyforma valószínűséggel lehet. Ha az 1. vagy a 8192. cilindren áll, akkor az átírt sávok átlagos száma $(1 + 2 + \dots + 8191)/8192$, azaz körülbelül 4096. Ha a 4096. cilindren, azaz pont középen áll a fej, akkor egyforma valószínűséggel fog a lemezen kifelé vagy befelé mozdulni, így átlagosan a sávok negyedét, azaz 2048 sávot fog megtenni. Egy kis számolással belátható, hogyha a kezdő pozíció az első cilindertől a 4096. cilindertig változik, akkor a fej által szükségeszerűen megtett átlagos távolság négyzetesen csökken 4096-tól 2048-ig. Hasonlóan látható, hogy ha a kezdő pozíció 4096-tól 8192-ig változik, akkor a megtett átlagos távolság 4096-ig fog négyzetesen növekedni. Mindezt a 2.9. ábra szemlélteti.

Ha a 2.9. ábrán szemléltetett mennyiségeket integráljuk az összes pozícióra, akkor azt kapjuk, hogy az átlagos megtett távolság pont a lemez egyharmadát, azaz 2730 cilindert tesz ki. Emiatt az átlagos keresési idő egy milliszekundum plusz az az idő, ami a 2730 cylinder megfelleléséhez kell, azaz $1 + 2730/500 = 6,5$ milliszekundum.⁶ Az átl-

⁶ Vegyük észre, hogy a számítás nem veszi figyelembe azt a valószínűséget, mikor a fejet egyáltalán nem kell elmozdítani, de ez az eset mindössze egyszer fordul elő a 8192 esetből,

gos kérésre vonatkozó becslésünk így $6,5 + 7,8 + 0,5 = 14,8$ milliszekundum; ahol az egyes tagok rendre az átlagos keresési idő, az átlagos rotációs késés és az átlagos átviteli idő.

2.2.5. Blokkok írása

Egy blokk írási folyamata a legegyszerűbb formában teljesen hasonló ahhoz, ahogy egy blokkot olvasunk. A lemez feje a megfelelő cilinderen áll, megvárjuk, hogy a megfelelő szektor vagy szektorok forgás közben a feji alá kerüljenek, és most ahelyett, hogy olvasnánk a feji alatti adatot, arra használjuk a fejet, hogy új adatot írjon. Ekkor az íráshoz szükséges minimális, maximális és átlagos idő pontosan ugyanannyi, mint az olvasás esetében.

Bonyolódik a helyzet, ha azt is ellenőrizni akarjuk, hogy a blokkot helyesen írtuk-e ki. Ekkor meg kell várnunk még egy körfordulást, és vissza kell minden kiírt szektort olvasnunk, hogy ellenőrizhessük, hogy azt tároljuk ott, amit ki akartunk írni. A helyesség ellenőrzésének egy egyszerűbb módja, mikor ellenőrző összegeket használunk. Erről a 2.5.2. részben lesz szó.

2.2.6. Blokkok módosítása

Egy blokkot nem lehet közvetlenül a lemezen módosítani. Még akkor is, ha csak néhány bájt (például a blokkon tárolt néhány sor egyikének egy komponensét) kívánunk módosítani, akkor is a következőképpen kell eljárniunk:

1. Beolvassuk a blokkot a központi memóriába.
2. A központi memóriában a blokk másolatán elvégezzük a kívánt változtatást.
3. A blokk új tartalmát visszatérítjük a lemezre.
4. Ha szükséges, akkor ellenőrizzük, hogy az írás helyesen történt meg.

Ezek szerint egy blokk módosításához szükséges idő kiszámításához össze kell adni az olvasási időt, a központi memóriában a változtatás végrehajtásához szükséges időt (ez rendszertől elhanyagolható a lemez írási vagy olvasási idejéhez képest), az írási időt, és ha ellenőrzés is van, akkor a lemez még egy körfordulásához szükséges időt.⁷

feléve, hogy a kívánt blokkot véletlenszerűen adják meg. Másrészt az is igaz, hogy az a feltevés, miszerint a blokkot véletlenszerűen választják, valójában nem is mindig tekinthető helyesnek, ahogy ezt majd a 2.4. részben látni fogjuk.

⁷ Először ránézésre meglepőnek tűnhet, hogy miért tart ugyanannyi ideig írni egy éppen beolvasott blokkot, mint egy olyan blokkot, amit véletlenszerűen jelöltek ki írásra. Ha a fejek márának ott, ahol voltak, akkor tudjuk, hogy az íráshoz egy teljes körfordulást kell várnunk. Viszont a keresési idő zéró. Ezzel szemben az igaz, hogy a lemezvezető nem tudja, hogy egy alkalmazás mikor fejeződik be a blokk új értékének visszatérítésével, így megőrizhet, hogy a fejek közben elmozdulnak egy másik sávra, hogy végrehajtsanak valamilyen másik lemez I/O-műveletet, mielőtt a lemezvezető megkapja azt az igényt, hogy a blokk új értékét vissza kell írni.

2.2.7. Feladatok

2.2.1. feladat: A *Megarion 777* lemeznek a következő paraméterei ismertek:

1. Tíz felülete van, egyenként 10 000 sávval.
2. A sávok átlagosan 1000 darab egyaránt 512 bájt nagyságú szektort tartalmaznak.
3. Minden sáv 20%-át a házágok töltik ki.
4. A lemez forgási sebessége 10 000 fordulat percenként.
5. Ahhoz, hogy a fej n sávot mozduljon el, $1 + 0,001n$ milliszekundumra van szükség.

Válaszoljunk meg a *Megarion 777*-re vonatkozó következő kérdéseket.

- * a) Mekkora a lemez kapacitása?
- b) Ha minden sáv ugyanannyi szektort tartalmaz, akkor mekkora egy sáv szektorában a bitsűrűség?
- * c) Mekkora a maximális keresési idő?
- * d) Mekkora a maximális rotációs késés?
- e) Ha egy blokk 16 384 bájt nagyságú (32 szektor), mekkora egy blokk átviteli ideje?
- f) Mekkora az átlagos keresési idő?
- g) Mekkora az átlagos rotációs késés?

! 2.2.2. feladat: Tegyük fel, hogy a *Megarion 747* lemez feje a 1024. sávon helyezkedik el, vagyis a sávok 1/8 részénél áll. Tegyük fel, hogy a következő igény egy olyan blokkra vonatkozik, amely egy véletlenszerűen választott sávon helyezkedik el. Számítsuk ki ennek a blokknak az olvasásához szükséges átlagos időt.

***! 2.2.3. feladat:** A 2.3. példa végén azt számoltuk ki, hogy mekkora az az átlagos távolság, amit a fejnek kell megtennie, ha egy véletlenszerűen választott sávról egy másik véletlenszerűen választott sávra kell eljutnia, és azt kaptuk, hogy ez a távolság a sávok 1/3 része. Tegyük fel, hogy egy sáv szektorainak száma fordítottan arányos a sáv hosszával (vagy a sugárral), így a bitsűrűség minden sávra megegyezik. Tegyük fel azt is, hogy a fejet egy véletlenszerűen választott szektorról egy másik véletlenszerűen választott szektorra kell mozgítani. Mivel a lemez külső része felé haladva a sávokban egyre több szektor gyűlik össze, ezért azt várhatjuk, hogy a fej átlagos mozgásánál a sávok kevesebb mint egy harmadát kell csak megtennie. Tegyük fel, hogy a *Megarion 747*-hez hasonlóan a sávok olyan körökön vannak, melyek sugara 0,75 és 1,75 hüvelyk közé esik. Számoljuk ki, hogy átlagosan mennyi sávot kell elmozdultatnia a fejnek, ha két véletlenszerűen választott szektor közötti távolságot kell megtennie.

!! 2.2.4. feladat: A 2.1. példa végén azt mondtuk, hogy a maximális sávűrésű csökkenhető, ha a sávokat három tartományba soroljuk úgy, hogy a szektorok száma tartományonként elérhető. Most ne követeljük meg, hogy egyformák legyenek a tartományok, azaz a három tartományt elválasztó két határoló kör tetszőleges sugárú lehet, továbbá a tartományokba eső szektorok száma is változhat azzal a megkötéssel, hogy egy felületen a 8192 sávhoz tartozó bájtok száma összesen 1 gigabájt. Ekkor az öt pa-

ránméter (a tartományok közti két beosztás sugara) és a sávokra eső szektorok száma a három tartományban) milyen választása esetén lesz minimális egy tetszőleges sáv maximális sűrűsége?

2.3. A másodlagos tárolók hatékony használata

Az algoritmusokról szóló tanulmányok legtöbbszörben azt szokták feltenni, hogy az adatok a központi memóriában helyezkednek el és bármely két adat eléréséhez ugyannyi idő szükséges. Ezt a számítási modellt gyakran „RAM-modellnek” vagy másképpen véletlen hozzáféréstű számítási modellnek nevezik. Ezzel szemben mikor egy adatbázis-kezelő rendszert implementálunk, akkor azt kell feltennünk, hogy az adatok *nem* férnek el a központi memóriában. Emiatt a hatékony algoritmusok tervezésénél számításba kell venni a másodlagos, sőt esetleg a harmadlagos tárolók használatát is. Ebből következik, hogy a nagyon nagy mennyiségű adatokat feldolgozó legjobb algoritmusok gyakran különböznek az ugyanarra a problémára vonatkozó, de csak a központi memóriát használó legjobb algoritmusoktól.

Ebben a részben elsődlegesen a központi memória és a másodlagos tárolók közti kölcsönhatással fogunk foglalkozni. Különösképpen előnyös olyan algoritmusokat tervezni, melyek korlátozzák a lemezhozzáférések számát, még akkor is, ha az algoritmus során a központi memóriában az adatokon végzett műveletek nem a lehető legjobban használják ki a központi memóriát. Hasonló elvet alkalmazunk a memória-hierarchia minden szintjén. Egy központi memóriára vonatkozó algoritmuson is lehet javítani azzal, ha figyelembe vesszük a cache méretét, és az algoritmusunkat olyannak tervezzük, hogy a cache-be átmozgatót adatokat lehetőleg minél többször használjuk fel. Hasonlóan egy harmadlagos tárolót használó algoritmusnak is figyelembe kell vennie a másodlagos és harmadlagos memória között mozgó adatmennyiséget, és érdemes ezt a mennyiséget minimalizálni még annak az árán is, hogy a hierarchia alacsonyabb szintjein többelmunkát kell végezni.

2.3.1. A számítás I/O-modellje

Képzeljünk el egy számítógépet, melyen fut egy adatbázis-kezelő rendszer. A számítógép megpróbál kiszolgálni bizonyos számú felhasználót. A felhasználók az adatbázist különböző módon akarják elérni: lekérdezések és adatbázis-módosítások révén. Pillanatnyilag tegyük fel, hogy a számítógépnek egy processzora, egy lemezvezérlője és egy lemeze van. Maga az adatbázis sokkal nagyobb annál, hogy beférjen a központi memóriába. Jóllehet az adatbázis lényeges részzeit pufferelehetjük a központi memóriában, mégis az az általános, hogy az adatbázisnak minden egyes olyan darabját, amit egy felhasználó el akar érni, először vissza kell nyerni a lemeztől.

Fel fogjuk tenni, hogy a lemez *Megatron 747* típusú, 4 Kb-át a blokkméret, és az időtényezők megegyeznek a 2.3. példában meghatározott értékekkel. Speciálisan az átlagos blokkirási vagy olvasási idő körülbelül 15 milliszekundum. Mivel sok fel-

használó van, és minden felhasználó rendszeresen lemez I/O-igényeket ad ki, ezért a lemezvezérlőnek gyakorta az igények sorát kell kielégíteni. A kiszolgálásról kezletben azt tesszük fel, hogy mindig az elsőnek beérkezett igényt szolgálja ki először a lemezvezérlő. Ennek a stratégiának az a következménye, hogy egy adott felhasználó minden igénye véletlenszerűen fog tűnni (vagyis a lemez feje véletlenszerű helyen fog állni az igény előtt), még akkor is, ha ez a felhasználó csak egyetlen relációhoz tartozó blokkokat olvas, és ez a reláció ráadásul a lemez egyetlen cilindrán van elhelyezve. Ebben a fejezetben azt is megvizsgáljuk, hogyan lehet különböző módszerekkel a rendszer működését javítani. A *számítás I/O-modelljére* vonatkozó következő szabályt azonban végig igaznak tételezzük fel:

Az I/O-költség dominanciája: Ha egy blokkot kell mozgatni a lemez és a központi memória között, akkor az íráshoz és olvasáshoz szükséges idő sokkal nagyobb annál, amennyi a központi memóriában az adatkezeléssel szükséges. Így az algoritmusokhoz szükséges idő értékére jó becslést ad a blokkhozzáférések (írások és olvasások) száma, vagyis ezt az értéket kell minimalizálni.

2.4. példa: Tegyük fel, hogy az adatbázisunkban van egy R reláció, és egy lekérdezés az R reláció bizonyos k kulcsértékű sorát keresi. Látni fogjuk majd, hogy igen célszerű az R táblán egy indexet létrehozni, és ennek a segítségével azonosítani azt a lemezblokkot, amelyen a k kulcsértékű sor van. Ezzel szemben általában az már nem fontos, hogy az index azt is megmondja nekünk, hogy ez a sor a blokkon hol helyezkedik el.

Ennek az az oka, hogy ezt a 4 Kb-át méretű blokkot mindössze körülbelül 15 milliszekundum alatt lehet teljesen beolvasni, és ez alatt az idő alatt egy modern mikroprocesszor akár utasítások millióit is végre tudja hajtani. Ha már egyszer a blokk a memóriában van, akkor a k kulcsérték kereséséhez még a legbuzább lineáris keresési módszerrel is elegendő néhány ezer utasítás. Emiatt az a pluszidő, amivel ez a központi memóriában végrehajtott keresés jár, kevesebb, mint a blokkhozzáférési idő 1%-a, és így ezt nyugodtan el lehet hanyagolni. □

2.3.2. Adatok rendezése a másodlagos tárolóban

Lássunk egy bővebb példát arra, hogyan kell az algoritmusokat megváltoztatni a számítási költség I/O-modellje esetén. Tekintsük az adatok rendezését abban az esetben, amikor olyan nagyon sok az adat, hogy nem fér el a központi memóriában. Azzal kezdjük, hogy egy speciális rendezési problémát vezetünk be, és valamennyire részletezzük azt a gépet, amelyen a rendezés történik.

2.5. példa: Tegyük fel, hogy egy nagy R reláció 10 000 000 sort tartalmaz. Minden relációs sort egy rekord reprezentál, melynek több mezője is lehet, és a mezők között van egy *rendezési kulcs* mező. Ezt egyszerűen „kulcsmezőnek” fogjuk hívni, ha nem téveszthető össze másféle kulccsal. Egy rendezési algoritmus célja az, hogy rendezze a rekordokat a rendezési kulcsok értékeinek növekvő sorrendjében.

Egy rendezési kulcs lehet is, meg nem is „kulcs” az SQL *elsődleges kulcsának* szo-

káros értelmeben, ahol a rekordok garantáltan egyedi értékekkel rendelkeznek az elsődleges kulcsukban. Ha a rendezési kulcs értékei ismétlődhetnek, akkor az egyenlő rendezési kulccsal rendelkező rekordok bármelyik sorrendje elfogadható. Az egyszerűség kedvéért fel tesszük, hogy a rendezési kulcsok egyediek. Szintén az egyszerűség végett azt is feltételezzük, hogy a rekordok állandó hosszúságúak, nevezetesen minden rekord 100 bájttal hosszú. Így a teljes reláció egy gigabájttal foglal el.

Az a gép, amelyen rendezni szeretnénk, egy *Megarion 747* lemezzel rendelkezik, és olyan 50 megabájttal méretű memóriája is van, amely alkalmas arra, hogy a reláció blokkjait puffelje. A központi memória ténylegesen 64 Mbájttal, de a központi memória többi részét a rendszer használja.

Fel tesszük, hogy a lemez blokkjainak mérete 4096 bájttal. Így 40 darab 100 bájttal méretű sort vagy rekordot tudunk egy blokkba tenni, és még marad 96 bájttal, ami vagy bizonyos adminisztrációs célra használható, vagy nem használjuk fel egyáltalán semmire. A reláció így 250 000 blokkot foglal el. Az 50 Mbájttal (amely, mint tudjuk 50×2^{20} bájttal) méretű memóriában egyszerűen $50 \times 2^{20} / 2^{12}$, azaz 12 800 blokk fér el. \square

Ha az összes adat elfér a központi memóriában, akkor a számos jól ismert algoritmus bármelyike tökéletesen megfelel,⁸ így például használhatjuk a „Gyorsrendezés” (Quick-sort) algoritmusnak valamelyik variánsát, amit általában a leggyorsabbnak tartanak. Ezenfelül olyan stratégiát követhetünk, ahol csak a kulcsmezőket kell rendezni, pontosabban a kulcsmezőkhöz olyan mutatók is hozzá vannak csatolva, amelyek a megfelelő teljes rekordokra mutatnak. Csak ha a kulcs és mutató párok már sorrendben állnak, akkor használjuk fel a mutatókat arra, hogy behozzunk minden rekordot a neki megfelelő helyre.

Sajnos, ezek az elvek nem működnek túl jól, ha az adatok tárolásához másodlagos memóriára is szükség van. Amikor az adatok nagy része a másodlagos memóriában található, akkor jobban kedvelt a rendezésnek az a megközelítés, hogy valamilyen szabályos mintát követe minden blokkot csak néhányszor mozgatunk a központi és a másodlagos memória között. Ezek az algoritmusok gyakorta kisszámú *futamor* (pass) hajtások végére, egy futamor során minden rekordot egyszer olvasunk be a központi memóriába, és egyszer írunk ki a lemeze. A következő részben megnézzük egy ilyen algoritmust.

2.3.3. Az összfésülő rendezés (Merge-Sort)

Lehet, hogy az olvasó már találkozott az összfésülő rendezés nevű rendező algoritmusmal, mely azon az elven működik, hogy rendezett listákat nagyobb rendezett listákká fésül össze. A rendezett listák *összfésülése* úgy történik, hogy ismételtlen összehasonlítjuk minden lista legkisebb megmaradt kulcsértékét, és a kisebb kulcsú rekordot át tesszük az eredménybe, és ezt ismételtjük addig, amíg csak egy lista marad. Ekkor a választott rendezés szerint az eredmény mögé kell tenni a ki nem ürült lista maradék rekordjait, és ez adja az összes rekord kívánt sorrend szerint rendezett halmazát.

⁸ Lásd D. E. Knuth, *The Art of Computer Programming*, 3. kötet, Addison-Wesley, Reading MA, 1998, 2. kiadás, *Rendezés és keresés* című fejezetét. (Magyarul *A számítógép programozásának művésze*, Műszaki Könyvkiadó, Budapest, 1987.)

2.6. példa: Tegyük fel, hogy két rendezett listánk van, négy-négy rekorddal. Ahhoz, hogy még egyszerűbbé tegyük a dolgunkat, a rekordokat egyedül a kulcsukkal reprezentáljuk, a többi adatra nincs szükségünk, továbbá a kulcsok legyenek egész számok. Az egyik rendezett lista az (1, 3, 4, 9) és a másik a (2, 5, 7, 8). A 2.10. ábrán az összfésülés folyamatainak állapotait figyelhetjük meg.

Lépés	1. lista	2. lista	Eredmény
start	1, 3, 4, 9	2, 5, 7, 8	semmi
1)	3, 4, 9	2, 5, 7, 8	1
2)	3, 4, 9	5, 7, 8	1, 2
3)	4, 9	5, 7, 8	1, 2, 3
4)	9	5, 7, 8	1, 2, 3, 4
5)	9	7, 8	1, 2, 3, 4, 5
6)	9	8	1, 2, 3, 4, 5, 7
7)	9	semmi	1, 2, 3, 4, 5, 7, 8
8)	semmi	semmi	1, 2, 3, 4, 5, 7, 8, 9

2.10. ábra. Két rendezett listának összfésülése egy rendezett listává

Az első lépésben a két listának a sorban elsőként álló elemeit, azaz az 1-et és a 2-t hasonlítjuk össze. Mivel $1 < 2$, ezért az 1 értéket el távolítjuk az első listából és betesszük az eredménybe, ez lesz az eredmény első eleme. A 2. lépésben maradék listák legkisebb elemeit, azaz most a 3-at és a 2-t hasonlítjuk össze; a 2 nyer, így aztán őt tesszük be az eredménybe. Az összfésülés a 7. lépésig folytatódik, mikor is a második lista kiürül. Ekkor az első lista maradékát, ami most csak egyetlen elemből áll, hozzácsapjuk az eredményhez, és ezzel kész is az összfésülés. Vegyük észre, hogy az eredmény úgy van rendezve, ahogy lennie kell, mivel, minden lépésben a maradék elemek közül a legkisebbet választottuk. \square

A központi memóriában végrehajtott összfésüléshez szükséges idő a listák hosszának összegében lineáris. Ennek az a magyarázata, hogy az adott listák rendezettek, így csak a két lista mindenkor első elemét köztölt lehet a legkisebb ki nem választott elem, és az összehasonlítás konstans időt vesz igénybe. Az összfésülő rendezés klasszikus algoritmus a rekurzív módon rendez, és ha n elemet kell rendezni, akkor ezt $\log_2 n$ fázisban teszi, ahogy ez a következőkben látható.

Indukciós alap: Ha egy lista egy elemet tartalmaz, akkor semmit sem kell tenni, mivel ez már így is egy rendezett lista.

Indukció: Ha egy egynél több elemből álló listát kell rendezni, akkor tetszőleges módon osszuk fel a listát két egyenlő (vagy ha az eredeti lista páratlan hosszú, akkor majdnem egyenlő) hosszú részre. Rekurzívan rendezzük a két részlistát, majd az eredményül kapott rendezett listákat fésültük össze egy rendezett listává.

Ennek az algoritmusnak a részletes elemzése meglehetősen közsímet, és számunkra most nem is túl lényeges. Röviden összefoglalva, ha $T(n)$ -nel jelöljük az n elem

rendezéséhez szükséges időt, akkor ez egy összegként írható fel. Az egyik tag az n időnek egy konstansszorosa (ami egyébként a lista szétvágásából és a rendezett listák összehajlesztéséből ered), a másik tag pedig az az idő, ami két $n/2$ méretű lista rendezéséhez szükséges. Így a következő egyenletet kapjuk: $T(n) = 2T(n/2) + an$, ahol a valamilyen konstans. Ennek a rekurzív függvényegyenletnek a megoldása $T(n) = O(n \log n)$, vagyis $n \log n$ kifejezéssel arányos.

2.3.4. Kétfázisú, többutas, összefésülő rendezés

Ahhoz, hogy a 2.5. példában megadott gépen egy relációt rendezzünk, nem az összefésülő algoritmust fogjuk használni, hanem annak egy változatát, melyet *kétfázisú, többutas, összefésülő rendezésnek* hívunk. Az adabázis-alkalmazások leggyakrabban ezt a rendező algoritmust szereti használni. Ez az algoritmus röviden a következőkből áll:

- *1. fázis:* Készítsünk az adatainkból központi memória méretű rendezett darabokat, vagyis minden rekord legyen része egy olyan rendezett listának, amely éppen befér a rendelkezésre álló központi memóriába. Így valahány, de már *rendezett részlistákat* kapunk, melyeket a következő fázisban összefésülünk.
- *2. fázis:* Az összes rendezett részlistát fésüljük össze egyetlen rendezett listává.

Az első megállapításunk ezzel kapcsolatban, hogy ha az adatok a másodlagos tárolón helyezkednek el, akkor nem akarjuk azzal indítani a rekurziót, mint az előbb, azaz nem egy, esetleg néhány rekordból indulunk ki. Ennek az az oka, ha a rendezésre váró rekordok kitöltik a memóriát, akkor az összefésülő rendezés nem olyan gyors, mint más algoritmusok. Tehát azzal kezdjük a rekurziót, hogy a teljes központi memóriát kitöljünk rekordokkal, és a gyorsrendezéssel vagy bármilyen alkalmas központi memóriás rendező algoritmussal rendezzük a rekordokat. Ezt aztán annyiszor ismételjük, amennyiszor szükséges:

1. Töltsük ki a teljes hozzáférhető központi memóriát a rendezésre szánt eredeti reláció blokkjaival.
2. Rendezzük a rekordokat a központi memóriában.
3. Írjuk ki a rendezett rekordokat a központi memóriából a másodlagos tároló új blokkjaiba. Ezzel egy rendezett részlistát kapunk.

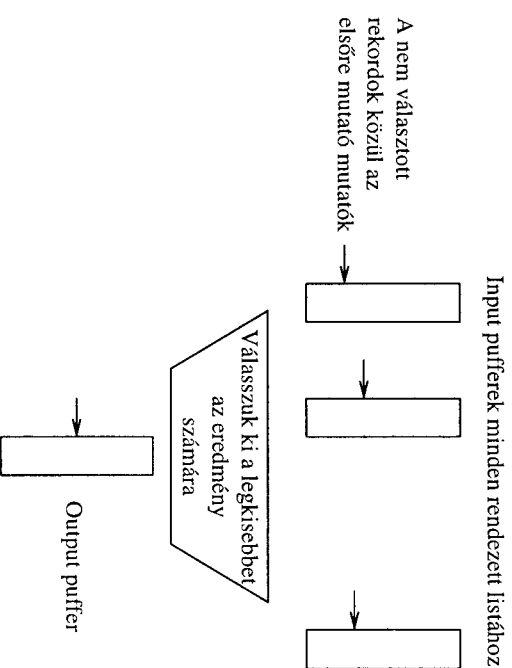
Az ily módon megadott *első fázis* végén az eredeti reláció összes rekordját egyszer olvastuk be a központi memóriába, minden rekord egy központi memória méretű rendezett részlistának lesz része, és ezeket a rendezett részlistákat kírítuk a lemezre.

2.7. példa: Tekintsük a 2.5. példában leírt relációt. Már meghatároztuk, hogy a 250 000 blokkból egyszerre 12 800 fér el a memóriában. Tehát 20-szor töltsük ki a memóriát, rendezzük a rekordokat a központi memóriában, és írunk rendezett részlistákat a lemezre. A 20 részlista közül az utolsó rövidebb, mint a többi, mivel csak 6800 blokkot foglal el, míg a többi 19 részlista egyaránt 12 800 blokk méretű.

Mennyi időt vesz ez a fázis igénybe? A 250 000 blokk mindegyikét pontosan egyszer olvassuk el, és 250 000 új blokkot írunk összesen. Ez félmillió lemez I/O-műveletet eredményez. Píllananyitlag fételezzük fel, hogy a blokkok véletlenszerűen helyezkednek el a lemezen. A 2.4. részben látni fogjuk, hogy ez egy olyan feltevés, amin jócskán lehet javítani. Mindenesetre a véletlenszerűégi feltevésünk mellett bármely blokk írása vagy olvasása egyformán 15 miliszekundumig tart. Ebből tehát az első fázisra vonatkozó I/O-idő összesen 7500 másodperc, ami 125 perc. Nem nehéz végiggondolni, hogy egy olyan processzorral, amely másodpercenként több tízmillió utasítást képes végrehajtani, a 10 000 000 rekord 20 rendezett részlistába sorolása az I/O-időnél sokkal kevesebb ideig tart. Tehát az első fázishoz szükséges összes idő 125 perc. □

Most térjünk rá arra, hogy hogyan fejezzük be a rendezést a rendezett részlisták összefésülésével. Megtehetünk, hogy páronként fésüljük össze, ahogy a klasszikus összefésülő rendezés esetében, de ekkor n rendezett részlista esetén $2 \log_2 n$ -szer kellene a memóriából ki-be olvasni az összes adatot. Például a 2.7. példa 20 rendezett részlistáját be kellene olvasni a másodlagos tárolóról, majd kírni ahhoz, hogy összefésüljük őket 10 rendezett részlistává, ezután újból egy teljes olvasást és írást kellene végrehajtanunk, hogy most már csak 5 rendezett lista maradjon, ebből aztán 4-nek az olvasása és írása után már csak 3 rendezett listánk marad és így tovább.

Jobbban járunk, ha minden egyes rendezett részlista első blokkját egy központi memóriapufferbe olvassuk be. Nagyon nagy relációk esetén az is előfordulhat, hogy az első fázisban túl sok rendezett részlistát kapunk, így aztán nincs annyi hely a központi memóriában, hogy minden rendezett listából egy blokkot be tudjunk olvasni. Ezzel a problémával a 2.3.5. részben fogunk foglalkozni. A 2.5. példához hasonló adatok esetében viszonylag kevés listát kapunk, a példában pont húszat, így minden listából egy blokkot könnyen be tudunk egyszerre tölteni a központi memóriába.



2.11. ábra. A központi memória szervezése a többutas összefésüléshez

Kijelölünk még egy puffert az eredményblokk számára is, melyet output puffernak nevezünk. Ez a puffer a teljes rendezett listának annyi első elemét fogja tartalmazni, amennyi csak befér. Kezdetben az output puffer üres. A pufferek elrendezését a 2.11. ábrán láthatjuk. A rendezett részlistákat a következő módon tudjuk egy olyan rendezett listába fésülni, amely már az összes rekordot tartalmazza.

1. Az összes lista megmaradt elemei közül válasszuk ki a legkisebb kulccsal rendelkező elemet. Mivel az összehasonlítás a központi memóriában történik, ezért elég egy lineáris keresést használnunk, amely a részlisták számával arányos gépi utasítást hajt végre. Természetesen használhatunk a legkisebb elem megtalálására jobb módszert is, például azt az eljárást, mely a „prioritásos sorban álláson”⁹ alapul. Ez utóbbi esetén a legkisebb elemet a részlisták számának logaritmusával arányos idő alatt lehet megtalálni.
2. Vegyük a legkisebb elemet az output blokk első elérhető helyére.
3. Ha az output blokk már megtelt, akkor írjuk ki a lemezre a tartalmát, és inicializáljuk újra a központi memóriának ugyanezt a puffert a következő output blokk tárolásához.
4. Ha az a blokk, amelyben a legkisebb elemet vetítük, ezáltal kiürül, akkor ugyanezen rendezett részlistából vegyük a következő blokkot, és olvassuk be ugyanebbe a pufferbe. Ha pedig már nincs több blokk, akkor hagyjuk ezt a puffert üresen, és ennek a listának az elemeit már nem kell figyelembe vennünk a továbbiakban, mikor is a maradék elemek közül a legkisebbet keressük.

Az első fázissal ellentétben a második fázisban nem lehet előre megmondani, hogy a blokkokat milyen sorrendben fogjuk beolvasni, mivel nem tudjuk megmondani, hogy az input blokk mikor ürül ki. Azt azonban megfigyelhetjük, hogy bármelyik rendezett lista rekordjait tartalmazó blokkot pontosan egyszer olvasunk be a lemezről. Emiatt a második fázisban összesen 250 000 blokkolvasást hajtunk végre, ugyanannyit, mint az első fázisban. Hasonlóan minden rekordot egyszer beteszünk egy output blokkba, és ezeket a blokkokat a lemezre küldjük. Tehát a második fázisban a blokkolvasók száma szintén 250 000. Mivel a második fázisban a központi memóriában végzett számítás ismét csak el lehet hanyagolni az I/O-költséghez viszonyítva, így arra következtethetünk, hogy a második fázis költsége szintén 125 perc, vagyis a teljes rendezés költsége 250 perc.

2.3.5. A többutas összefűlés kiterjesztése nagyobb relációkra

Az előbb leírt kétfázisú, többutas, összefűlítő rendezést nagyon nagy rekordhalmazok rendezésére is használhatjuk. Ahhoz, hogy lássuk, mekkora lehet ez a „nagyon nagy”, tételezzük fel a következőket:

⁹ Lásd Aho, A. V., J. D. Ullman *Foundations of Computer Science*. Computer Science Press, 1992.

Milyen nagyoknak kell egy blokknak lennie?

A *Megatron 747* lemezt használó algoritmusaink elemzése során azt tételítettük fel, hogy 4 Kb-ait mérjük minden blokk. Ezzel szemben indokolható az is, hogy ennél nagyobb blokkméret előnyösebb lenne. Emlékezzünk vissza a 2.3. példára, ahol kiszámoltuk, hogy körülbelül fél milliszekundum a 4 K méretű blokk átviteli ideje, és körülbelül 14 milliszekundum az átlagos keresési idő és rotációs késés együttesen. Ha megdupláznánk a blokkok méretét, akkor feleannyi lemez I/O-műveletre lenne szükség egy olyan algoritmus során, mint amilyen az itt leírt többutas, összefűlítő rendezés. Ezzel szemben a blokkhosszátérési időben az egyetlen változás az lenne, hogy az átviteli idő 1 milliszekundumra növekedne. Ezzel tehát a rendezést hozzávetőleg feleakkora idő alatt tudnánk végrehajtani, mint az eredeti blokkméret esetén.

Ha most ismét megdupláznánk a blokkméretet 16 Kb-átra, akkor az átviteli idő 2 milliszekundumra emelkedne, míg 64 K blokkméret esetén 8 milliszekundum lenne. Ekkor már az átlagos blokkhosszátérési idő 22 milliszekundum lenne, de csak 62 500 blokkhosszátérésre lenne szükségünk ahhoz, hogy a rendezést 10-szeresére gyorsítsuk.

Külfülbőlő okai vannak annak, hogy a fentiek ellenére a blokkméret általában meglehetősen kicsi. Először is nem tudjuk hatékonyan használni az olyan blokkokat, amelyek több sávon helyezkednek el. Másodsor a kis relációk egy blokknak csak a töredékét foglalják el, így sok elpazarolt hely lenne a lemezen. Aztán a másodlagos adattároló szervezésére léteznek olyan adatstruktúrák is, melyek jobban szeretik, ha az adatok sok blokkba vannak szétosztva, és ezért ezek kevésbé jól működnek, ha a blokkméret túl nagy. Tulajdonképpen a 2.3.5. részben látni fogjuk, hogy minél nagyobb a blokkméret, annál kevesebb rekordot tudunk az itt leírt kétfázisú, többutas módszerrel rendezni. Mindazonáltal a gépek sebességének és a lemezek kapacitásának növekedésével megfigyelhető a blokkok méretének megnövekedése irányuló tendencia.

1. A blokkok mérete B bájtt.
2. M bájtt használható a központi memóriában a blokkok pufferelésére.
3. A rekordok mérete R bájtt.

A központi memóriában ezért összesen M/B számú puffer képezhető. A második fázisban ezek közül a pufferek közül egy kivételével mindegyik hozzá van rendelve valamelyik rendezett részlistához. A fennmaradó puffer pedig az output blokkhoz kell. Emiatt $(M/B) - 1$ rendezett részlistát lehet készíteni ebben a fázisban. Ez a szám egyezik azzal a számmal, mely megmondja, hogy hányszor kell feltölteni a központi memóriát a rendezni kívánt rekordokkal. A központi memória minden feltöltése során összesen M/R rekordot rendezünk. Így az összes rekord, amit rendezni tudunk $(M/R)((M/B) - 1)$, azaz körülbelül M^2/RB rekord.

2.8. példa: Ha a 2.5. példában vázolt paramétereket használjuk, akkor $M = 50\,000\,000$, $B = 4096$ és $R = 100$. Tehát összesen $M^2/RB = 6,1$ milliárd rekordot tudunk rendezni, amely összesen egy terabájt hajtuzedét foglalja el.

Vegyük észre, hogy egy ekkora méretű reláció nem is fér el egy *Megatron 747* lemezen, sőt még gyakorlati szempontból elfogadható számú lemezen sem. Valószínűleg a rekordok tárolására egy harmadlagos tárolóeszközt kellene használni, és a harmadlagos tárolóról kellene a rekordokat a lemezeire vagy lemezekre átmozgassunk egy többitas, összefésült rendezéshez hasonló stratégiával, csak most a harmadlagos és másodlagos tároló játsza azt a szerepet, amit előzőleg a másodlagos tároló és a központi memória játszott. \square

Ha még több rekordot kell rendeznünk, akkor kiegészítjük egy harmadik menettel. A kétfázisú, többitas, összefésült rendezéssel M^2/RB rekordcsoportokat rendezünk rendezett részlistákká. Ezután egy harmadik fázisban ezek közül a listák közül legfeljebb $(M/B) - 1$ listát összefésülünk egy többitas összefésüléssel.

A harmadik fázisban nagyjából M^3/RB^2 rekord rendezését teszi lehetővé, melyek M^3/B^2 blokkot foglalnak el. A 2.5. példa paramétereit használva körülbelül 75 trillió rekordot kapunk, melyek összesen 7500 petabájtot foglalnak el. Ez akkora mennyiség, amiről ma még hallani sem lehet. Mivel még a kétfázisú, többitas, összefésült rendezésre adott 0,61 terabájt limit is valószínűtlenül nagyobb, mint amekkora mennyiséget a másodlagos tárolóval kell kezelnünk, ezért azt mondhatjuk, hogy a többitas, összefésült rendezésnek a kétfázisú változata valószínűleg minden gyakorlati célhoz megfelel.

2.3.6. Feladatok

2.3.1. feladat: A 2.5. példa relációját mennyi idő alatt lehet a kétfázisú, többitas, összefésült rendezéssel rendezni, ha a *Megatron 747* lemezt lecseréljük a 2.2.1. feladatban leírt *Megatron 777* lemeze, de egyébként a gépnek és az adatoknak minden más jellemzője változatlanul maradt?

2.3.2. feladat: Tegyük fel, hogy a kétfázisú, többitas, összefésült rendezést akarjunk használni a 2.5. példában megadott gépre és R relációra, de a paramétereken változtatunk. Számoljuk ki, hogy a rendezéshez mennyi lemez I/O-műveletre van szükség, ha az R reláció és/vagy a gép jellemzőit a következőkre változtatjuk:

- * a) Az R relációban a sorok számát megduplázzuk (minden más változatlan marad).
- b) A sorok hosszát megduplázzuk, azaz 200 bájtira változtatjuk (minden más meg-egyezik a 2.5. példában szereplő értékekkel).
- * c) A blokkok méretét duplázzuk meg, azaz 8192 bájtira növeljük (minden más paraméter, mint ebben a feladatban végig, változatlan marad).
- d) A hozzáférhető memória méretét duplázzuk meg, azaz 100 megabájtira növeljük.

! 2.3.3. feladat: Tegyük fel, hogy a 2.5. példa R relációja olyan nagyra nő, hogy már annyi sorral rendelkezik, amennyit maximálisan rendezni lehet a kétfázisú, többitas, összefésült rendezéssel a példában szereplő gépen. Azt is tegyük fel, hogy a lemez is akkora, hogy az R relációt be tudja fogadni. Mennyi ideig tart az R rendezése akkor, ha a lemez, gép és az R reláció összes többi jellemzőjét változtatlanul hagyjuk?

* **2.3.4. feladat:** Tekintsük újra a 2.5. példa R relációját, de most azt tételezzük fel, hogy a rendezési kulcs (ami szokásos értelemben egy kulcs, azaz egyértelműen azonosítja a rekordokat) alapján rendezve tároljuk. Továbbá még azt is tegyük fel, hogy az R relációt olyan blokkok sorozatán tároljuk, melyek elhelyezkedését pontosan ismerjük, azaz bármilyen i esetén az R_i blokkját egyetlen lemez I/O-művelettel lehet elérni. Ha adott egy K kulcsérték, akkor az ezzel a kulccsal rendelkező sort a szabványos bináris kereséssel találhatjuk meg. Maximálisan hány lemez I/O-műveletre van szükség ahhoz, hogy megtaláljuk a K kulccsal rendelkező sort?

!! **2.3.5. feladat:** Tegyük fel, hogy ugyanaz a helyzet, mint a 2.3.4. feladatban, de most 10 előre adott kulcsértéket keresünk. Maximálisan hány lemez I/O-művelet szükséges ahhoz, hogy megtaláljuk mind a 10 sort?

* **2.3.6. feladat:** Tegyük fel, hogy van egy relációnk n sorral és mindegyik sor R bájt hosszú. Adott továbbá egy gép olyan M méretű központi memóriával és B méretű lemezblokkokkal, amely éppen elég ahhoz, hogy az n sort a kétfázisú, többitas, összefésült rendezéssel rendezni lehessen. Hogyan változik a maximális n , ha a paramétereken a következő változtatásokat tesszük?

- a) Megduplázzuk B -t.
- b) Megduplázzuk R -et.
- c) Megduplázzuk M -et.

! **2.3.7. feladat:** Ismételjük meg a 2.3.6. feladatot, de most olyan paraméterekkel, amelyek mellett még lehetséges a háromfázisú, többitas, összefésült rendezés.

* **! 2.3.8. feladat:** Határozzuk meg, hogy k fázisú (k egy egész szám), többitas, összefésült rendezés esetén hány rekordot lehet maximálisan rendezni. Az eredményt a 2.3.6. feladatban használt R , M , B és k függvényében adjuk meg.

2.4. A másodlagos tároló hozzáférési idejének javítása

A 2.3.4. rész elemzése során feltettük, hogy az adatokat egyetlen lemezen tároljuk, és a blokkokat véletlenszerűen választjuk a lemez lehetséges helyei közül. Ez a feltevésünk egy olyan rendszer esetében alkalmazható, amely nagyon sok, de kis lekérdezést hajt végre szimultán módon. Ezzel szemben, ha a rendszernek nincs más feladata,

mint hogy egy nagy relációt rendezzen, akkor jelentős időt meg tudunk takarítani, ha jobban megfontoljuk, hogy a rendezésbe bevont blokkokat hová helyezzük. Ezáltal ki tudjuk használni azt is, hogy hogyan működik a lemez. Tulajdonképpen még az előző esetben is, vagyis mikor a rendszeren a legnagyobb terhet sok egymással kapcsolatlan nem álló lekérdezés okozza azáltal, hogy „véletlen” blokkokat kell elérni a lemezen, még akkor is sok mindent tehetünk annak érdekében, hogy a lekérdezések sokkal gyorsabban fussanak, és/vagy a rendszer megengedje, hogy egyszerre több lekérdezést is végre lehessen hajtani (azaz növeljük az „*tejtímsímenyt*”). Ezek közül a stratégiák közül ebben a részben a következőket tekintjük át:

- Helyezzük ugyanarra a cilindere azokat a blokkokat, amelyeket egyszerre kell elérni. Ezzel gyakran megússzuk a keresési időt és lehet, hogy még a rotációs késési is.
- Ahelyett, hogy egy nagy lemezt használhánk, inkább osszuk szét az adatunkat több kisebb lemezre. Mivel a több fejszerelvény a blokkokat egymástól függetlenül keresheti, ezért az egysegényi időre eső blokkhozzáférések száma ezzel megnövekedhet.
- „Tükrözzük” a lemezt: Készítsünk kettő vagy több másolatot a lemez adatáról. Amellett, hogy lemezhiba esetén mentésünk marad az adatunkról, ez a stratégia arra is jó, hogy egyszerre több blokkhoz tudunk hozzáférni, ahogy az előző pontban, azaz mikor az adatokat több lemezre osztottuk szét.
- Használjunk valamilyen lemezütemező algoritmust. Ez lehet az operációs rendszerben, az adatbázis-kezelő rendszerben vagy a lemezvezérlőben, és ez adja meg, hogy ha több blokkolvasási vagy írási igény érkezik be, akkor ezeket milyen sorrendben kell végrehajtani.
- Hozzuk be előre azokat a blokkokat a központi memóriába, amelyekről előre látható, hogy a későbbiekben használni fogjuk őket.

A vizsgálatunk során hangsúlyozni fogjuk, hogy javulás akkor várható, ha a rendszer legalább egy bizonyos ideig egy speciális feladattal foglalkozik, amely lehet például a 2.5. példában bevezetett rendezési művelet. Van még legalább két másik nézőpont is, amelynek segítségével mérhetjük a rendszer működését, illetve a másodlagos memória használatát:

1. Mi történik abban az esetben, mikor nagyon sok folyamatot kell a rendszernek egy időben támmogatnia? Például egy repülőgépes helyfoglalási rendszernek egyszerre nagyon sok ügynököt kell kiszolgálnia, akik a járatokról lekérdezéseket tesznek fel, vagy helyet foglalnak le.
2. Mit csináljunk, ha előre megadott fix költségből kell kiegészítenünk egy számítógépes rendszert, illetve mit csináljunk, ha mindenféle lekérdezést kell végrehajtanunk egy olyan rendszeren, amely már adott és nem könnyen változtatható meg?

Ezeket a kérdéseket a 2.4.6. részben nézzük majd meg, miután megvizsgáltuk a fenti lehetőségeket.

2.4.1. Az adatok cilinderes szervezése

Mivel a keresési idő általában az átlagos blokkhozzáférési időnek a felét teszi ki, ezért sok alkalmazásban lehet értelme annak, hogy azokat az adatokat, amelyeket egyszerre kell majd elérni, például a relációkat, egy cilinderen tároljuk. Ha nem lenne elég hely, akkor néhány szomszédos cilindert használhatunk fel erre a célra.

Valójában, ha egyszerre beolvassuk a sávon vagy a cilinderen tárolt összes blokkot, akkor lehet, hogy csak az első keresési idő (ami a cilindere álláshoz szükséges) és az első rotációs idő (ami az első blokknak a fej alá kerüléséhez kell) marad meg, minden más késleltetéstől el lehet tekinteni. Ebben az esetben az adatok lemezről olvasásánál, írásánál meg tudjuk közelíteni az elméletileg elérhető átviteli gyorsaságot.

2.9. példa: Nézzük meg újra a 2.3.4. részben leírt kétfázisú, többutas, összefésülő rendezés működését. A 2.3. példában meghatároztuk az átlagos blokkátviteli időt, keresési időt és rotációs késést. Ezekre rendre a 0,5 milliszekundum, 6,5 milliszekundum és 7,8 milliszekundum értékeket kaptunk *Megatron 747* lemez esetén. Azt is megállapítottuk, hogy az egy gigabájtot elfoglaló 10 000 000 rekord rendezése 250 percet vesz igénybe. Ezt az időt négy nagy műveletre osztottuk: egy olvasás és egy írás tartozott az algoritmus két fázisának mindegyikéhez.

Nézzük meg, hogy az adatok cilinderes szervezése tud-e javítani ezeken a műveleti időkön. Az első művelet az volt, hogy az eredeti rekordokat beolvastuk a központi memóriába. A 2.7. példa szerint 20 alkalommal töltöttük fel a központi memóriát, meghozzá minden esetben 12 800 blokkal.

Az eredeti adatokat egymás utáni cilindereken is tárolhatjuk. A *Megatron 747* lemez 8192 cilindereinek mindegyike körülbelül egy megabájtot tud tárolni; valójában ez az érték csak egy átlag, mivel a belső sávok kevesebbet, a külső sávok többet tudnak tárolni, de az egyszerűség kedvéért fellesszük, hogy minden sáv és minden cylinder kapacitása az átlaggal egyezik meg. Így tehát a kezdeti adatainkat 1000 cilinderen tudjuk tárolni, amiből 50 cilindert olvasunk be a központi memóriába. Emiatt egy cilindert egy keresési idővel olvashatunk be. Még azt sem kell megvárunk, hogy a cilindernek egy speciális blokkja kerüljön a fej alá, mivel ebben a fázisban a rekordok sorrendje nem játszik szerepet. Összesen 49 alkalommal kell a fejet a szomszédos cilindere mozgatnunk. A 2.3. példa paramétereit mellett kiszámoltuk, hogy egy sávnyi elmozduláshoz csak egy milliszekundum szükséges. Ekkor tehát a központi memória kitértéséhez szükséges összes idő:

1. Az átlagos kereséshez 6,5 milliszekundum kell.
2. A 49 egycilinderes mozgáshoz 49 milliszekundum szükséges.
3. A 12 800 blokk beolvasásához 6,4 másodperc kell.

Az utolsó érték kivételével a többi el lehet hanyagolni. Mivel 20-szor töltjük ki a memóriát, így az első fázishoz tartozó teljes beolvasás idő körülbelül 2,15 perc. Ha hasonlítunk össze ezt az időt azzal az egy órával, amint a 2.7. példában kaptunk az első fázis olvasási részére, igaz, akkor azt tételeztük fel, hogy a blokkok eloslása a lemezen véletlenszerű. Az első fázis írási részénél a rekordok 20 rendezett részlistájának

tárolására is hasonlóan használhatunk szomszédos cilindreket. Ezeket a rendezett listákat szintén ki lehet írni másik 1000 cilindrere. Ehhez ugyanolyan fejmozgatásokra van szükség, mint az olvasás esetén: egy véletlen keresés és 49 egycylinderes keresés szükséges a 20 lista mindegyikénél. Így az első fázis írási ideje szintén körülbelül 2,15 perc, azaz 4,3 percet kapunk a teljes első fázisra, szemben a véletlen eloszlású blokkok esetén kiszámolt 125 perccel.

Másrészt a cilinderes tárolás nem segít a rendezés második fázisában. Emlékezzünk vissza arra, hogy a második fázisban a 20 rendezett részlista elejétől olvassuk be a blokkokat, méghozzá olyan sorrendben, amit az adatok határoznak meg, illetve az, hogy melyik lista aktuális blokkja kerül ki legközelebb. Hasonlóan a végleges rendezett listát tartalmazó output blokkokból is időnként egyet ki kell írniuk, így időnként egy blokk-írás tartkítja a blokkolvasások sorozatát. Tehát a második fázis változatlanul körülbelül 125 percig tart. Következésképpen a rendezési időt majdnem a felére tudtuk csökkenteni, de csupán a cilindreket átgondolt használatával nem lehet jobb eredmény elérni. □

2.4.2. Több lemez használata

A rendszertünk sebességén gyakran javítani tudunk azzal, hogy ahelyett, hogy egy lemezt használhánk, sok egymáshoz kapcsolódó fejjel, inkább több lemezt használjunk független fejekkel. Egy ilyen elrendezést mutatott be a 2.6. ábra, ahol három lemez kapcsolódott egyetlen vezérlőhöz. Ha a lemezvezérlő, az adatsín és a központi memória az átvitt adatokat nagyon gyorsan tudja kezelni, akkor nagyjából az lesz a hatás, mintha a lemezolvasási és -írási sebességeket osztanánk a lemezek számával. Egy példán mutatjuk be az így keletkezett változást.

2.10. példa: A *Megarion 737* lemez összes jellemzője megegyezik a 2.1. és 2.3. példában megadott *Megarion 747* lemezzel, de a 737 esetében csak egy tányér van két felülettel. Így módon egy *Megarion 737* lemez 2 gigabájtot képes tárolni. Tegyük fel, hogy egy *Megarion 747* lemezünket lecseréljük négy darab *Megarion 737* lemezre. Tekintsük át, hogy a kétfázisú, többutas, összefésülő rendezést hogyan lehet végrehajtani.

Először szétoszadjuk az adott rekordokat a négy lemez között. Az adatok minden egyes lemezen 1000 szomszédos cylindert fognak elfoglalni. Mikor az első fázis során meg akarjuk tölteni a lemeztől a központi memóriát, akkor minden lemeztől a központi memória 1/4 részét töltjük ki. A 2.9. példa alapján megint kihasználhatjuk, hogy a keresési idő és a rotációs késés lényegében tart a nullához. Viszont ahhoz, hogy a memória negyedét elég blokkal töltsük ki, be kell olvasnunk 3200 blokkot a lemeztől, amit körülbelül 1600 miliszekundum, azaz 1,6 másodperc alatt tehetünk meg. Ha rendszer ilyen sebességgel tudja a négy lemeztől egyszerre jövő adatokat is kezelni, akkor a központi memória 50 megabájt méretű részét 1,6 másodperc alatt tudjuk megöltöni, szemben az egy lemez használatánál kiszámolt 6,4 másodperccel.

Hasonlóan, mikor az első fázisban sor került arra, hogy az adatokat ki kell írni a központi memóriából, akkor a rendezett részlistákat is szétoszthatjuk a négy lemez között, amiből minden lemezen körülbelül 50 szomszédos cylindert foglalunk le. Tehát az első fázis írási részének sebességét is négyeszeresére tudjuk gyorsítani. Így a

teljes első fázis egy percig tart szemben a 2.4.1. részben leírt, csak cilinderes javítással elért 4 perccel, nem beszélve az eredeti, véletlenszerűégi feltevésnél kapott 125 percről.

Most vegyük szemügyre a kétfázisú, többutas, összefésülő rendezés második fázisát. A kiülbözőző listák elejétől látszólag véletlenszerű, adatfüggő sorrendben kell meg beolvasni a blokkokat. A második fázis algoritmusának maga azt követeli meg, hogy mind a 20 listának megfelelő blokkteljesen be legyen töltve a központi memória-ba, ahhoz, hogy a 20 részlista megmaradó elemei közül a legkisebbet kiválasszassuk. Emiatt nem tudjuk kihasználni, hogy 4 lemezünk van. Minden esetben, mikor kimerül egy blokk, várnunk kell, amíg egy új blokkot teljesen beolvassunk ugyanarról a listáról, hogy ezzel helyettesítsük az előzőt. Így egyszerre mindig csak egy lemezt használunk.

Ha azonban ügyesebben írjuk meg az algoritmus kódját, akkor a 20 legkisebb elem összehasonlítását már abban a pillanatban folytathatjuk, mikor az új blokk első eleme megjelenik a központi memóriában.¹⁰ Ha így teszünk, akkor egyszerre több lista is betölthető a blokkjait a központi memóriába. Amikor ezek a blokkok kiülbözőző lemezekben helyezkednek el, akkor minden blokkolvasást egy időben tudunk elvégezni, és ezáltal a 2. fázis olvasási részének sebességét egy potenciális 4-szeres faktorral tudjuk növelni. A blokkolvasások véletlen sorrendje azért továbbra is korlátoz bennünket; ugyanis, ha a következő két blokk történetesen ugyanazon a lemezen helyezkedik el, akkor az egyiknek meg kell várni a másikat, és a teljes központi memória áll adóig, amíg legalább a második blokk eleje meg nem érkezik a központi memóriába.

A 2. fázis írási részét könnyebb felgyorsítani, ugyanis használhatunk négy output puffert, és sorba mindegyiket kitöltjük. Ha valamelyik puffert megtelít, akkor rögtön kinyit az egyik meghatározott lemezre úgy, hogy a cilindreket sorba töltjük fel. Ezáltal a pufferek közötti egyet mindig fel lehet tölteni, amíg a többi hármat éppen kinyitjuk.

Mindazonáltal nem lehet gyorsabban kiírni a teljes rendezett listát annál, ahogy a 20 köztes listáról az adatokat beolvastuk. A fentiekben láttuk, hogy nem lehet azt elérni, hogy mind a négy lemez egyszerre és folyamatosan hasznos munkát végezzen, ezért a második fázis valószínűleg csak 2-3-szorosára gyorsítható, de még a kétszeres tényező is egy órái takarít meg nekünk. Összegezve: a cilinderes szervezés a 4 lemezes adattárolással együtt a rendezéses példánk idejét mindkét fázisban csökkenteni tudja 125 percről 1 percre az első fázis esetében, és 1 órára, a második fázis esetében. □

2.4.3. Lemezek tükrözése

Előfordulhatnak olyan helyzetek, amikor értelmesebb tűnik, hogy két vagy több lemez ugyanazoknak az adatoknak a másolatait tartalmazza. Ekkor azt mondjuk, hogy ezek a lemezek egymás *tükrözései*. Egyik fontos készletésünk lehet, hogy ilyen módon bár-

¹⁰ Azt azonban hangsúlyozni kell, hogy ez a megközelítés rendkívül finom implementálást követel meg, és csak akkor kell ezzel megpróbálkozni, ha fontos előny származik belőle. Jelenlét, ugyanis a kockázata annak, hogy ha nem vagyunk elég előrelátóak, akkor egy rekord már azelőtt próbálunk elolvasni, hogy az megérkezett volna a központi memóriába.

melyik lemez fejének meghibásodását az adataink túlélik, mivel a meghibásodott lemez egyik tükrözéséről még mindig beolvashatók. Azoknál a rendszereknél, amelyek két úgy terveztek, hogy támogassák a megbízható működést, gyakran lemezpárokat használnak, melyek egymás tükröképei.

Ettől túlmenően a lemezek tükrözése is felgyorsíthatja az adatelelést. Emlékezzünk vissza, hogy a 2.10. példában a kétfázisú, összetévesztő rendezés 2. fázisának elemzése során észrevettük, hogyha nagyon előrelátóak vagyunk az időzítésekkel, akkor azt is el tudnánk érni, hogy a négy különböző rendezett listáról négy blokkot töltsünk fel egyszerre, ha az előző blokkjuk már kimerült. Azt viszont nem tudjuk előre kiválasztani, hogy melyik négy listának lesz szüksége új blokkra. Ha nagyon szerencsétlenek vagyunk, akkor például azt találjuk, hogy az első két lista ugyanazon a lemezen van, vagy az első három listából kettő van ugyanazon a lemezen.

Ha hajlandóak vagyunk arra, hogy egy nagy lemeztől 4 másolatot készítsünk, és ezáltal pazaroljuk a lemeztérületet, akkor cserébe a rendszer mindig garantáltan viszszanyerhet négy blokkot egy időben. Vagyis, ha mindegy, hogy melyik négy blokkra van szükség, akkor mindegyiket hozzáférdejük a négy lemez valamelyikéhez, és arról a lemeztől olvastatjuk be a blokkot.

Általánostva, ha n másolatot készítünk egy lemeztől, akkor tetszőleges n blokkot olvashatunk párhuzamosan. Ha n blokknál kevesebbet kell egyszerre olvasni, akkor gyakra növelni tudjuk a sebességet azáltal, hogy megfontoltan választjuk ki azt a lemezt, amelyről olvassuk. Ugyanis vehetjük a rendelkezésre álló lemezek közül azt, amelyiknek a feje a legközelebb esik ahhoz a cilinderrhez, amelyet olvasni akarunk.

A tükrözött lemezek használata az egylemez használatához képest nem gyorsítja fel az írást, de szerencsére nem is lassítja le. Ugyanis mikor egy blokkot kell írunk, akkor igaz, hogy minden olyan lemeze ki kell ezt írni, amelyiken másolat található, de mivel az írás párhuzamosan történik, így az eltele idő körülbelül ugyanakkora, mintha egy lemezt írunk volna. Valójában a különböző tükrözött lemezek esetében picit eltérhet az íráshoz szükséges idő, mert nem lehetünk biztosak abban, hogy a forgásuk pontosan szinkronizáltak történiük. Így lehet, hogy az egyik lemez feje éppen leköszik egy blokkot, míg a másik lemez feje lehet, hogy éppen most készül elhaladni ugyanazon blokk fölött. Azonban ezek a rotációs késési időkre vonatkozó eltérések átlagosan kiegyenlítődnék, és ha a 2.4.1. rész cilinderalapú stratégiáját használjuk, akkor a rotációs késés teljesen elhanyagolható.

2.4.4. A lemez ütemezése és a lift algoritmus

Bizonyos esetekben a lemezhozzáférést más hatékony módon is lehet gyorsítani. Például azzal, hogy a lemezvezérlővel választjuk ki, hogy több igény közül melyiket hajtsa végre először. Igaz, hogy ezt a lehetőséget nem nagyon tudjuk kihasználni abban az esetben, mikor a rendszernek bizonyos adott sorrendben kell a lemezblokkokat olvasnia vagy írnia, mint például az összefésült rendezésünk egyes részében. Abban az esetben viszont, amikor a rendszer sok kis folyamatot támogat, melyek mindegyike néhány blokkot akar csak elérni, akkor növelni lehet a teljesítményt azzal, hogy kiválasztjuk, hogy mely folyamattípus kapja meg az elsőiséget.

Egy egyszerű és hatékony ütemező módszer nagyon sok blokkhozzáférés esetére az úgynevezett *lift algoritmus*. Képzeli el a lemez fejt, amint pásztázza a lemezt, a legbelső cilindertől a legkülsőig, aztán újból vissza. Úgy is gondolhatunk a lemez fejére, mint egy liftre, amely függőlegesen mozog egy épület aljától a tetejéig, aztán vissza. Amint a fej egy cilindert halad keresztül, megáll, ha egy vagy több igény vonatkozik olyan blokkokra, melyek ezen a cilindert találhatók. Minden ilyen blokkot az igény szerint olvassuk vagy írunk. A fejek ezután ugyanabban az irányban haladnak tovább, mint eddig. Egészen addig haladnak, amíg el nem érnek a következő olyan cilinderrhez, amelyen olyan blokkok vannak, amiket el akarunk érni. Amikor a fejek olyan helyre érnek, hogy az eddig mozgási irányjukban már nincs több keresett blokk előtűk, akkor irányt váltanak, és az ellenkező irányba folytatják a keresést.

2.11. példa: Tegyük fel, hogy a *Megatron 747* lemezt akarjuk ütemezni. Emlékezzünk vissza, hogy a lemez átlagos keresési ideje, rotációs késése és átviteli ideje rendre 6,5, 7,8 és 0,5. Ebben a példában minden idő milliszekundumban értendő. Tegyük fel, hogy valamikor olyan blokkokat akarunk elérni, melyek a 1000., 3000. és a 7000. cilindert találhatók. A fejek az 1000. cilindert helyezkednek el. Továbbá később még befut három blokkhozzáférési igény, ahogy ezt a 2.12. ábra szemlélteti. Például egy 2000. cilindert találhatók blokk elérését igényeljük a 20. milliszekundumban.

Fellesszük azt is, hogy minden blokkhozzáférés esetén 0,5 jut az átvételre és 7,8 az átlagos rotációs késésre, vagyis a blokkhozzáféréshez összesen 8,3 plusz annyi milliszekundum kell, amennyi a keresési idő. Ezt a keresési időt a *Megatron 747* lemeze a 2.3. példában megadott szabállyal lehet kiszámolni, azaz a sávok számához hozzáadunk egyet, és osztjuk 500-zal. Nézzük meg, hogy mi történik, ha a lift algoritmussal ütemezzük a végrehajtást. Az 1000. cilindert vonatkózó első igényhez nem kell keresési idő, mivel már ott vannak a fejek. Tehát az első igény kielégítéséhez szükséges idő 8,3. Mivel a 2000. cilindert vonatkózó igény ekkor még nem jött be, így a fejeket továbbvisszük a 3000. cilindert, amely a következő igényelt megállás a magasabb sorszámú sávok irányába. Az 1000.-tól a 3000. cilindert tartó mozgás keresési ideje 5 milliszekundum, így 13,3-kor érkezünk oda, és a blokkelérést 8,3 milliszekundum múlva fejezzük be. Tehát a második eléréssel is végzünk 21,6-kor. Ekkorra befut az igény a 2000. cilindert, de mi már túlhaladtunk ezen a cilindert 11,3-kor, és nem is jövünk vissza a következő menetre.

Tehát továbbmegyünk a 7000. cilindert. A keresési idő 9, a rotációs és átviteli idő 8,3, tehát a harmadik eléréssel 38,9-kor végzünk. Most már megérkezett a 8000. cilindert vonatkózó igény is, tehát továbbmegyünk ugyanebbe az irányba. A keresési idő most 3 milliszekundum, így az elérést 38,9 + 3 + 8,3 = 50,2-kor fejezzük be. Ekkor már az 5000. cilindert vonatkózó igényt is meglették, így ez és a 2000. cilindert máradt még hátra. Így visszafelé, azaz a lemez belseje felé indulunk, hogy ezt a két igényt is kielégítsük.

Hasonlítsuk össze a lift algoritmus végrehajtását egy sokkal naivabb megközelítéssel, például azzal, hogy mindig az elsőre bejövő igényt szolgáljuk ki (first-come-first-served). Az első három igényt pontosan ugyanúgy hajtuk végre, mint az előbb, feltéve, hogy az első három igény beérkezési sorrendje 1000, 3000, 7000. Einnél a pontnál

viszont vissza kell menni a 2000. cilinderré, mivel ez volt a negyediknek beérkező igény. Az ehhez az igényhez tartozó keresési idő most 11,0, mivel majdnem a fél lemezt megesszük, míg a 7000.-ról a 2000. cilinderré jutunk. A 8000. cilinderré vonatkozó ötödik igény 13 milliszekundum keresési időt jelent, az utolsó, azaz az 5000. cilinderré tartozó keresési idő pedig 7. A 2.14. ábra összegzi, hogy ez a megközelítés milyen tevékenységekkel járt. A két algoritmus között 14 milliszekundum az eltérés, amely ugyan nem tűnik jelentősnek, de ne feledjük, hogy ebben az egyszerű példában az igények száma kicsi volt, és az algoritmusok a hat igény közül a negyedikig nem is térnek el.

<i>Igényelt cylinder</i>	<i>Az igénylés ideje</i>
1000	0
3000	0
7000	0
2000	20
8000	30
5000	40

2.12. ábra. Hat blokkhozjárési igény érkezési sorrendje

<i>Igényelt cylinder</i>	<i>A befejezés ideje</i>
1000	8,3
3000	21,6
7000	38,9
8000	50,2
5000	65,5
2000	80,8

2.13. ábra. A blokkhozjárások befejezési időpontjai a lift algoritmus használata esetén

<i>Igényelt cylinder</i>	<i>A befejezés ideje</i>
1000	8,3
3000	21,6
7000	38,9
2000	58,2
8000	79,5
5000	94,8

2.14. ábra. A blokkhozjárások befejezési időpontjai az „első érkezés első kiszolgálás” algoritmus használata esetén

Ha a lemezre váró igények átlagos száma növekszik, akkor a lift algoritmus tovább javítja a teljesítményt. Például, ha a várakozási igények száma megegyezik a cilindretek számával, akkor néhány cylinder kivételével mindegyiket meg kell keresni, és ekkor az átlagos keresési idő megközelíti a minimumot. Ha több lekérdezésünk van, mint amennyi cylinder, akkor tipikusan egynél több igény jut egy cylinderre. Ekkor a

A lift algoritmus tényleges késése

Bár a 2.11. példában azt láttuk, hogy a lemezlejáréshez szükséges átlagos idő csökkenthető, de a nyereség nem egyforma minden igényre. Például a 2.13. és 2.14. ábrákat megvizsgálva észrevehetjük, hogy a 2000. cilinderré vonatkozó igényt az első érkezés első kiszolgálás algoritmus 58,2-kor elégti ki, ezzel szemben a lift algoritmus 80,8-kor. Mivel az igényt 20-kor adták ki, ezért a lemez látszólagos késése az igénylés folyamatára vonatkozóan 38,2-ről 60,8 milliszekundumra változik.

Ha sokkal több lemezlejárési igény várakozik, akkor a lift algoritmus alatti fejpasztázások nagyon hosszú ideig fognak tartani. Ha egy igény éppen lekészte a liftet, akkor a látszólagos késés ebben az esetben valójában rendkívül magas lesz. Viszonyosképpen viszont, ha nem használjuk a lift algoritmust vagy egy másik jó titemező módszert, akkor a teljesítmény csökken, és a lemez nem tudja olyan sebességgel kielégíteni az igényeket, amilyen gyorsan azok generálódnak. A rendszerben végül tetszőleges hosszú késéseket tapasztalunk, vagy egy másodperc alatt csak kevesebb folyamatot lehetne kiszolgálni.

lemezvezérlő rendezheti az egy cilinderré tartozó igényeket, és ezzel csökkenteni tudja az átlagos rotációs késést, és ezzel együtt az átlagos keresési időt is. Vigyázzunk arra, hogy ha az igények száma nagyon nagyra nő, akkor bármelyik igény kiszolgálásához szükséges idő különösen nagy lesz. A következő példa mutatja be ezt az esetet.

2.12. példa: Tegyük fel, hogy megint a *Megatron 747* lemezzel dolgozunk, melynek 8192 cylindere van. Képzeljünk el, hogy 1000 lemezlejárési igény várakozik. Az egyszerűség kedvéért tegyük fel, hogy minden igényelt blokk különböző cylindereken van, 8 cylinderenként. Ha a lemez egyik végétől indulunk és végighaladunk a lemezen, akkor az 1000 igény mindegyikéhez valamilyen több, mint 1 milliszekundum keresési idő, 7,8 milliszekundum rotációs késés és 0,5 milliszekundum átviteli idő tartozik. Így minden 9,3 milliszekundumban ki tudunk elégíteni egy igényt, ami körülbelül a 60%-a a véletlen blokklejáréshez tartozó átlagos 14,4 milliszekundumnak. A teljes ezer igény kielégítése így 9,3 másodpercig tart. Emiatt az egy igény kielégítéséhez szükséges átlagos késés ennek a fele, azaz 4,65 másodperc, ami már számottevő késést jelent.

Most tegyük fel, hogy nagyon sok igényt kell kielégíteni, mondjuk 16 384-et, és az egyszerűség kedvéért feltelesszük, hogy minden cylinderre pontosan két eléjárési igény esik. Ebben az esetben minden keresési idő egy milliszekundum, és az átviteli idő igényenként természetesen fél milliszekundum. Mivel minden cylinderen két blokkot kell elérni, ezért a 2 blokk közül a távolabbi 2/3 útnyira helyezkedik a lemezen, mikor a fejek ehhez a sávhoz érkeznek. Ennek a becslésnek a bizonyítása trükkös, ezt fogjuk elmagyarázni a „Várakozás két blokk közül az utolsóira” keretes részben.

Tehát ennek a két blokknak az átlagos késése a 2/3 körülfordulási időnek a fele

Várakozás két blokk közül az utolsóra

Tegyük fel, hogy egy cilinderen véletlenszerűen van két blokkunk. Legyen x_1 és x_2 a két pozíció a teljes kör törtreszeként megadva, azaz mindkét szám 0 és 1 közé esik. A nagyobbik érték várható értéke vagyunk kíváncsiak. Annak a valószínűsége, hogy a nagyobbik szám kisebb, mint egy 0 és 1 közé eső y szám, megegyezik azzal, hogy mindkét szám, egymástól függetlenül kisebb ennél az y -nál, ami y^2 -tel egyenlő. Így a nagyobbik számhoz tartozó sűrűségfüggvény az y^2 deriváltja, azaz $2y$, ami egy lineáris függvény. A nagyobbik szám várható értékét úgy kapjuk, hogy a sűrűségfüggvénynek és y -nak a szorzatát integráljuk 0 és 1 között, és $\int_0^1 2y^2 dy = 2/3$. Azaz a távolabbi blokk átlagosan a $2/3$ lemezkerületére helyezkedik el.

lesz, azaz $0,5 \times \frac{2}{3} \times 15,6 = 5,2$ milliszekundum. Ezzel az egy blokk eléréséhez szükséges átlagos időt lecsökkentettük $1 + 0,5 + 5,2 = 6,7$ milliszekundumra, ami kevesebb, mint a fele az első érkezés első kiszolgálás ütemezésnél. Kapott átlagos időnek. Másrészt viszont a 16 384 elért összesen 102 másodpercig tart, így egy igény átlagos késése 51 másodperc. \square

2.4.5. Korai beolvasás és nagy léptékű pufferezés

Az utolsó javaslatunk egy másodlagos memória algoritmus felgyorsítására az úgynevezett *korai beolvasás* (prefetching) vagy másképpen *dupla pufferezés*. Bizonyos alkalmazásokban előre meg lehet mondani, hogy a lemeztől milyen sorrendben igényeljük a blokkokat. Ilyen esetben betölthetjük ezeket a blokkokat a központi memória puffereibe, mielőtt szükségünk lenne rájuk. Az ebből származó egyik előny az, hogy csökkenteni tudjuk a blokkeléréshez szükséges átlagos időt úgy, hogy jobb lemeztemezést használunk, például ahogy a lift algoritmus esetében tettük. A 2.12. példában látott blokkhozzáférések felgyorsítását is elérhetjük anélkül, hogy az igények kielégítése során a példában is bemutatott nagy késés lépne fel.

2.13. példa: Ahhoz, hogy egy példát lássunk a dupla pufferezés használatára, nézzük meg újra a kétfázisú, többutas, összefésülő rendezés második fázisát, amit a 2.3.4. részben vázoltunk. Emlékezzünk vissza arra, hogy úgy fésültünk össze 20 rendezett részlistát, hogy minden listától egy blokkot hoztunk be a központi memóriába. Ha annyi rendezett részlistát kell összefésülni, hogy a listákról behozott blokkok teljesen kitöltik a központi memóriát, akkor nem tudunk semmit sem jobbá tenni. Igen ám, de példánkban bőséges mennyiségű memória maradt meg. Például meglehetjük, hogy

minden listához kétblokkos puffert rendelünk, és az egyik puffert feltöltjük, amíg a másiktól az összefésüléshez válogatjuk a rekordokat. Ha kimerül az egyik puffert, akkor késlekedés nélkül átkapcsolunk ugyanannak a listának a másik puffere. \square

Ennek ellenére a 2.13. példa sémája még mindig annyi időt vesz igénybe, amennyi ahhoz kell, hogy a rendezett listák összes (250 000) blokkját beolvassuk. A 2.4.1. rész cilindert alapú stratégiáit és a korai beolvasást kombinálhatjuk is:

1. Ha a rendezett listákat teljes egészében egymás utáni cilindereken tároljuk, még hozzá úgy, hogy minden sávon a blokkok a rendezett lista egymás utáni blokkjai.
2. Ha egy adott listáról bizonyos rekordokra van szükségünk, akkor az egész sávot vagy cilindert beolvassuk.

2.14. példa: Ahhoz, hogy megértsük, miért előnyösebb a sáv vagy cilindert méretű olvasások, vegyük elő megint a kétfázisú, többutas, összefésülő rendezés második fázisát. A központi memóriában annyi hely van, hogy mind a 20 listához két sáv méretű pufferek tartozhassanak. Emlékezzünk vissza arra, hogy a *Megatron 747* lemezen egy sáv 128 Kb-ot tartalmaz, így a teljes puffertírelehez körülbelül 5 megabájt központi memória szükséges. Egy sáv olvasását tetszőleges szektorától kezdhetjük, így egy sáv olvasásához szükséges idő lényegében az átlagos keresési időnek és a lemez egy-szeri körülfordulási idejének az összegével egyezik meg, azaz $6,5 + 15,6 = 22,1$ milliszekundum. Mivel az 1000 cilindert, azaz 8000 sáv összes blokkját be kell olvasni ahhoz, hogy mind a 20 rendezett részlistát elolvassuk, így az összes adat elolvasásához szükséges teljes idő körülbelül 2,95 perc.

Még jobban járunk, ha két cilindert méretű puffereket használunk minden rendezett listához, és az egyiket feltöltjük, amíg a másikat használjuk. Mivel a *Megatron 747* lemeznek 8 sávja van minden cilindern, így összesen 40 egy megabájt méretű puffert fogunk használni. Amennyiben 50 megabájt használható a rendezéshez, akkor van elég hely a központi memóriában ehhez a módszerhez. Cilindert méretű pufferek esetén cilindertként csak egy keresést kell elvégezni. A keresési idő és az egy cilindert tartalmazó 8 sáv olvasási ideje így $6,5 + 8 \times 15,6 = 131,3$ milliszekundum. Ahhoz, hogy mind az 1000 cilindert elolvassuk, 1000-szer több időre van szükség, azaz körülbelül 2,19 percre. \square

A most bemutatott ötlet nem csak az olvasáshoz, hanem analóg módon az íráshoz is használható. A korai beolvasás szellemében a puffertelt blokkok kiírását késleltetjük addig, amíg a közeljövőben már nem kell újból felhasználni a blokkot. Ez a stratégia megőv bennünket a késéstől, noha várakozunk addig, amíg egy blokkot ki lehet írniuk.

Még jobb az a stratégia, amely nagy (sáv vagy cilindert méretű) output puffereket használ. Ha az alkalmazás megengedi, hogy ilyen nagy tömbökben írjunk, akkor a keresési idővel és a rotációs késéssel nem kell számolnunk, így a lemeze írás a maximális lemez átviteli sebességgel végezhető el. Például, ha úgy módosítjuk a rendező algoritmusunk második fázisát, hogy két egy megabájtos output puffert használunk,

akkor feltölthetjük rendezett rekordokkal az egyik puffert, és kiírjuk az egyik cilinderré, miközben feltöljtük a másik output puffert a következő rendezett rekordokkal. Így az íráshoz szükséges idő 2,15 perc lenne, hasonlóan a 2.14. példa olvasási idejéhez, és a teljes 2. fázis 4,3 percig tartana, éppen annyi, amint a 2.9. példa javított 1. fázisában. Összességében a cylinder stratégia, a cylinder méretű pufferezés és a korai beolvasás trükkök kombinációjával a rendezést 8,6 perc alatt is el lehet végezni, szemben a naiv lemezkezelő stratégiával kapott 4 órával.

2.4.6. A stratégiák előnyeinek és hátrányainak összegzése

Az előzőekben öt különböző trükköt láttunk, melyekkel néha javítani lehet egy lemez rendszer működésén. Ezek a következők:

1. Az adatokat cylinderként szervezzük.
2. Egy helyett több lemezt használunk.
3. Tükrözzük a lemezeket.
4. Az igényeket a lift algoritmussal ütemezzük.
5. Sáv vagy cylinder méretű tömbökben előre behozzuk az adatokat.

Megnéztük a fentiek hatását két esetben, amelyek a lemezelési igényeknek két szélsőséges esetét reprezentálják:

- a) A legszabályosabb esetben, amit a kétfázisú, többutas, összefésülítő rendezés első fázisával szemléltettünk, a blokkokat előre megíósolt sorrendben lehet előre beolvasni vagy kiírni, és egyszerre egy folyamat használja a lemezt.
- b) Kisebbl folyamatok gyűjteménye esetén, mint amilyenek a repülőgégy-foglalások, vagy a bankszámlák változtatásai, a folyamatokat párhuzamosan lehet végrehajtani, megoszthatnak ugyanazon a lemezen vagy lemezeken, és előre nem tudunk semmit megíósolni. A kétfázisú, többutas, összefésülítő rendezés második fázisa rendelkezik ezek közül bizonyos jellemzőkkel.

Az alábbiakban ezeknek a módszereknek az előnyeit és hátrányait összegezzük a fenti kétféle alkalmazásra és azokra, amelyek ezek közé esnek.

Cylinder alapú szervezés

- **Előny:** Kiváló az a) típusú alkalmazásokra, ahol a hozzáféréseket előre meg lehet mondani, és csak egy folyamat használja a lemezt.
- **Hátrány:** Nem segít a b) típusú alkalmazásoknál, ahol a hozzáféréseket nem lehet előre megíósolni.

Több lemez használata

- **Előny:** Mindkét típusú alkalmazásra növeli az írási/olvasási igények kielégítési sebességét.
- **Probléma:** Ugyanarra a lemezre vonatkozólag egy időben nem lehet egyszerre több olvasási vagy írási igényt kielégíteni, így a gyorsulási tényező kisebb, mint az a tényező, amennyivel a lemezek számát növeltük.
- **Hátrány:** Több kis lemez költsége meghaladja azt a költséget, amennyibe egy ugyanolyan összkapacitású lemez kerül.

Tükrözés

- **Előny:** Mindkét típusú alkalmazásra növeli az írási/olvasási igények kielégítési sebességét. A több lemez használatánál említett összetükröző elérések problémája nem fordul elő.
- **Előny:** Minden alkalmazásra növeli a hibatűrést.
- **Hátrány:** Két vagy több lemez árat kell megfizetni, de csak egynek megfelelő tárolási kapacitást kapunk érte.

Lift algoritmus

- **Előny:** Csökkenti a blokkok átlagos írási/olvasási idejét abban az esetben, mikor a blokkelérésekről előre nem tudunk semmit mondani.
- **Probléma:** Az algoritmus akkor a leghatékonyabb, mikor sok lemezelési igény várakozik, vagyis mikor az igénylési folyamat átlagos késése magas.

Korai beolvasás/Dupla pufferezés

- **Előny:** Felgyorsítja az elérést, mikor a szükséges blokkokat ismerjük, de az igények időzítése adatfüggő, mint a többutas, összefésülítő rendezés 2. fázisában.
- **Hátrány:** Újabb puffereket igényel a központi memóriában. Nem segít abban az esetben, mikor az elérések véletlenszerűek.

2.4.7. Feladatok

- 2.4.1. feladat:** Tegyük fel, hogy egy *Megatron 747* lemez I/O-igényei akarjuk ütemezni. Az igények a 2.15. ábrán láthatók. Kezdetben a fej a 4000. sávon tartózkodik. Mikorra fogjuk mindegyik igényt teljesen kielégíteni, ha:

Igényelt cylinder	Az igény érzékelési ideje
1000	0
6000	1
500	10
5000	20

2.15. ábra. Négy blokkelérési igény beérkezési ideje

- a) a lift algoritmust használjuk (Kezdetben bármelyik irányba megengedett az indulás),
b) az „első érkezés, első kiszolgálás” algoritmust használjuk.

*! 2.4.2. feladat: Tegyük fel, hogy két *Megatron 747* lemezt használunk, melyek egymásnak tükröképei. Most ne engedjük meg a két lemez bármelyik blokkjának olvasását. Tegyük fel, hogy az első lemez fejét mindig a lemez belső felén találhatjuk. Tegyük fel, hogy az első lemez fejét a külső fél cilindereire korlátozzuk. Tegyük még azt is fel, hogy az olvasási igények véletlenszerűen választott sávokra vonatkoznak, és soha semmit sem kell írunk.

- a) Milyen átlagos sebességgel tudja ez a rendszer olvasni a blokkokat?
b) Hogy viszonyul ez a sebesség a fenti megszorítás nélküli, tükrözött *Megatron 747* lemezekre vonatkozó átlagos sebességhez?
c) Milyen hátránya látszik előre ennek a rendszernek?

! 2.4.3. feladat: Vizsgáljuk meg a kapcsolatot az igények átlagos beérkezési sebessége, a lift algoritmus teljesítménye és az igények átlagos késése között. A probléma egyszerűsítése érdekében tegyük fel a következőket:

1. A lift algoritmusban egy menet mindig a legbelső sávtól a legkülső sávig tart, illetve fordítva, még akkor is, ha a legszálsó cilindrekre nincs is igény.
2. Amikor egy menet elkezdődik, akkor csak azokat az igényeket kell figyelembe venni, amelyek a kezdéskor éppen várakoztak, és semmilyen olyan igényvel nem kell foglalkozni, amely a menetet indulása után érkezett be, még akkor sem, ha a fej éppen egy ilyen cilinderekhez érkezik.¹¹
3. Egy menet alatt egy cilinderre legfeljebb egy blokkigény vonatkozhat.

Jelölje A a beérkezési sebességet, vagyis azt az időt, amennyi két blokkigény beérkezése között telik el. Tegyük fel, hogy a rendszer stabil, kiegyensúlyozott állapotban van, azaz már hosszú idő óta fogadja és megválaszolja az igényeket. Az A függvényében számoljuk ki egy *Megatron 747* lemeze az alábbiakat:

¹¹ Ennek a feltevésnek az a célja, hogy ne kelljen foglalkozni azzal a ténnyel, hogy a lift algoritmusban egy tipikus menet először gyorsan halad, mivel kevés várakozó igény vonatkozik arra a helyre, ahol a fej éppen tartózkodik, és felgyorsul, amikor olyan helyre érkezik a lemezen, ahol mostanáig még nem járt. Önmagában érdekes feladat annak az elemzése, hogy egy menetet alatt hogyan változik az igényűrtés.

- * a) Egy menetet végrehajtása átlagosan mennyi ideig tart?
b) Egy menetet alatt mennyi igényt fogunk kielégíteni?
c) Egy igénynek mennyit kell átlagosan várnia a kiszolgálásra?

*! 2.4.4. feladat: A 2.10. példában láttuk, hogy ha a rendezni kívánt adatokat szétosztjuk négy lemez között, akkor egy időben egymél több blokkot is lehet olvasni. Tegyük fel, hogy az összetéstitő fázis során olyan blokkokat kell olvasni, melyek véletlenszerűen választott lemezekre helyezkednek el. Azt is tételezzük még fel, hogy minden olvasási igényt csak akkor veszünk figyelembe, ha nem olyan lemeze vonatkozik, amely éppen egy másik igényt szolgál ki. Határozzuk meg, hogy az igények kiszolgálását egyszerre átlagosan hány lemez végzi. Megjegyzés: a következő két észrevételt egyszerűsíti a problémát:

1. Amint egy blokkolvasási igényt nem lehet végrehajtani, akkor az összetéstitésnek meg kell állnia, és nem generál több igényt, mivel a kimerült listában már nincs olyan adat, amely a kielégítetlen olvasási igényt generálta.
2. Amint az összetéstitést folytatni lehet, akkor egy olvasási igényt fog generálni, mivel a központi memóriában az összetéstités elhanyagolható ideig tart összevetve az olvasási igény kielégítéséhez szükséges idővel.

! 2.4.5. feladat: Ha egy cilindertől k darab véletlenül választott blokkot akarunk beolvasni, akkor átlagosan mekkora körülfordulást kell tennünk a cilinderen ahhoz, hogy találjassunk mind a hat blokkon?

2.5. Lemezhibák

Ebben és a következő fejezetben megnézzük, hogyan hibázhatnak a lemezek, és miképpen lehet mérsékelni az ilyen meghibásodásokat.

1. A legszokásosabb formája a meghibásodásnak az *ideiglenes meghibásodás*. Ez azt jelenti, hogy amikor megpróbálunk írni vagy olvasni egy szektorot, először nem sikerül, de ismételt próbálkozásokkal végül sikerül írni vagy olvasni.
2. Sülyosabb formájú az a meghibásodás, mikor egy vagy több bit végtrvényesen elromlik, és emiatt lehetetlenné válik egy szektor olvasása, mindegy hányszor próbáljuk újra. Ezt a meghibásodási formát *eszközhibának* hívjuk.
3. Ezzel rokon hibatípus az *írási hiba*, ami azt jelenti, hogy megpróbálunk írni egy szektorot, de az írás nem sikerül, és a szektor korábban írt tartalmát sem lehet már visszanyerni. Ennek egy lehetséges oka, hogy áramkimaradás történt, mialatt a szektorot írtuk.
4. A legkomolyabb formája a lemez meghibásodásának a *lemezhiba*, amikor hirtelen a teljes lemez végtrvényesen olvashatatlaná válik.

Ebben a fejezetben a lemezmegehibásodásnak egy egyszerű modelljét vesszük figyelembe. Áttekinjtük a paritás-ellenőrzést, ami egy lehetséges módszer arra, hogy felderítsük az ideiglenes megehibásodásokat. Megvizsgáljuk még a „stabil tárolást” is. Ez egy olyan lemezszervezési technika, amely arra jó, hogy eszközhiba vagy hibás írások sem eredményeznek végleges adatvesztést. A 2.6. részben megnézzük azokat a technikákat, melyek „RAID” gyűfőnéven ismeretesek. Ezek segítenek abban, hogy megbirkózzunk a lemezhibákkal.

2.5.1. Ideiglenes megehibásodás

A lemezszektorokat általában redundáns bitekkel együtt szokták tárolni. Erről a 2.5.2. részben lesz majd részletesebben szó. Ezeknek a biteknek az a célja, hogy segítségükkel meg tudjuk mondani, hogy amit beolvastunk egy szektorból jó-e vagy hibás, vagy ha kirtunk egy szektort, akkor az írás helyesen történt-e meg.

Egy jól használható modell a lemezolvasásokra a következő: Az olvasási függvény egy (w, s) párt ad vissza, ahol w a szektorból beolvasott adat, és s egy olyan *státuszbit*, amely azt mondja meg, hogy az olvasás sikeres volt-e vagy nem, vagyis megbízhatunk-e abban, hogy w a szektor igazi tartalma. Egy ideiglenes megehibásodás esetén többször is a „rossz” státust kaphatjuk, de ha elégszer (tipikusan legfeljebb 100-szor) ismételtük az olvasási függvényt, akkor végül meg fogjuk kapni a „jó” státust, és biztosak lehetünk abban, hogy azok az adatok, amiket ezzel a státusszal kaptunk vissza, a lemezszektor valódi tartalmával egyeznek meg. A 2.5.2. részben látni fogjuk, hogy előfordulhat, hogy mégis be leszünk csapva, azaz a státus „jó”, de a visszakapott adatok valójában rosszak. Azonban az ilyen eset előfordulásának valószínűségét tetszőlegesen kicsivé tehetjük, ha még több redundanciát adunk a szektorokhoz.

A szektorok írásánál is előfőnyünkre szolgálhat, ha megnézzük annak a státusát, amit írunk. Ahogy a 2.2.5. részben említettük, megehetjük, hogy minden szektort a kírás után megpróbálunk beolvasni, hogy megehatározzuk, vajon sikeres volt-e az írás. Egy nyilvánvaló módszer az ellenőrzésre, hogy beolvassuk a szektort, és összehasonlíjuk azzal a szektórral, amit ki akartunk írni. Azonban ahelyett, hogy a teljes összehasonlítást a lemezvezérlővel végeztetnénk el, egyszerűbb, ha megpróbáljuk beolvasni a szektort, és megenézzük, hogy a státusa „jó”. Ha igen, akkor fel tesszük, hogy helyes volt az írás, ha a státus „rossz”, akkor az írás láthatólag sikertelen volt, és meg kell ismételnünk. Vegyük észre, hogy az olvasáshoz hasonlóan itt is előfordulhat, hogy be leszünk csapva, azaz a státus „jó”, de az írás valójában sikertelen volt. Az olvasásnál említett lehetőség most is rendelkezésre áll, azaz, ha akarjuk, akkor az ilyen tévesztés valószínűségét tetszőlegesen kicsivé tehetjük.

2.5.2. Ellenőrző összegek

Első pillanatra rejtélyesnek tűnhet, hogyan tudja megehatározni az olvasási művelet egy szektor jó/rossz státusát. Ennek ellenére a modern lemezmegehibajótkban használt technikák teljesen egyszerűek: minden szektornak vannak még további bitei, ezeket

hívjuk *ellenőrző összegeknek*, és ezeknek a beállított értéke az ebben a szektorban tárolt adatok értékkeitől függ. Ha olvasáskor azt találjuk, hogy az ellenőrző összeg nem helyes az adatbitekre, akkor „rossz” státust adunk vissza, különben pedig „jó”-t. Annak is van egy kis valószínűsége, hogy az adatbitekkel rosszul olvastuk be, de a téves bitekkel is ugyanazt az ellenőrző összeget kapjuk, mint a helyes bitekkel (és emiatt a rossz bitek is „jó” státust fognak kapni), de elég sok ellenőrző bit használatával ezt a valószínűséget tetszőlegesen kicsivé tehetjük.

Az ellenőrző összegek egyik egyszerű formája a szektor összes bitejének *paritásán* alapul. Ha egy bitekből álló halmazban páratlan sok egyes található, akkor azt mondjuk, hogy a bitek paritása *páratlan*, vagyis a bitekhez tartozó paritás bit értéke 1. Hasonlóan, ha egy bitekből álló halmazban páros sok egyes található, akkor azt mondjuk, hogy a bitek paritása *páros*, és a bitekhez tartozó paritás bit értéke 0. Ebből következők, hogy:

- Ha egy bitekből álló halmazhoz hozzávesszük a nekik megfelelő paritás bitet, akkor az így kapott számok között mindig páros sok egyes szerepel.

Amikor egy szektort írunk, akkor a lemezvezérlő ki tudja számolni a paritás bitet, és hozzáfeszí ahhoz a bisorozatához, amit a szektorba írunk. Emiatt minden szektornak páros a paritása.

2.15. példa: Ha egy szektorban 01101000 volt a bitek sorozata, akkor páratlan sok 1 szerepel, ezért a paritás bit értéke 1. Ha a sorozat végére tesszük a paritás bitet, akkor 011010001 sorozatot kapjuk. Ha a sorozatunk az 11101110 volt, akkor páros sok 1 szerepel, így a paritás bit értéke 0. Ha a sorozat végére tesszük a paritás bitet, akkor az 111011100 sorozatot kapjuk. Vegyük észre, hogy a paritás bit hozzáadásával keletkezett kilenc bit hosszú sorozat mindegyikének páros a paritása. □

Ha a paritás bittel kiegészített bitsorozat olvasása vagy írása során pontosan egy bit hiba keletkezik, akkor páratlan lenne a paritás, vagyis páratlan sok 1 szerepelne. A lemezvezérlő könnyen össze tudja számolni az egyesek számát, és meg tudja határozni a hibát azáltal, ha a szektor paritására páratlant kap.

Természetesen a szektorban egynél több bit is elromolhat. Ha ez történt, akkor 50% a valószínűsége annak, hogy az 1 értékű bitek száma páros, és ekkor nem fogjuk a hibát észrevenni. Ha több paritás bitet használunk, akkor nagyobb lesz az esélyünk arra, hogy sok hibát fel tudunk ismerni. Például használhatunk 8 paritás bitet, egyet minden bájt első bitejére, egyet minden bájt második bitejére és így tovább, egészen a nyolcadikig, amely a bájtok utolsó bitejére vonatkozik. Nagy tömegű hiba esetén 50% annak a valószínűsége bármelyik paritás bit esetére, hogy hibát jelez, és annak az esélye, hogy a nyolcból egyik sem jelez hibát, csak egy a 2^8 -hoz, azaz 1/256. Általában, ha n független bitet használunk ellenőrző összegeként, akkor annak az esélye, hogy nem vesszük észre egy hibát, mindössze $1/2^n$. Például, ha 4 bájt számunk egy ellenőrző összegre, akkor annak az esélye, hogy nem fogunk felismerni egy hibát, csak 1 a 4 milliódhoz.

Igaz, hogy az ellenőrző összegek majdnem biztosan felismerik, ha az eszköz elromlott, vagy az olvasás vagy írás nem sikerült hibátlanul, de nem segítenek kijavítani a hibát. Továbbá, írás esetén abba a helyzetbe kerülhetünk, hogy már áfirtuk egy szektor korábbi tartalmát, de mégsem tudjuk elolvasni az új tartalmat. Ez a helyzet nagyon komoly is lehet. Képzeliük el például, hogy a folyószámához egy kis összeget akarunk adni, és most elvesztettük a folyószámja eredeti egyenlegét és az újat is. Ha biztosak lehetnénk abban, hogy a szektor tartalma vagy az új vagy a régi egyenleg, akkor csak azt kellene meghatározni, hogy az írás sikerült-e vagy sem.

Ahhoz, hogy ezt a problémát kezelni tudjuk egy vagy több lemezen, egy olyan elv megvalósítását használjuk, amit *stabli tárolásnak* hívunk. Az alapötlet, hogy a szektorból párokat képezzünk, és minden pár egy X szektor tartalmát reprezentálja. Az X -et reprezentáló szektorhoz tartozó pár tagjaira bal (X_L) és jobb (X_R) másolatként hívhatunk. Feltelesszük, hogy a másolatok írásakor elég sok paritás-ellenőrző bitet használunk, így kizárhatjuk annak az esélyét, hogy egy rossz szektor jónak látszik a paritás-ellenőrzéskor. Tehát feltelesszük, hogy ha az olvasási függvény a (w , j) értéket adja vissza az X_L vagy az X_R esetében, akkor a w az X valódi értéke. A stabli tárolás írási elve a következő:

1. Írjuk az X értékét X_L -be. Ellenőrizzük, hogy az érték státusa „jó”; vagyis a paritás-ellenőrző bitek helyesek a kiírt másolatban. Ha nem, akkor ismételiük meg az írást. Ha egy meghatározott számú próbálkozás után sem sikerül X értékét X_L -be írni, akkor azt tesszük fel, hogy ebben a szektorban megsértült az eszköz. Ekkor valamilyen javítási elvet kell alkalmazni, például egy másik szabad szektorral helyettesítjük X_L -t.
2. Ismételiük meg 1.-t az X_R -re.

A stabli tárolás olvasási elve a következő:

1. Ahhoz, hogy X értékét visszanyerjük, olvassuk be X_L -t. Ha a „rossz” státust kapjuk vissza, akkor ismételiük meg az olvasást valamilyen előre megadott számszor. Ha végül kapunk egy értéket, amelynek „jó” a státusa, akkor ezt az értéket vesszük X -nek.
2. Ha nem tudjuk X_L -t elolvasni, akkor 1.-t ismételiük X_R -rel.

2.5.4. A stabli tárolás hibakezelő képessége

A 2.5.3. részben leírt elv több különféle hibatípus ellen használható, melyeket az alábbiakban sorolunk fel.

1. *Eszközhiba*. Tegyük fel, hogy az X -et már letároltuk az X_L és X_R szektorokban. Ekkor, ha valamelyik a kettő közül eszközhiba miatt állandó jelleggel olvashatat-

lanná válik, akkor az X -et a másikkól tudjuk kiolvasni. Ha az X_R hibás, de az X_L nem, akkor az olvasási elv alapján, az X_R -re ügyet semelve az X_L -t helyesen fogjuk elolvasni. Akkor fogjuk észrevenni, hogy az X_R hibás, mikor legközelebb megpróbálunk új X értéket írni. Ha csak az X_L sértült meg, akkor nem sikerül „jó” státust kapni az X -re, bárhol is próbáljuk az X_L -t beolvasni. (Emlékezzünk vissza arra, hogy a fellelelézésünk szerint egy rossz szektor mindig „rossz” státust ad vissza, még ha a valóságban van is egy csöpp esély arra, hogy „jó” lesz a státus, ha véletlenül az összes paritás-ellenőrző bit megfelelő értékű.) Tehát végrehajtuk az olvasási algoritmus 2. lépését és helyesen beolvassuk X -et az X_R -ből. Vegyük észre, hogy ha az X_L és X_R mindegyike hibás, akkor az X értéket nem lehet elolvasni, de az egyidejű két meghibásodás valószínűsége rendkívül kicsi.

2. *Íráshiba*. Tegyük fel, hogy az X írása közben rendszerhiba, például áramkimaradás történt. Lehet, hogy el fog veszni az X a központi memóriából, és ráadásul az X -nek az a másolata, amelyet éppen kiírtunk, összezavarodik. Például az X új értékek egy részével már megírtuk a szektor felét, de a szektor másik fele még változatlan. Mikor a rendszer újra működőképessé válik, megvizsgáljuk az X_L -t és X_R -t, és biztonsággal meg tudjuk határozni az X régi vagy új értékét. A lehetséges esetek a következők:

- a) Akkor történt a hiba, mikor az X_L -t írtuk. Ekkor azt fogjuk kapni, hogy az X_L státusa „rossz”, viszont mivel az X_R -t nem kellett írni, ezért az X_R -nek „jó” a státusa (kivéve, ha éppen egy X_R -re vonatkozó eszközhiba is történt, aminek olyan kicsi az esélye, hogy nyugodtan elvethetjük). Így megkaphatjuk az X régi értékét. Az X_L sértültségnek kijavításához megtehetjük, hogy X_R -t X_L -be másoljuk.
- b) A hiba az X_L írása után történt. Ekkor várhatóan az X_L -nek „jó” lesz a státusa, és így kiolvashatjuk az X_L új értékét az X_L -ből. Megjegyezzük, hogy az X_R -nek lehet, hogy „rossz” a státusa, és ekkor az X_L -t X_R -be kell másolni.

2.5.5. Feladatok

2.5.1. feladat: Számoljuk ki a következő bisorozatokat paritás bitjét:

- * a) 00111011.
- b) 00000000.
- c) 10101101.

2.5.2. feladat: Egy bináris sorozat végére kétféle paritás bitet tesszünk, az első a páratlan pozíciókhoz tartozó paritás bit, a második pedig a páros pozíciókhoz tartozó. A 2.5.1. feladatban szereplő összes sorozatra határozzuk meg ezeket a paritás biteket.

2.6. Lemezhiba helyreállítás

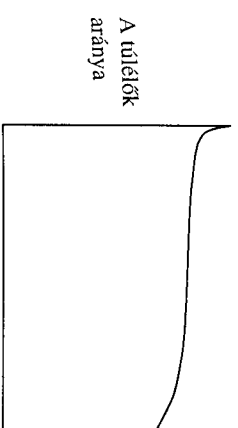
Ebben a fejezetben a lemezhibák közül a leg súlyosabbat fogjuk tárgyalni, az úgynevezett „fejsérülést” (head crash), amikor az adatok végérvényesen tönkrementek. Egy ilyen esemény bekövetkezésakor az adatokból semmi sem lehet helyreállítani, csak abban az esetben, ha az adatokat egy másik eszközzel, például szalagos archíváló eszközre vagy a 2.4.3. részben vizsgált tükrözött lemezre előtte lementettük. Ha nem lenne mentésünk, akkor egy ilyen helyzet katasztrofális lenne az adatbázis-kezelő rendszerre épülő legfőbb alkalmazásokra (banki, pénzügyi alkalmazásokra, repülőgépjegyek vagy más módon helyfoglalási rendszerekre, raktárnyilvántartó rendszerekre és egyebekre) nézve.

Sok különböző sémát fejlesztettek ki, hogy csökkentse a lemezhibából származó adatvesztést. Ezek általában redundanciát okoznak, például a 2.5.2. részben tárgyalt paritás-ellenőrzés ötletének kiterjesztésével, vagy a 2.5.3. részben leírtakhoz hasonló megduplázott szektorok használatával. Ezeknek a stratégiáknak az osztályát közös kifejezéssel RAID¹²-nek, *független lemezek redundáns tömbjének* (Redundant Arrays of Independent Disks) hívjuk. Elsődlegesen három sémát fogunk megnevezni közelebbről, melyeket 4, 5 és 6 szintű RAID-nek hívunk. Ezek a RAID-sémák arra is jók, hogy kezeljék a 2.5. részben tárgyalt meghibásodásokat, azaz az eszközhibát és egy szektor adatainak ideiglenes rendszerhibából származó sérüléseit.

2.6.1. A lemezek meghibásodási modelljei

Ahhoz, hogy elkezdjük a lemezérülések tárgyalását, először az ilyen hibák statisztikáit kell megvizsgálnunk. A legegyszerűbben úgy írhatjuk le az ilyen hibák viselkedését, hogy megadunk egy mértéket, a *meghibásodás várható idejét*. Ez a szám annak az időnek a hosszát jelenti, amennyi idő elteltével a lemezek egy előre adott sokaságának az 50%-a katasztrofálisan tönkremegy, vagyis olyan lemezhiba történik, hogy soha többé nem lesz már olvasható a lemez. A modern lemezek esetében a várható meghibásodási idő körülbelül 10 év.

A legegyszerűbben úgy használhatjuk fel ezt a számot, hogy fel tesszük, hogy a meghibásodások lineárisan történnek, vagyis ha 50% 10 év alatt megy tönkre, akkor 5% megy tönkre az első, a második stb. évben 20 éven keresztül. A valóságban a lemezek túljelési százaléka inkább a 2.16. ábrán megadott grafikonot követi. Az elektromos berendezések legtöbb típusára igaz, hogy az életciklus korai szakaszában sok lemezhiba jelenkezik, amelyek főleg a lemez gyártása közben keletkezett apró hibákból származnak. Ezeknek a gyártási hibáknak a legtöbbit remélhetőleg észre vesszük, mielőtt a lemez elhagyja a gyárat, de néhánynál még hónapok elteltével sem jelentkeznek. Egy olyan lemez, amelynél nem jelenkezik hamar semmilyen hiba, valószínűleg sok évig fog jól működni. Az életciklus későbbi szakaszában több tényező (a haszná-



2.16. ábra. Lemezekre vonatkozó meghibásodási ráta görbéje

latból eredő kopás, a parányi porsszemcsék összegzett hatása) növeli a meghibásodás esélyét.

Valójában egy lemez meghibásodásának várható ideje nem kell, hogy megegyezzen az adatvesztés várható idejével. Ennek az az oka, hogy sok olyan séma áll a rendelkezésünkre, melyek biztosítják, hogyha egy lemez meghibásodik, akkor más lemezek segítenek helyreállítani a sérült lemez adatait. Ebben a fejezetben át fogjuk tekinteni a legáltalánosabb sémákat.

Ezeknek a sémáknak mindegyike azzal kezdődik, hogy egy vagy több lemez adatakat tartalmaz (ezeket hívjuk majd *adatlemezeknek*), és ezekhez hozzávesszünk egy vagy több olyan lemezt, amelyeken a tárolt információt teljesen az adatlemezek tartalma határozza meg. Ez utóbbiakat nevezzük *redundáns lemezeknek*. Amikor egy adatlemez vagy egy redundáns lemez megsérül, akkor a többi lemez használható a sérült lemez visszaállítására, és emiatt nem lesz maradandó információvesztés.

2.6.2. A tükrözés mint redundanciatechnika

A legegyszerűbb séma, ha minden lemezt úgy tükrözzünk, ahogy a 2.4.3. részben ezt leírtuk. Az egyik lemezt *adatlemeznek*, a másik lemezt *redundáns lemeznek* hívjuk, hogy melyik melyik, nem játszik szerepet ebben a sémában. A tükrözést mint az adatvesztés elleni egyik lehetséges védekezést, gyakran *I. szintű RAID-nek* hívják. Ezáltal a memóriavesztés várható ideje sokkal nagyobb, mint a lemezhiba várható ideje, ahogy ezt a következő példa szemlélteti. Lényegében a tükrözés és más redundancia-sémák esetén adatvesztés csak akkor történhet, ha a második lemez is tönkremegy, miáltal az első sérülését javítják.

2.16. példa: Tegyük fel, hogy mindegyik lemeznek 10 év a várható meghibásodási ideje. A meghibásodásokra a 2.6.1. részben leírt lineáris modellt használjuk, amely azt jelenti, hogy annak az esélye, hogy egy lemez tönkremegy, 5% évente. Ha a lemezeket tükrözzük, akkor egy lemezérülés esetén csak ki kell cserélnünk egy jó lemezzel, és a (tükröképet tartalmazó) tükrölmezt az új lemezre kell másolni. Végül újra két lemezünk lesz, melyek tükröképei egymásnak, és ezzel a rendszer visszaállt az előző állapotára.

¹² Korábban a RAID betűszóban az I beír az olcsó (Inexpensive) szónak a kezdőbetűjét jelentette, és ezzel az értelmezéssel még mindig találkozhatunk a szakrodalomban.

Az egyetlen baj, ami történet, hogy a másolása közben a tükörlemez is tönkremegy. Mivel most legalább az adatok egy részének mindkét másolata elveszett, ezért nincs mód a helyreállításra.

De vajon milyen gyakran fordul elő az eseményeknek ez a láncolata? Tegyük fel, hogy a megsérült lemez cserélésének folyamata 3 órát vesz igénybe, amely 1/8 nap, azaz 1/2920 része az évnek. Mivel 5% éves meghibásodási ráfát tételeztünk fel, ezért annak a valószínűsége, hogy a tükörlemez megsérül a másolás alatt $(1/20) \times (1/2920)$, azaz 1 az 58 400-hoz. Ha egy lemez 10 évenként megy tönkre, akkor két lemezből egy átlagosan 5 évenként hibásodik meg. Minden 58 400 ilyen hiba közül egy fog adatvesztést eredményezni, azaz másképpen fogalmazva az adatvesztést eredményező meghibásodás várható ideje $5 \times 58\,400 = 292\,000$ év. \square

2.6.3. Paritásblokkok

Bár a lemezek tükörözése hatékonyan képes csökkenteni az adatvesztéssel járó lemezhiba valószínűségét, de ehhez annyi redundáns lemezre van szükség, mint amennyi az adatlemezek száma. Egy másik lehetséges megközelítés, hogy az adatlemezek számától függetlenül csak egyetlenegy redundáns lemezt használjunk. Ezt gyakran *4. szintű RAID-nek* nevezik. Fellesszük, hogy a lemezek egyformák, azaz minden lemezen 1-től n -ig számozhatjuk meg a blokkokat. Természetesen minden lemez minden blokkjában ugyanannyi bit van, például az alappéldánkban szereplő *Megatron 747* lemezeinkön minden blokk 4096 bájt méretű, azaz bármelyik blokkban $8 \times 4096 = 32\,768$ bit van. A redundáns lemezen az i . blokk paritás-ellenőrzéseket tartalmaz az összes adatlemez i . blokkjaira nézve. Vagyis az adatlemezek és a redundáns lemez i . blokkjainak j . bitei között páros sok 1-nek kell szerepelnie, és a redundáns lemez biteit mindig úgy választjuk meg, hogy ez a feltétel igaz legyen.

A 2.15. példában láttuk, hogy lehet elérni, hogy igazza vájon ez a feltétel. A redundáns lemezen a j . bitet 1-nek választjuk, ha páratlan sok adatlemezen szerepel 1 ezen a biten. A redundáns lemezen a j . bitet 0-nak választjuk, ha páros sok adatlemezen szerepel 1 ezen a biten. Ennek a számtásnak a neve *modulo-2 összegzés*. Tehát adott bitek modulo-2 összege 0, ha páros sok 1 szerepel a bitek között, és 1, ha páratlan sok 1 szerepel.

2.17. példa: Vegyük a lehető legegyszerűbb esetet, vagyis mikor a blokkok csak egyetlenegy bájtól, azaz 8 bitet tartalmaznak. Legyen három adatlemezünk, nevezzük őket 1., 2. és 3. lemeznek, és a 4. lemez legyen a redundáns lemez. Vegyük szemügyre, mondjuk, az összes lemez első blokkját. Ha az adatlemezek első blokkjaiban a következő bitsorozatok találhatók:

1. lemez: 11110000
2. lemez: 10101010
3. lemez: 00111000

akkor a redundáns lemez 1. blokkjában a paritás-ellenőrző bitek a következők:

4. lemez: 01100010

Vegyük észre, hogy a négy darab 8 bit hosszú sorozat közül minden pozícióban páros számúnál látunk 1 értéket. Két 1 érték szerepel az 1., 2., 4., 5. és 7. helyen, négy 1 van a 3. helyen, és zéró darab 1-et találunk a 6. és 8. helyen. \square

Olvasás

Egy adatlemezről ugyanúgy olvasunk be blokkokat, mint ahogy bármilyen más lemezről. Általában semmi okunk, hogy a redundáns lemezt olvassuk, de ha akarjuk, megtehetjük. Bizonyos körülmények között a redundáns lemezt arra is használhatjuk, hogy a segítségével az adatlemezek egyikének két különböző blokkját tudjuk egy kis ritkkel egy időben beolvasni, bár várhatóan ritkán teljesülnek az ehhez szükséges feltételek.

2.18. példa: Tegyük fel, hogy az első adatlemezen éppen olvasunk egy blokkot, mikor ugyanennek a lemeznek egy másik blokkjára, mondjuk az elsőre beérkezik egy másik olvasási igény. Közönséges esetben meg kellene várnunk, hogy az első igény befejeződik. Viszont, ha éppen ekkor a többi lemez egyikét sem használjuk, akkor ez-alatt beolvashatjuk ezekről az első blokkot, és a modulo-2 összeadás segítségével ki-számolhatjuk az első lemez első blokkját.

Speciálisan legyenek a lemezek és az első blokkjaik ugyanazok, mint a 2.17. példában, azaz a 2., 3. és a redundáns lemezt beolvassva a következő blokkokat kapjuk:

2. lemez: 10101010
3. lemez: 00111000
4. lemez: 01100010

Ha most minden oszlopban vesszük a bitek modulo-2 összegét, akkor az alábbiakat kapjuk:

1. lemez: 11110000

ami pontosan megegyezik az első lemez első blokkjával. \square

Írás

Amikor egy adatlemezen egy új blokkot akarunk írni, akkor nem csak ezt a blokkot kell megváltoztatni, hanem a redundáns lemez megfelelő blokkját is, hogy továbbra is az összes adatlemez megfelelő blokkjainak paritás-ellenőrzéseit tartalmazza. Egy naiv megoldás lenne, ha az n adatlemez megfelelő blokkjait beolvassánk, vennénk a modu-

10-2 összegzést, és a redundáns lemez blokkját újraíránk. Ez a módszer $n - 1$ olyan adatblokkot olvas be, amit nem is írunk át, kirírja az újraírt adatblokkot, és egy blokkot ír újra a redundáns lemezre. Ez összesen $n + 1$ lemez I/O-műveletet jelent.

Jobb az a megközelítés, mely szerint csak az újraírt i . adatblokknak nézzük meg a régi és az új értékét. Ha vesszük ezek modulo-2 összegét, akkor tudni fogjuk, hogy az összes lemez i . blokkjában melyik pozícióban változott meg az 1 értékek száma. Mivel egy ilyen változás csak eggyel különbözhet az előző értéktől, ezért páros számú egységből páratlan számú egyes lesz. Ha most a redundáns lemezre is megváltoztatjuk az értéket ugyanezen a pozíción, akkor az egyesek száma újból minden helyen páros lesz. Ezekhez a számításokhoz négy lemez I/O-műveletre van szükség.

1. Beolvassuk a változani kívánt adatblokk régi értékét.
2. Beolvassuk a redundáns lemezről a megfelelő blokkot.
3. Kirírjuk az új adatblokkot.
4. Újraszámoljuk és kirírjuk a redundáns lemez blokkját.

2.19. példa: Tegyük fel, hogy három adatlemezen az első blokkok ugyanazok, mint a 2.17. példában:

1. lemez: 11110000
2. lemez: 10101010
3. lemez: 00111000

Tegyük most fel, hogy a második lemezen a blokkot 10101010-ről 11001100-re változtatjuk. Ha most a 2. lemez blokkjának régi és új értékének vesszük a modulo-2 összegét, akkor a 01100110 sorozatot kapjuk. Ebből kiolvashatjuk, hogy a redundáns lemez első blokkjában a 2., 3., 6. és 7. helyen kell változtatni. Beolvassuk ezt a blok-

A modulo-2 összegzések algebraja

Abhoz, hogy jobban megértsük a paritás-ellenőrzésekhez használt trükköket, hasznos lehet, ha ismerjük, hogy milyen algebrai szabályok vonatkoznak a bitvektorok modulo-2 összeadási műveletére. Ezt a műveletet \oplus jellel jelöljük. Például $1100 \oplus 1010 = 0110$. Az alábbiakban megadunk néhány hasznos szabályt a \oplus műveletre:

- A *kommutatív törvény*: $x \oplus y = y \oplus x$.
- Az *asszociatív törvény*: $x \oplus (y \oplus z) = (x \oplus y) \oplus z$.
- A megfelelő hosszú csupa 0 vektor, melyet $\bar{0}$ -val jelölünk, a \oplus művelet *egységeleme*, vagyis $x \oplus \bar{0} = \bar{0} \oplus x = x$.
- A \oplus műveletnek saját maga az inverze, azaz $x \oplus x = \bar{0}$. Ennek fontos következménye, hogyha $x \oplus y = z$, akkor hozzáadhatjuk x -et mindkét oldalhoz, amivel azt kapjuk, hogy $y = x \oplus z$.

kot: 01100010. Ezt a blokkot azzal az új blokkal helyettesítjük, amelyet úgy kapunk, hogy a megfelelő pozíciókban változtatunk, ami hatását tekintve megegyezik azzal, mintha vennénk a redundáns lemezről beolvasott blokknak és a 01100110-nak a modulo-2 összegét, azaz 00000100-t. Egy másik mód arra, hogy kifejezzük az új redundáns blokkot az, hogy modulo-2 összeadjuk az újraírt blokk régi és új értékét és a redundáns lemez régi értékét. Ezzel a példánkban a 4. lemez (három adatlemez és egy redundáns lemez) első blokkja a második lemez blokkjának és a redundáns blokkon végrehajtott szükséges újraszámolás-kirírása után a következő lesz:

1. lemez: 11110000
2. lemez: 11001100
3. lemez: 00111000
4. lemez: 00000100

Vegyük észre, hogy a fenti blokkok esetében minden oszlopban továbbra is páros sok egyes szerepel.

Vegyük észre, hogy példánkban egy adatblokk újraírásához 4 lemez I/O-műveletre volt szükség, annyira, amennyit a fenti séma alapján az adatblokkok írásához kell használni. Ez a műveletszám most történetesen megegyezik a naiv módszer műveletigényével. A naiv séma esetén az újraírt blokk kivételével beolvassuk a többi blokkot, és közvetlenül újraszámoljuk a redundáns blokkot, azaz beolvassuk az adatokat az első és a harmadik lemezről; ez két művelet, és két művelet kell még a második lemez és a redundáns lemez írásához, azaz összesen négy műveletre van szükség. Természetesen, ha háromnál több adatlemeziünk lenne, akkor a naiv séma I/O-száma az adatlemek számával lineárisan nőne, míg az itt javasolt séma költsége továbbra is négy művelet maradna. \square

Hibajavítás

Most nézzük meg mit tennénk, ha az egyik lemez megsérülne. Ha ez a redundáns lemez, akkor kicseréljük egy új lemezre, és újraszámoljuk a redundáns blokkokat. Ha a tönkrement lemez egy adatlemez, akkor ezt kell kicserélni egy jó lemezre, és a többi lemez adataiból újra kell számolni az adatait. Bármely hiányzó adat újraszámolásához használhatunk egy tulajdonképpen egyszerű szabályt, amely nem függ attól, hogy melyik lemez, adat vagy redundáns lemez ment tönkre. Mivel tudjuk, hogy az összes lemezen a megfelelő bitek között páros számú egyes szerepel, ebből az következik, hogy:

- Bármelyik bit az összes többi lemez megfelelő helyén álló bitek modulo-2 összege.

Aki kétfelkudne ebben a szabályban, annak csak két esetet kell végiggondolnia. Ha a kérdéses bit 1, akkor a megfelelő bitek között páratlan sok egyesnek kell szerepelnie, azaz a modulo-2 összegük 1. Ha a kérdéses bit 0, akkor a megfelelő bitek között páros sok egyesnek kell szerepelnie, azaz a modulo-2 összegük 0.

2.20. példa: Tegyük fel, hogy a második lemez megy tönkre. Ki kell számolnunk a leserélt lemez minden blokkját. Követve a 2.17. példát, nézzük meg, hogy lehet újra-számolni a második lemez első blokkját. Adottak az első, harmadik és a redundáns lemez megfelelő blokkjai, azaz a következő a helyzet:

1. lemez: 111110000
2. lemez: ??????????
3. lemez: 001111000
4. lemez: 011100010

Ha most minden oszlopban vesszük a modul-2 összeget, akkor arra következtethetünk, hogy a hiányzó blokk 10101010, ahogy azt a 2.17. példa kiindulásánál láttuk. \square

2.6.4. Egy továbbfejlesztés: az 5. szintű RAID

A 2.6.3.-ban leírt 4. szintű RAID eredményesen használható adatok megőrzésére, kivéve, ha majdnem egyszerre két lemez megy tönkre. Hogy hol van a legnagyobb baj ezzel a módszerrel, az rögtön látható, ha újra megvizsgáljuk, hogy mi történik egy új adablokk írásánál. Bármilyen sémát is használunk a lemezek módosítására, a redundáns lemez blokkját olvasni és írni is kell. Ha n darab adattlemeznél van, akkor a redundáns lemezre írások száma az n -szerese annak, mint ahányszor átlagosan írunk egy tetszőleges adattlemre.

Viszont a 2.20. példában megfigyeltük, hogy a javítási szabály ugyanaz az adattlemezekre és a redundáns lemezre, azaz venni kell a többi lemez megfelelő bitjeinek modulo-2 összegét. Emiatt nem kell külön foglalkoznunk azzal, hogy melyik a redundáns lemez és melyik az adattlem. Sőt minden lemezt úgy tekinthetünk, mintha bizonyos blokkokhoz tartozó redundáns lemez lenne. Ezt a továbbfejlesztési gyakorlatom 5. szintű RAID-nek hívják.

Például ha $n + 1$ lemeznél van, melyeket 0-tól n -ig számozzunk, akkor a j . lemez i . cilindertét redundánsnak tekintjük, ha az $i-t-n + 1$ -gyel osztva j -t kapunk maradékul.

2.21. példa: Vegyük az alappéldánkat, azaz $n = 3$, így 4 lemeznél van. Az első, azaz a 0 sorszámú lemeznek a 4, 8, 12 stb. sorszámú cilinderei lesznek redundánsak, mert ezek a számok adnak 0-t, ha négyvel osztjuk őket. Az első lemezen az 1, 5, 9 stb. blokkok lesznek redundánsak, a 2. lemezen a 2, 6, 10, ... és a 3. lemezen a 3, 7, 11, ... sorszámú blokkok.

Ennek eredményeképpen minden lemez olvasási és írási terhelése megegyezik. Ha minden blokkot egyenlő valószínűséggel írunk, akkor egy írás esetén minden lemezen $1/4$ az esély, hogy a blokk azon a lemezen van. Ha nincs rajta, akkor $1/3$ az esélye, hogy ennek a blokknak a redundáns lemeze. Tehát a négy közül mindegyik lemez az

$$\frac{1}{4} + \frac{3}{4} \times \frac{1}{3} = \frac{1}{2}$$

részt vesz részt. \square

2.6.5. Mi a teendő, ha több lemez is tönkremehet?

A hibajavító kódok elmélete lehetővé teszi, hogy tetszőleges számú lemez (adat vagy redundáns lemez) tönkremenését kezelni tudjuk, ha elég sok redundáns lemezt használunk. Ez a stratégia jelenti a legmagasabb RAID-fokozatot, a 6. szintű RAID-et. Adni fogunk egy egyszerű példát, amelyben két egyidejű tönkremenés kijavítható. Az ehhez használt stratégia a legegyszerűbb hibajavító kódon, a *Hamming-kódon* alapul.

A leírásunkban egy olyan rendszerrel foglalkozunk, amelynek 7 lemeze van, 1-től 7-ig számozva. Az első négy adattlemez, a maradék három pedig redundáns lemez. Az adattlemezek és a redundáns lemezek közti kapcsolatot a 2.17. ábrán látható 0-kbóli és 1-ekből álló 3×7 -es mátrix írja le. Vegyük észre, hogy:

- a) Minden lehetséges 0-kbóli és 1-ekből képezhető oszlop megjelenik a 2.17. ábra mátrixában, kivéve a csupa nullából álló oszlopot.
- b) A redundáns lemezek oszlopaiban csak egy egyes szerepel.
- c) Az adattlemezek oszlopaiban legalább két egyes szerepel.

Lemez száma	Adat			Redundáns			
	1	2	3	4	5	6	7
	1	1	1	0	1	0	0
	1	1	0	1	0	1	0
	1	0	1	1	0	0	1

2.17. ábra. Redundancia minta egy olyan rendszerre, amely két lemez egyidejű tönkremenését is helyre tudja hozni

Ennek a nullákból és egyesekből álló három sornak a következő az értelme. Ha megnézzük mind a hét lemez megfelelő bitjeit, akkor azokra a lemezekre modulo-2 összegezve a biteket, ahol a sorban egyes szerepel, nullát kapunk. Másiképpen elmondva, azok a lemezek, amelyekhez a sorban egyes tartozik, együttesen 4. szintű RAID-sémájú lemezhalmozatot alkotnak. Tehát egy redundáns lemez bitjeit úgy számolhatjuk ki, hogy megnézzük melyik sorban szerepel egyes ennek a redundáns lemeznek az oszlopában, kiválasztjuk azokat a lemezeket, amelyeknél szintén egyes szerepel ebben a sorban, és ezeknek a lemezeknek a megfelelő bitjeit modulo-2 összeadjuk. Ebből a szabályból a 2.17. ábra mátrixára az alábbiak következnek:

1. Az 5. lemez bitjeit úgy kapjuk, hogy az 1., 2. és a 3. lemez megfelelő bitjeit modulo-2 összeadjuk.
2. Az 6. lemez bitjeit úgy kapjuk, hogy az 1., 2. és a 4. lemez megfelelő bitjeit modulo-2 összeadjuk.
3. Az 7. lemez bitjeit úgy kapjuk, hogy az 1., 3. és a 4. lemez megfelelő bitjeit modulo-2 összeadjuk.

Mindjárt látni fogjuk, hogy a mátrixban a biteknek ez a speciális választása egy egyszerű szabályt ad a kezünkbe, amellyel két lemez egyidejű tönkremenetét is helyre tudjuk hozni.

Olvasás

Bármelyik adatlemezről normálisan olvashatjuk be az adatokat. A redundáns lemezeket figyelmen kívül hagyhatjuk.

Írás

Hasonló elvet követünk, mint a 2.6.4. részben vázolt írásra vonatkozó stratégiában, csak most lehet, hogy több redundáns lemezzel kell törődnünk. Ha egy adatlemezben írni akarunk egy új blokkot, akkor kiszámoljuk a blokk régi és új változatának a modulo-2 összegét. Megnézzük, hogy melyek azok a redundáns lemezek, amelyekhez van olyan sor a mátrixban, hogy az adatlemez oszlopában és a redundáns lemez oszlopában is egyes szerepel. Ezeknek a redundáns lemezeknek megfelelő blokkjához modulo-2 hozzáadjuk az előbb kiszámolt biteket.

2.22. példa: Megint tegyük fel, hogy a blokkok csak 8 bit hosszúak, és fordítsuk a fi-gyelmünket a 6. szintű RAID példában szereplő 7 lemez első blokkjaira. Először tegyük fel, hogy az első adat és redundáns blokkok a 2.18. ábrán látható módon vannak megadva. Vegyük észre, hogy az 5. lemez blokkja az első három lemez blokkjainak modulo-2 összege, a 6. sor az 1., 2. és 4. sorok modulo-2 összege, és az utolsó sor az 1., 3. és 4. sorok modulo-2 összege.

Tegyük fel, hogy a 2. lemez első blokkját át akarjuk írni 00001111-re. Ha ezt

Lemez	Tartalom
1)	11110000
2)	10101010
3)	00111000
4)	01000001
5)	01100010
6)	00011011
7)	10001001

2.18. ábra. A hét lemez első blokkja

Lemez	Tartalom
1)	11110000
2)	00001111
3)	00111000
4)	01000001
5)	11000111
6)	10111110
7)	10001001

2.19. ábra. A lemezek első blokkja azután, hogy a 2. lemezt újraírjuk és a redundáns lemezeket is megváltoztatjuk

modulo-2 hozzáadjuk a blokk régi értékéhez, azaz 10101010-hoz, akkor az 10100101 sorozatot kapjuk. Ha egy pillantást vetünk a 2.17. ábra 2. lemezének oszlopára, akkor azt látjuk, hogy az első két sorban szerepel egyes, a harmadikban nem. Mivel az első két sorban az 5. és a 6. redundáns lemez oszlopán látunk szintén egyest, ezért ezeknek az első blokkjához kell modulo-2 hozzáadni a most kiszámolt 10100101 sorozatot. Emiatt ezekben a blokkokban az 1., 3., 6. és 8. pozíciókban szereplő értékeket az ellenkező értékre állítjuk. A 2.19. ábrán látható az összes lemez változás utáni első blokkjának tartalma. Vegyük észre, hogy továbbra is érvényben marad a 2.17. ábrán megfigyelt megszorítás, azaz a 2.17. ábra mátrixának bármelyik sora alapján kiválasztva azokat a lemezeket, amelyekre egyes szerepel a sorban, és ezeknek a lemezeknek a megfelelő blokkokat modul-2 összeadva, eredményül mindig a csupa 0 sorozatot kapjuk.

Hibajavítás

Most nézzük meg, hogyan lehet az előbb vázolt redundanciasémát felhasználni arra, hogy két egyidejűleg tönkrement lemezt helyreállítsunk. Legyen a és b a két tönkrement lemez. Mivel a 2.17. ábra mátrixának minden oszlopa különböző, ezért lehet találni egy olyan r sort, amelynek az a és b oszlopa különbözők. Tegyük fel, hogy az r sorban az a értéke 0, míg a b értéke 1.

Ekkor a b helyes értékét kiszámolhatjuk úgy, hogy vesszük a b -n kívül az összes olyan lemezt, amelynél egyes szerepel az r sorban, és ezeknek a lemezeknek a megfelelő biteit modulo-2 összeadjuk. Vegyük észre, hogy az a nem szerepel ezek között a lemezek között, így nem fordulhat elő, hogy nem tudjuk elvégezni a modulo-2 összeadást. Miután ezt megtettük, ki kell számolnunk az a -t is a többi lemez felhasználásával. Mivel a 2.17. ábra mátrixának minden oszlopban legalább egy egyes szerepel, ezért vehetünk egy olyan sort, amelynek az a oszlopában egyes szerepel. Ha most vesszük az összes a -tól különböző lemezt, amelyre egyes szerepel ebben a sorban, és ezeknek a megfelelő biteit modulo-2 összeadjuk, akkor sikeresen újra kiszámolunk az a lemez tartalmát.

2.23. példa: Tegyük fel, hogy egyszerre tönkrement a 2. és az 5. lemez. A 2.17. ábra mátrixát megvizsgálva észrevehetjük, hogy az ennek a két lemeznek megfelelő oszlopok különbözőnek a 2. sorban, ahol a 2. lemeznél egyes szerepel, míg az 5. lemeznél nulla. Tehát rekonstruálhatjuk a 2. lemezt úgy, hogy vesszük az 1., 4. és 6. lemezeket, vagyis, ahol szintén egyes szerepel ebben a sorban, és ezeknek a lemezeknek a megfelelő biteit modulo-2 összeadjuk. Vegyük észre, hogy a fenti három lemez között nem szerepel tönkrement lemez. Folytassuk például a 2.19. ábrának megfelelő helyzetet az első blokkokra vonatkozóan, azaz a 2. és 5. lemez tönkremenése után legyenek adva kezdetben a 2.20. ábra adatai:

Ha most az 1., 4. és 6. lemezek blokkjainak tartalmát modulo-2 összeadjuk, akkor a második lemeznek erre a blokkjára a 00001111-et kapjuk, amelyről a 2.19. ábrán ellenőrizhetjük, hogy tényleg ez a helyes blokk.

Most vegyük észre, hogy a 2.17. ábrán az 5. lemez oszlopában az első sorban egyes szerepel, ezért újra kiszámolhatjuk az 5. lemezt úgy, hogy vesszük a többi olyan lemezt, melyre szintén egyes szerepel az első sorban, így kapjuk az 1., 2. és 3. lemezeket, aztán ezek megfelelő bitjeit modulo-2 összeadjuk. Az első blokkra így 11000111 lesz az összeg. A számítás helyességéről meggyőződhetünk a 2.19. ábra segítségével. □

Lemez	Tartalom
1)	11110000
2)	?? ?? ?? ??
3)	00111000
4)	01000001
5)	?? ?? ?? ??
6)	10111110
7)	10001001

2.20. ábra. A 2. és 5. lemez tönkremenése utáni helyzetet

2.6.6. Feladatok

2.6.1. feladat: Tegyük fel, hogy a 2.16. példához hasonlóan tükrözzük a lemezeket. A meghibásodási arány 4% évente. 8 órát vesz igénybe egy lemez cseréje. Mekkora a várható értéke az olyan lemezhibának, amely adatvesztéssel is jár?

*1 **2.6.2. feladat:** Tegyük fel, hogy a lemezeknek a meghibásodási rátája egy F törtszám évente és H óra kell egy lemez cseréjéhez.

a) Ha tükrözött lemezeket használunk, akkor F és H függvényében mennyi az adatvesztési idő várható értéke?

b) Ha 4. vagy 5. szintű RAID-sémát használunk N számú lemezzel, akkor mennyi az adatvesztési idő várható értéke?

Lemez	Tartalom
1)	11110000
2)	00001111
3)	00111000
4)	01000001
5)	?? ?? ?? ??
6)	10111110
7)	10001001

2.21. ábra. A 2. lemez helyreállítása után

!! **2.6.3. feladat:** Tegyük fel, hogy három lemezt használunk tükrözött csoportként, azaz mindhárom lemez azonos adatokat tartalmaz. Tegyük fel, hogy egy lemez meghibásodási rátája egy F törtszám évente és H óra kell egy lemez helyreállításához. Az F és H függvényében mennyi az adatvesztési idő várható értéke?

További észrevételek a 6. szintű RAID-del kapcsolatban

1. Az 5. és 6. szintű RAID elvét össze is kombinálhatjuk, azaz a redundáns lemezeket a blokk vagy cilinder sorszáma szerint változtatjuk. Ha így tesszük, akkor elkerüljük azt az írásnál jelentkező problémát, hogy a redundáns lemezeket többet használjuk, mint az adatlemezeket, ugyanis a 2.6.5. részben leírt sémának a szűk keresztmetszetét a redundáns lemezek használata jelenti.

2. A 2.6.5. részben leírt séma nem csak négy lemezre használható. A lemezek száma egy 2-hatványnál eggyel kisebb szám lehet, azaz $2^k - 1$. Ekkor a lemezek közül k darab redundáns és a többi, azaz $2^k - k - 1$ darab pedig adatlemez. Emiatt a redundancia durván a lemezek számának logaritmusával arányosan nő. Tetszőleges k -ra egy 2.17. ábrának megfelelő mátrixot készíthetünk úgy, hogy vesszük az összes nullát és egyest tartalmazó lehetséges k dimenziós oszlopot, kivéve a csupa nullából álló oszlopot. Azok az oszlopok, amelyek egyetlen egyest tartalmaznak a redundáns lemezeknek felelnek meg, az egy-nél több egyest tartalmazó oszlopok pedig az adatlemezeknek.

2.6.4. feladat: Tegyük fel, hogy 4. szintű RAID-sémát használunk négy adatlemezzel és egy redundáns lemezzel. A 2.17. példához hasonlóan tegyük fel, hogy a blokkok mérete egy bájt. Adjuk meg a redundáns lemez blokkját, ha a megfelelő blokkok az adatlemezeken a következők:

- * a) 01010110, 11000000, 00111011 és 11111011.
- b) 11110000, 11111000, 00111111 és 00000001.

2.6.5. feladat: Használjuk ugyanazt a 4. szintű RAID-sémát, mint a 2.6.4. feladatban, és tegyük fel, hogy tönkremegy az 1. lemez. Hozzuk helyre az elromlott lemez blokkját a következő körülmények mellett:

- * a) A 2., 3. és 4. lemez tartalma 01010110, 11000000 és 00111011, a redundáns lemez tartalma pedig 11111011.
- b) A 2., 3. és 4. lemez tartalma 11110000, 11111000 és 00111111, a redundáns lemez tartalma pedig 00000001.

2.6.6. feladat: Tegyük fel, hogy a 2.6.4. feladatban az első lemez blokkját 10101010-ra változtatjuk. Milyen változtatást kell tenni a többi lemez megfelelő blokkjain?

2.6.7. feladat: Tegyük fel, hogy a 2.22. példában szereplő 6. szintű RAID-sémával dolgozunk. A négy adatlemez blokkjai rendre 00111100, 11000111, 01010101 és 10000100.

a) Mik a redundáns lemezek megfelelő blokkjai?

Hibajavító kódok és a 6. szintű RAID

A redundáns lemezek tartalmának meghatározásához egy megfelelő mátrixot kell választanunk, amely olyan, mint a 2.17. ábra mátrixa. A mátrix kiválasztásában egy alaposan kidolgozott elmélet nyújt segítséget. Egy n hosszú bitvektorból (*kódszavakból*) álló halmazt n hosszú *kódnak* hívunk. Két kódszó közti *Hamming-távolság* az a szám, ahány pozícióban a két kódszó különbözik. Egy kód *minimális távolsága* két tetszőleges, különböző kódszó közti legkisebb Hamming-távolság.

Legyen C egy tetszőleges n hosszú kód. Megkövetelhetjük, hogy n lemez megfelelő bitjei olyan sorozatokat alkossanak, amelyek a kód elemei. Vegyünk egy egyszerű példát, melyben egy lemezt és a tükröképét használjuk, azaz $n = 2$. Ekkor használhatjuk a $C = \{00, 11\}$ kódot, mert ez azt jelenti, hogy a két lemezek a megfelelő bitjei meg kell hogy egyezzenek. Egy másik példaként nézzük a 2.17. ábra mátrixát. Ez a mátrix egy olyan kódot definiál, amely 16 darab 7 hosszú bitvektorból áll, melyek első négy bitje tetszőlegesen választható, de a maradék három bitet a három redundáns lemezre vonatkozó szabály határozza meg.

Legyen egy kód minimális távolsága d . Legyenek a lemezeink olyanok, melyek tartalmára igaz, hogy a lemezek megfelelő bitjeiből álló sorozatok mindig a kódnak valamelyik vektorával egyeznek meg. Ekkor ez a rendszer $d - 1$ lemez egyidejű meghibásodását is helyre tudja hozni. Ennek az az oka, hogy ha egy kódszóban kiürültünk $d - 1$ helyet, és feltesszük, hogy kétféleképpen is ki lehetne tölteni ezeket a pozíciókat úgy, hogy kódszavakat kapjunk, akkor az így kapott két kódszó legfeljebb $d - 1$ pozícióban különbözne, de akkor nem lehetett volna d a kód minimális távolsága. Például vegyük a 2.17. ábra mátrixát, mely a jól ismert *Hamming-kódot* definiálja, és amelynek 3 a minimális távolsága. Emiatt ezzel két lemez meghibásodását lehet kezelni.

b) Ha a harmadik lemez blokkját átírjuk 10000000-ra, akkor milyen lépéseket kell tennünk a többi lemez megváltoztatásának érdekében?

2.6.8. feladat: Legyen adott egy 6. szintű RAID-séma hét lemezzel. Írjuk le, hogy milyen lépéseket kell végrehajtani a helyreállításához, ha a következő lemezek sérültek meg:

- * a) 1. és 7. lemez.
- b) 1. és 4. lemez.
- c) 3. és 6. lemez.

2.6.9. feladat: Keresünk olyan 6. szintű RAID-sémát, amely 15 lemezre használható. A 15 lemezből 4 redundáns lemez. *Segítség:* Általánosítsuk a 7 lemezes Hamming-mátrixot.

2.6.10. feladat: Adjuk meg a 7 hosszú Hamming-kódnak mind a 16 kódszavát, azaz melyik az a 16 megengedett bitsorozat, amely előfordulhat 7 lemez megfelelő bitsorozatként, ha a 2.17. ábra mátrixán alapuló 6. szintű RAID-sémát használjuk?

2.6.11. feladat: Tegyük fel, hogy négy lemezünk van, az első kettő adatlemez, a 3. és 4. pedig redundáns. A 3. lemez az 1. lemez tükröképe. A 4. lemez a 2. és 3. lemez megfelelő bitjeinek paritás-ellenőrző bitjeit tartalmazza.

a) Fejezzük ki ezt a helyzetet azzal, hogy megadjuk egy 2.17. ábrához hasonló paritás-ellenőrző mátrixot.

!! b) Bizonyos esetekben, de nem az összes lehetséges esetben, ez a rendszer is képes a helyreállításra, ha két lemez egyszerre megy tönkre. Határozzuk meg, hogy milyen párokra lehetséges a helyreállítás, milyen párokra nem lehetséges.

***2.6.12. feladat:** Tegyük fel, hogy nyolc adatlemezünk van, 1-től 8-ig számozva, és három redundáns lemezünk, a 9., 10. és 11. lemez. A 9. lemez az első négy adatlemez paritás-ellenőrzése, a 10. lemez pedig a további négy lemez paritás-ellenőrzése. Tegyük fel, hogy bármelyik két lemez egyforma valószínűséggel megy egyszerre tönkre. Mely lemezeknek legyen a 11. lemez a paritás-ellenőrzése, ha maximalizálni akarjuk annak a valószínűségét, hogy két lemez egyidejű meghibásodását is helyre tudjuk hozni?

!! **2.6.13. feladat:** Adjunk meg egy tíz lemezre vonatkozó 6. szintű RAID-sémát, mely-lyel lehetséges a helyreállítás akkor is, ha egyszerre három lemez megy tönkre. Használjunk annyi adatlemezt, amennyit csak lehet.

2.7. Összefoglalás

- *Memóriahierarchia:* Egy számítógépes rendszer többféle tárolóelemet használ. Ezek a sebesség, kapacitás és egy bite jutó költség tekintetében különböző nagyságrendűek lehetnek. A legkisebbtől/legdrágábbtól a legnagyobb/big/legolcsóbbig a sorrend a következő: cache, központi memória, másodlagos memória (lemez), harmadlagos memória.
- *Harmadlagos tárolás:* Az alapvető harmadlagos tárolóeszközök a következők: szalagos kazetták, szalagsírlók (szalagos kazettákat kezelő mechanikus eszközök), lemeztárak vagy „juke box”-ok (CD-ROM-lemezeket kezelő mechanikus eszközök). Ezeknek a tárolóeszközöknek a kapacitása sok terabájt, de ezek a leglassabb tárolóeszközök.
- *Lemezek/másodlagos tárolás:* A másodlagos tárolóeszközök alapvetően sok gigabájt kapacitással rendelkező mágneses lemezek. A lemezegységek több kör alakú, mágneses anyaggal bevont táánytérből állnak, melyek koncentrikus sávokban tárolják a biteket. A táánytér egy közepén elhelyezett tengely körül forognak. A középponttól azonos távolságra elhelyezkedő sávok alkotnak egy cilindert.

- **Blockok és szektorok:** A sávok szektorokra vannak felosztva, melyeket nem magnezett hézagok választanak el egymástól. A szektor a lemezről olvasás, illetve lemezre írás egysége. A blokkok a tárolás logikai egységei, melyeket alkalmazások, például adatbázis-kezelő rendszerek használnak. A blokkok tipikusan több szektorból állnak.
- **Lemezvezérlő:** A lemezvezérlő egy olyan processzor, mely egy vagy több lemezegységet vezérel. A lemezvezérlő felelős azért, hogy mikor olvassni vagy írni akarunk egy sávot, akkor a lemezfeleket a megfelelő cilinderré vigye. Feladata lehet még a versenyző igények ütemezése, és azoknak a blokkoknak a pufferezése, amiket írunk vagy olvasunk.
- **Lemezlelévi idő:** Egy lemez kérése az az idő, amely egy blokk olvasására vagy írására vonatkozó igény beérkezése és a blokkelérés végrehajtása között telik el. A kérést alapvetően három tényező okozza: a keresési idő, amely ahhoz kell, hogy a fejek a megfelelő cilinderrig eljussanak, a rotációs késés, amely ahhoz kell, hogy a kívánt blokk a lemez forgása közben a fej alá kerüljön, és végül az átviteli idő, amely ahhoz kell, hogy a fej alatti blokkot olvassuk vagy írjuk.
- **Moore törvénye:** Egy valóságággal összhangban lévő irányzat szerint az olyan paraméterek, mint a processzor sebessége, a lemez kapacitása és a központi memória kapacitása minden 18 hónapban megduplázódik. Ezzel szemben hasonló időszak alatt a lemezlelévi idő alig csökken, ha egyáltalán változik. Ennek egy fontos következménye, hogy a lemezlelés (relatív) költsége az évek során növekszik.
- **Másodlagos tárolókat használó algoritmusok:** Ha az adatmennyiség túl nagy, akkor nem fér el a központi memóriában. Ekkor olyan algoritmusokat kell használni az adatok kezelésére, melyek figyelembe veszik, hogy a lemez és memória között lezajló műveletek, lemezblokk olvasása, írása gyakran sokkal tovább tart, mint amennyi a központi memóriában szükséges ahhoz, hogy a memóriába bekerült adatokkal műveleteket végezzünk. A másodlagos memóriában tárolt adatokra vonatkozó algoritmusok értékelésénél ezért elsődlegesen azt kell vizsgálni, hogy mennyi lemez I/O-műveletet igényel az algoritmus.
- **Kétfázisú, többitas, összefésülő rendezés:** Ez a rendező algoritmus lemezen tárolt hatalmas adatmennyiséget képes rendezni. Minden adathoz csak két lemezolvasást és két lemezírást használ. A legtöbb adatbázis-alkalmazásban ezt a rendezési módszert választják.
- **A lemezlelés gyorsítása:** Több technika is használható arra, hogy bizonyos alkalmazásokhoz a lemezblokkokat gyorsabban lehessen elérni. Ezek közül a következőket emeljük ki: szétszórjuk az adatokat több lemezre (ezzel lehetőségessé válik a párhuzamos elérés), tükrözzük a lemezeket (az adatok több példányának kezelése szintén megengedi a párhuzamos hozzáférést), azonos sávra vagy cilinderré helyezjük azokat az adatokat, amelyeket együtt akarunk elérni, korai beolvasást, dupla pufferezést használunk, azaz együtt olvasunk vagy írunk teljes sávokat vagy cilindreket.
- **Lift algoritmus:** Azzal is fel lehet gyorsítani az elérést, ha az elérési igényeket sorba állítjuk, és olyan sorrendben dolgozzuk fel őket, hogy a fejek mindig oda-vissza, végigmegyjenek a teljes lemezen. Az igények kezelésére a fejek mindig megállnak, ha elérnek egy olyan cilinderré, amely egy vagy több olyan blokkot tartalmaz, amelyre várakozási listán szereplő igény vonatkozik.

- **Lemezhibák típusai:** A rendszereknek kezelni kell tudniuk a hibákat azért, hogy elkerüljék az adatvesztést. Az alapvető lemezhibatípusok a következők: ideiglenes (ha elégszer megismételjük a műveletet, akkor egy idő után az írási vagy olvasási hiba nem fog előfordulni), állandó (adatok romlottak el a lemezen, és nem lehet helyesen olvasni őket), a lemez tönkremenése (a teljes lemez olvashatatlanná vált).
- **Ellenőrző összegek:** Ha egy plusz paritás-ellenőrző bittel egészítjük ki a bitsorozatunkat úgy, hogy a bitsorozatban szereplő egyesek számát az extra bit párossá teszi, akkor az ideiglenes és állandó hibákat felismerhetjük, bár nem tudjuk őket kijavítani.
- **Stabil tárolás:** Minden adatból két másolatot készítünk, és ügyelünk arra, hogy milyen sorrendben írjuk ezeket a másolatokat. Ezáltal egyetlen lemez használható arra, hogy kivédjünk majdnem minden, egy szektorra vonatkozó állandó hibát.
- **RAID:** Többféle séma létezik arra vonatkozóan, hogy lehet egy vagy több lemez hozzáadásával elérni, hogy az adatok túlélhessenek egy lemeztönkremenést. Az 1. szintű RAID a lemezek tükrözését jelenti. A 4. szintű RAID egy plusz lemez hozzáadását jelenti, amely az összes többi lemez megfelelő bitjeinek paritás-ellenőrzésait tartalmazza. Az 5. szintű RAID változatja, hogy a paritás bit mikor melyik lemezre kerüljön, ezáltal megóv attól, hogy csak egyetlen paritáslemez legyen, ami az írás stűk keresztmetszetét jelentené. A 6. szintű RAID már hibajavító kódok használatát is magában foglalja, és lehetővé teszi, hogy az adatok több lemez egyidejű tönkremenését is túléljék.

2.8. Irodalomjegyzék

- A RAID-elv a [6]-ban szereplő lemezávozásra vezethető vissza. A hibajavító képesség neve az [5]-ből ered.
- A 2.5. részben szereplő lemezmegehibásodási modell Lampson és Sturgis meg nem jelent [4] munkájában olvasható.
- A fejezet lényeges elemeihez több hasznos áttekintés is létezik. A [2] a lemeztárolók és hasonló rendszerek irányzatait vizsgálja. A RAID-rendszerek összefoglalása az [1]-ben található. A [7] azokat az algoritmusokat tekinti át, melyek a számítás során másodlagos tárolási modellel (blokkmodellel) használnak.
- A [3] egy fontos tanulmány arról, hogy egy bizonyos feladat végrehajtásához hogyan lehet optimalizálni egy processzorral, memóriával és lemezzel ellátott rendszert.
1. P. M. Chen et al., „RAID: high-performance, reliable secondary storage”, *Computing Surveys*, 26:2 (1994), pp. 145–186.
 2. G. A. Gibson et al., „Strategic directions in storage I/O-issues in large-scale computing”, *Computing Surveys*, 28:4 (1996), pp. 779–793.
 3. J. N. Gray and F. Putzolo, „The five minute rule for trading memory for disk accesses and the 10 byte rule for trading memory for CPU time”, *Proc. ACM SIGMOD Intl. Conf. on Management of Data*, (1987), pp. 395–398.