

## Käyttöjärjestelmät II

LUENTO 10

### TIEDOSTOJEN HALLINTA Käytännön esimerkit

Ch 12.8-9 [Stal 05]  
Ch 10.6.4, 11.6-7 [Tane 01]  
Ch 20.7 [DDC 04]

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola

7.4.2006

1

## Mitä KJ-I:ssä / KJ-II:ssa?

### KJ-I

- n Tiedostojen organisointi, hakemistot
- n Tiedostojen yhteiskäyttö, tietueet ja lohkot
- n Levytilan hallinta
- n UNIX: tiedostojärjestelmä

### Seuraavaksi KJ-II:ssa [Stal 05] [Tane 01] [DDS 04]

- n Linux
  - u Virtual File System (Ch. 12.8 [Stal 05])
  - u ext2fs (Ch 6, Ch 11.6 [Tane 01], Ch 2.7 [DDS 04])
  - u NFS, Network File System (Ch 10.6 [Tane 01])
- n Windows
  - u Journaling File System
  - u NTFS - W2K File System (Ch 12.9 [Stal 05], Ch 11.7 [Tane 01])

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola

7.4.2006

2

## Käyttöjärjestelmät II



# LINUX

## Tiedostojärjestelmät

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola

7.4.2006

3

## Linux

- n **Tiedosto = tavujono**
  - u ei tietueita, ei jaksoja
  - u organisointi sovelluksissa
- n **Tiedostonimi ja attribuutit erillään** Fig. 6-16 [Tane 01]
  - u attribuutit = i-solmu (i-node, index node)
- n **Hakemisto**
  - u tiedosto, jossa pareja (tiedostonimi, i-solmunumero)
- n **Symbolinen linkki (soft)**
  - u tiedosto, jossa tiedoston polkunimi

Bovet D.P., Cesati M.:  
*Understanding the LINUX KERNEL*. O'Reilly, 2001.



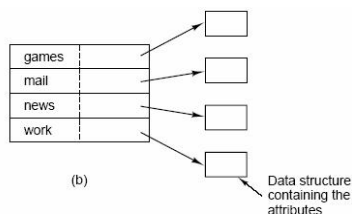
KJ-II K2006 / Auvo Häkkinen - Teemu Kerola

7.4.2006

4

games	attributes
mail	attributes
news	attributes
work	attributes

(a)



(b)

Fig. 6-16. (a) A simple directory containing fixed-size entries with the disk addresses and attributes in the directory entry. (b) A directory in which each entry just refers to an i-node.

[Tane01]

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola

7.4.2006

5

## Tiedostojen yhteiskäyttö

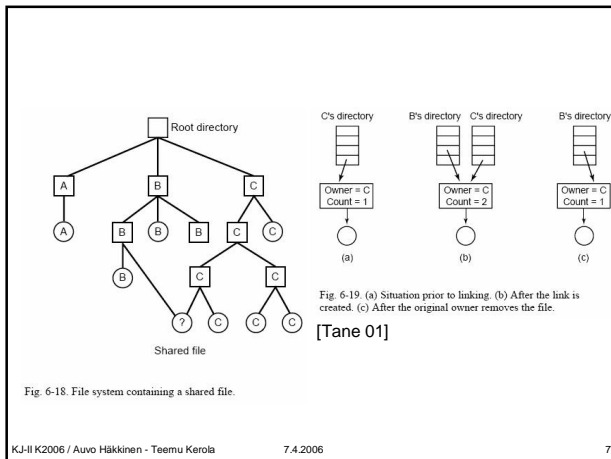
Fig. 6-18 [Tane 01]

- n **Sama tiedosto käytössä monta kertaa**
- n **Hard link** Fig. 6-19 [Tane 01]
  - u monta omistajaa (vai monella omistajan oikeudet?)
  - u kaikilla samat oikeudet
    - F omistaja poistaa → muilla silti käytössä normaalisti
- n **Soft link eli symbolinen linkki**
  - u tiedoston tyyppi: symbolinen linkki
  - u tiedoston sisältö: merkkijono, joka indikoi varsinaisen tiedoston
  - u vain yksi omistaja
    - F omistaja poistaa → muiden linkit epäkelvoja

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola

7.4.2006

6



## Linux VFS (virtual file system)

- n **Tuki useille tiedostojärjestelmille** Fig. 12.15 [Stal 05]
  - u register\_filesystem(), unregister\_filesystem()
  - u ext2fs, procs, FAT, NTFS, minix, NFS, smb, ... Fig. 12.16 [Stal 05]
  - u superlohko
    - F määrittelee tiedostojärjestelmän
    - F oma paikka levyllä
- n **Tarjoaa sovelluksille yhtenäisen rajapinnan**
  - u open(), read(), write(), seek(), close(), ...
  - u kaikki VFS:n kautta viitattu tieto **ei ole** levyllä talletettuja tiedostoja
    - F i-node
    - F KJ-oliot, laitteet, yhteiset muistialueet
    - F tiedon suojaus siitti tiedostojen tapaan



KJ-II K2006 / Auvo Häkkinen - Teemu Kerola 7.4.2006 8

## Linux VFS rakenne

Ch. 20.7 [DDC 04]

- n **VFS fd – file descriptor (file object)**
  - u aukiolevalle tiedostolle, missä kohtaa lukemassa/kirjoittamassa
  - u oikeudet (user ID, group ID)
  - u linkki tiedoston *dentry*yn
  - u sallitut operaatiot
- n **VFS dentry – directory entry**
  - u osoittaa hakemistopuussa kaikkiin lähisukulaisiin (niiden *dentry*)
    - F isä-hakemisto
    - F lapsi-hakemistot tai –tiedostot
    - F sisarus-hakemistot tai –tiedostot
  - u linkki tiedoston i-node:en
- n **VFS i-node (varsinainen tiedoston metatieto)**
  - u tiedostojärjestelmän tunniste ja superlohko (*superblock*)
  - u tiedostojärjestelmän sisäinen *i-node*



KJ-II K2006 / Auvo Häkkinen - Teemu Kerola 7.4.2006 9

## VFS metatiedon välimuistit

- n **VFS ja hakemistohierarkia hidastavat käyttöä**
  - u tiedot pitää yleensä hakea levyltä hakemisto kerrallaan
  - u tiedostojärjestelmäkohtainen *lookup()*
- n **dcache (dentry cache)**
  - u viimeksi viitattujen tiedostojen *dentry*
  - u nopea kuvaus *filename* → *i-node*
  - u tiedoston X *dentry* välimuistissa → myös kaikki tiedoston X esivanhempien *dentry* välimuistissa
- n **i-node cache**
  - u viimeksi käytössä olleiden tiedostojen *VFS i-nodet*
    - F näistä löytyy tiedostojärjestelmän *i-node*

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola 7.4.2006 10

## Linux tiedostojärjestelmät



- n **ext2fs (second extended file system)**
  - u Linuxia varten kehitetty tiedostojärjestelmä
  - u esikuvana BSD Fast File System (FFS)
    - F lohkokoryhmat
  - u tehokkuus, luotettavuus
- n **/proc**
  - u erikoistiedostot, luodaan 'lennosta'
  - u esim. ytimen parametrien kysely/asettaminen
  - u KJ-palvelut piilotettu tiedostojärjestelmän käytöksi
    - F käytön valvonta tiedostojärjestelmän suojauksen avulla
- n **ext3fs**
  - u journaling file system, log-structured file system (LFS)
  - u Red Hat Linux'issa

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola 7.4.2006 11

## Looginen levy – yleinen tapaus

- n **MBR (master boot record)** Fig. 6-11 [Tane 01]
  - u fyysinen sektori 0, jonka BIOS lukee
    - F Basic Input Output System Flash BIOS
    - F mitä tehdään ennen alustusta tai miten alustetaan
  - u partitiotaulu
    - F kunkin partition alku ja loppu
    - F tiedostojärjestelmän tyyppi
    - F yksi partitioista aktiivinen → bootti
      - voidaan valita alustuksen yhteydessä?

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola 7.4.2006 12

## Linux ext2fs levy

General: Fig. 6-11 [Tane 01]

Linux: Fig. 10-35 [Tane 01]

- n Lohkoryhmät (block groups)
  - n yhtenäisesti levyttä allokoitu alue
  - n datalohkot ja i-nodet fyysisesti lähellä toisiaan
  - n säästä hakuvarren siirroissa
- n Kaikki lohkot samankokoisia (1 KB)
- n Kaikki i-nodet 128B (tavallinen UNIX 64B)

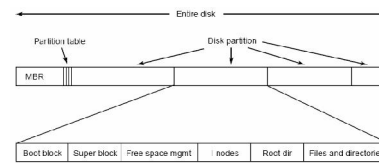


Fig. 6-11. A possible file system layout.

[Tane01]

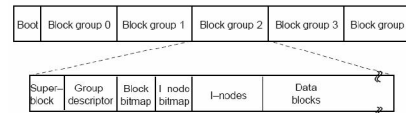


Fig. 10-35. Layout of the Linux Ext2 file system.

## ext2fs superlohko (superblock)

- n 1 lohko Fig. 10-35 [Tane 01]
- n Kuvaa koko ext2fs-partition rakenteen
- n Kopio jokaisen lohkoryhmän alussa
  - n luotettavuus, virheestä toipuminen
- n Ydin operoi vain lohkoryhmän 0 superblokilla ja ryhmäkuvaajilla
  - n muille käyttöä, jos superblock 0 'rikki'
  - n `/sbin/e2fsck` kopioi aika-ajoin muualle



## ext2fs superlohko (superblock)

	0	1	2	3	4	5	6	7
0	Number of i-nodes			Number of blocks				
8	Number of reserved blocks			Number of free blocks				
16	Number of free i-nodes			First data block				
24	Block size			Fragment size				
32	Blocks per group			Fragments per group				
40	i-nodes per group			Time of mounting				
48	Time of last write			Status		Max. mnt cnt		
56	Ext2signal.		Status		Error behav.		Pad word	
64	Time of last test			Max test interval				
72	Operating system			File system revision				
80	RESUID		RESGID		Pad word			
blocksize:	Pad words							

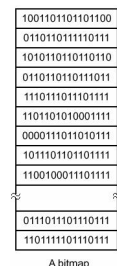
## ext2fs ryhmäkuvaaja (group descriptor)

- n lohkoa Fig. 10-35 [Tane 01]
  - n tietoa kaikista lohkoryhmistä
    - n käyttöä valittaessa ryhmää varaukselle
    - n ei perustu sylintereihin
  - n kopio jokaisen lohkoryhmän alussa
  - n yksi kuvaaja 24 B
- |                                   |                                   |
|-----------------------------------|-----------------------------------|
| <code>bg_block_bitmap</code>      | lohkobbittikartan lohkonro (4 B)  |
| <code>bg_inode_bitmap</code>      | i-solmubittikartan lohkonro (4 B) |
| <code>bg_inode_table</code>       | i-solmutaulun lohkonro (4 B)      |
| <code>bg_free_blocks_count</code> | vapaiden lohkojen lkm (2 B)       |
| <code>bg_free_inodes_count</code> | vapaiden i-solmujen lkm (2 B)     |
| <code>bg_used_dirs_count</code>   | hakemistojen lkm ryhmässä (2 B)   |
| <code>bg_pad, bg_reserved</code>  | tyhjää (6 B)                      |



## ext2fs lohkoryhmän bittikartat

- n 2 bittikarttaa Fig. 10-35 [Tane 01]
- n vapaat lohkot
- n vapaat i-nodet
- n Molemmissa 8192 bittiä (1 KB lohko)
- n 0 = vapaa, 1 = varattu



## ext2fs lohkokryhmän i-node taulu

- n lohkoa
- n *i-nodet* á 128 B
- n tiedoston attribuutit
- n tiedoston lohkojen numerot
  - n 12 suoraa viitettä lohkoihin, viite 4B
  - n 3 epäsuoraa viitettä lohkoihin
- n Esimerkki
  - n 1 KB bittikartta → ryhmässä 8192 lohkoa tai i-nodea
  - n i-noden lohkoviiite 13225
    - n lohkokryhmä 1, siirtymä 5033 (=13225-8192)
  - n root-hakemiston i-node = *i-node* #2

Fig. 10-35 [Tane 01]



## ext2fs i-node

Fig. 10-35 [Tane 01]

0	Mode	Uid	File size
8	Access time		Time of creation
16	Time of modification		Time of deletion
24	Gid	Link counter	No. of blocks
32	File attributes		Reserved (OS-dependent)
40	12 direct blocks		
88	One-stage indirect block	Two-stage indirect block	
96	Three-stage indirect block	File version	
104	File ACL	Directory ACL	
112	Fragment address		Reserved (OS-dependent)
120	Reserved (OS-dependent)		

Access Control List



## ext2fs datalohkot (data blocks)

- n Lohkon koko 1 KB
  - n 2 KB, 4 KB tai 8 KB?? ei.
- n Suurin tiedostokoko 2 GB
  - n kenttä "file size" rajoittaa (ylin bitti ei käytössä!)
  - n jos 64-b kone, max 4TB
- n Tiedosto voi jakautua useamman ryhmän alueelle

Fig. 10-35 [Tane 01]



## ext2fs hakemistoalkio

4 B      2 B    1 B    1 B                      1 - 255 B  
(1 - EXT2\_NAME\_LEN)

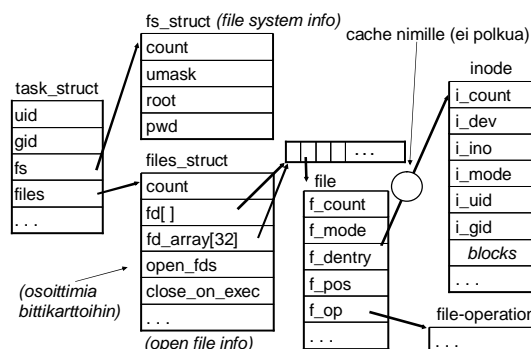
i-node number	entry len	name len	type	file name
---------------	-----------	----------	------	-----------

- n Hakemisto on tiedosto, joka kuvaa tiedostonimet *i-node*iksi
  - n peräkkäinen lista hakemistoalkioita (directory entry)
  - n voi olla myös B-puu, jos paljon tiedostoja
- n Hakemistoalkio on vaihtelevanpituinen
  - n pituus aina 4:n monikerta (lopussa /0-merkkejä)
- n Tyyppi
  - n 0 = tuntematon, 1 = tavallinen tiedosto, 2 = hakemisto
  - n 3 = merkkilaitte, 4 = lohkolaitte, 5 = nimetty putki
  - n 6 = pistoke, 7 = symbolinen linkki

Esim: TKTL mail server

Polku tässä tiedostossa tai hakemistoalkiossa (fast symbolic link)

## Linux: tiedoston käyttö (vfs)



## Linux: tiedostojen käyttö (vfs)

file\_operations

llseek( file, offset, whence)
read( file, buf, count, offset)
write( file, buf, count, offset)
readdir( dir, dirent, filldir)
poll( file, poll_table)
ioctl( inode, file, cmd, arg)
map( file, vma)
open( inode, file)
flush( file)
release( inode, file)
fsync( file, dentry)
fasync( file, on)
check_media_change(dev)
revalidate(dev)
lock( file, cmd, file_lock)

- n Jokaisella tiedostojärjestelmällä omat funktiot
- n File\_operations rakenteessa funktion osoite
- n Jos ei toteuta kyseistä operaatiota, osoitin NULL

Fig. 6-5 [Tane 01]



## Linux procfs tiedostojärjestelmä

- n **Process file system**
- n **Ei todellinen (fyysinen) tiedostojärjestelmä**
  - u kaikki keskusmuistissa, levyllä ei tiedostoja
- n **Käyttöliittymä prosessikuvaajiin**
  - u hakemistossa /proc
  - u jokainen /proc'in alihakemisto määrittelee omat read() ja write() operaationsa
    - F /proc/4321 on prosessin 4321 hakemisto
  - u KJ-tietojen lukeminen ja kirjoittaminen
  - u read() ja write() toteuttavat suojatun tietorakenteen
    - F käytön valvonta tiedostojärjestelmän avulla
    - F samanaikaisuuden hallinta

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola

7.4.2006

25

## Linux sysfs

- n **hakemisto /sys**
- n **käyttöliittymä laitekuvaajiin**
  - u unified device model
- n **väylät hakemistossa /sys/bus**
  - u pci laitekuvaaja hakemistossa /sys/bus/pci
- n **I/O laitteet laitetyypin mukaan**
  - u /sys/class/input
  - u laitetyypin nimi, numero, laitteet, ajurit
- n **pidetään kirjaa kaikista laitteista, jotka käytössä ja missä ne ovat**
- n **pollataan aika ajoin väyliä, jos uusia laitteita tulisi tai vanhoja poistuisi**
  - u hot swappable devices

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola

7.4.2006

26

## Käyttöjärjestelmät II

### NFS Network File System

Ks. esim. Ch 10.6.4 [Tane 01]

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola

7.4.2006

27

## NFS

- n **Etäkoneiden hakemistojen liittäminen omaan hakemistopuuhun**
  - u kehittäjä Sun Microsystems
- n **NFS-protokolla**
  - u pyyntö-vastaus protokolla
  - u ei ota kantaa siihen kuinka toteutetaan
    - F NFS-palvelija, NFS-asiakas
- n **Windowsin vastine SMB-protokolla**
  - u Server Message Block

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola

7.4.2006

28

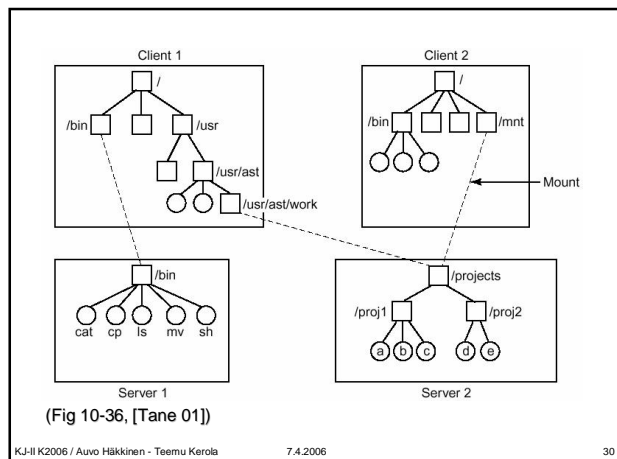
## NFS-arkkitehtuuri

- n **Etäkone (Palvelija)**
  - u suorittaa NFS-palvelijaa Fig. 10-36 [Tane 01]
  - u määrittelee hakemiston julkiseksi hakemistossa /etc/exports
    - F mm. käyttöoikeuksien rajaaminen
- n **Asiakas**
  - u suorittaa NFS-asiakasprosessia
  - u asemoi ("mounttaa") hakemiston omaan hakemistopuuhun
  - u mount-point määritelty tiedostossa /etc/fstab
- n **VFS huomaa milloin viitataan toisessa koneessa olevaan (mountattuun) tiedostoon**
  - u välitä pyyntö palvelijalle
  - u palvelijan tiedostojärjestelmä ei ole tärkeä

KJ-II K2006 / Auvo Häkkinen - Teemu Kerola

7.4.2006

29



KJ-II K2006 / Auvo Häkkinen - Teemu Kerola

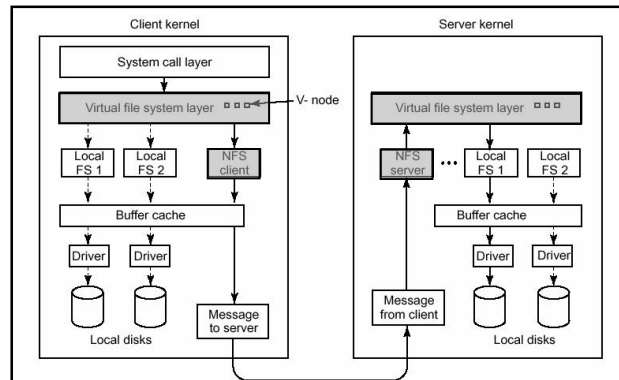
7.4.2006

30

## NFS-protokolla

Fig. 10-37 [Tane 01]

- n **mount**
  - u asiakas lähettää polkumimen palvelijalle
  - u palvelija palauttaa kahvan (file handle)
    - F tiedostojärj. tyyppi, laite#, inode#, oikeudet
    - F käytetään jatkossa kaikissa pyynnöissä
  - u voidaan tehdä alustuskripteissä (boot)
- n **automounting**
  - u mountataan, kun viitataan ens. kertaa
- n **pyynnöt**
  - u normaalit palvelupyynnöt sanomina
    - F read(), write(), ...
- n **tilaton protokolla**
  - u kaikki tarvittava tieto mukana pyynnössä
    - F kahva, lukupositio, paljonko, ...



(Fig 10-37, [Tane 01])

## Käyttöjärjestelmät II



### Windows 2000 Tiedostojärjestelmä (NTFS)

## LFS – Log-Structured File System

(Ch 6.3.8 [Tane 01])

- n "Uusi, parempi, fiksumpi, luotettavampi" – ihan totta!
- n **Usein tilanne**
  - u paljon päivityksiä, useimmat todelliset levyviitteet kirjoituksia
    - F luvut levyvälimuistista
  - u useimmat kirjoitukset pieniä päivityksiä
    - F levyn hakuvarsi liikkuu paljon, vähän dataa siirtyy
- n **Ongelma tavallisen tiedostojärjestelmän uuden tiedoston X luomisessa:**
  - u kirjoita hakemiston i-node, hakemisto, tiedoston i-node ja lopulta tiedosto
  - u virta poikki (rms vika) kesken kaiken? Oooops.
- n **Ratkaisu: tapahtumaloki, joka takaa tiedostojärjestelmän konsistenssisuuden**
  - u pidä lokia sekä metatiedosta (esim. inode) että itse datasta (sektorit)
- n **Journaling File System**
  - u pidä lokia vain metatiedosta – järjestelmä säilyy konsistenssina (data ei)
- n **Esim: Microsoft NTFS, Red Hat Linux ext3fs**

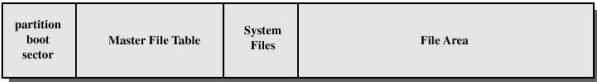
## LFS – alkuperäinen idea

- n **Koko levy on loki tapahtumista**
  - u uudet tapahtumat kirjoitetaan "loppuun" vapaaseen tilaan, peräkkäisiin lohkoihin
  - u nopeata, levyn täysi kapasiteetti hyödynnettävissä
- n **Uusi tiedosto X hakemistoon D**
  - u kirjoita X:n data-tapahtuma
  - u kirjoita X:n metadata (i-node?) -tapahtuma
  - u kirjoita X:n hakemistoalkio D:ssä –tapahtuma
- n **systemi koko ajan konsistenssissa tilassa**
- n **tiedon haku hidasta**
  - u ei niin paha, kun useimmat levytapahtumat kirjoittaa
  - u metadata (i-node) välimuistit: **i-node map, superblock**
  - u **cleaner** säie etsii tyhjää tilaa ja tiivistää lokeja

## NTFS: Piirteitä

- n **Kaatumisista ja levyvirheistä toipuminen**
  - u LFS lokitiedoston avulla
- n **Käyttöoikeudet**
  - u pääsilylistat (security descriptor)
- n **Sallii suuret levyt ja tiedostot**
  - u FAT32:ssa vain  $2^{32}$  lohkoa, suuri allokointitaulu
- n **Tiedosto-oliot ovat (arvo, attribuutti) -pareja**
- n **Mahdollisuus indeksointiin tiedoston käsittelyn nopeuttamiseksi**
- n **Lohko, cluster**
  - u yksi tai useampi peräkkäinen sektori (esim. 512 B - 4 KB)
    - F 32 GB levyllä 128 sektoria/lohko (→ lohko 64-512 KB)
  - u varauksen ja kirjanpidon perusyksikkö
- n **Partio, volume**
  - u fyysinen levyn looginen osa, jolla oma tiedostojärjestelmä

# NTFS-partitio



- n **Boottilohko** (Fig. 12.17 [Stal 05])
  - u partition ja tiedostojärj. rakenne, boottitietue ja -koodi
  - u MFT:n sijainti
- n **MFT**
  - u tietoa tiedostoista, hakemistoista (folders) ja vapaasta tilasta
- n **System Files (~ 1MB)**
  - u kopio MFT:n alkuosasta
  - u virheistäoipumisloki, bittikartta vapaat/varatut lohkot, attribuuttien kuvaustaulu
- n **File Area - tiedostojen lohkoille**

# NTFS – MFT

- n 1 KB:n kokoisia MFT-tietueita (Fig. 11-36 [Tane 01])
  - u jokainen kuvaa yhden taltiolla olevan tiedoston
    - F myös hakemisto on tiedosto
  - u vaihtelevanmittainen osa käytössä (Fig. 11-35 [Tane 01])
    - F (attribuutti, arvo) pareja (ei paikkasidonnainen!)
      - F data attribuutti, 'arvo' = lohkojen sijainti
- n 16 ensimmäistä tietuetta varattu ns. metadataalle (Fig. 11-34 [Tane 01])
- n Jos pieni tiedosto, tietue sisältää myös datan
- n Jos iso tiedosto, data erillisellä tallealueella
  - u MFT-tietuessa lohkonumeroita
  - u kuvaus voi jatkua useampaan MFT-tietueeseen (Fig. 11-37 [Tane 01])

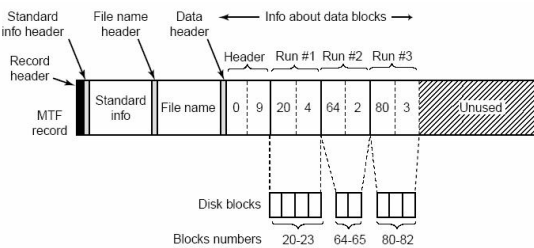


Fig. 11-36. An MFT record for a three-run, nine-block file. [Tane 01]

Attribute	Description
Standard information	Flag bits, timestamps, etc.
File name	File name in Unicode; may be repeated for MS-DOS name
Security descriptor	Obsolete. Security information is now in \$Extend\$Secure
Attribute list	Location of additional MFT records, if needed
Object ID	64-bit file identifier unique to this volume
Reparse point	Used for mounting and symbolic links
Volume name	Name of this volume (used only in \$Volume)
Volume information	Volume version (used only in \$Volume)
Index root	Used for directories
Index allocation	Used for very large directories
Bitmap	Used for very large directories
Logged utility stream	Controls logging to \$LogFile
Data	Stream data; may be repeated

Fig. 11-35. The attributes used in MFT records. [Tane 01]

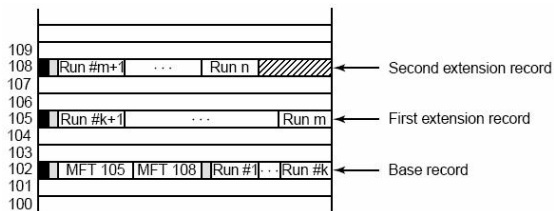
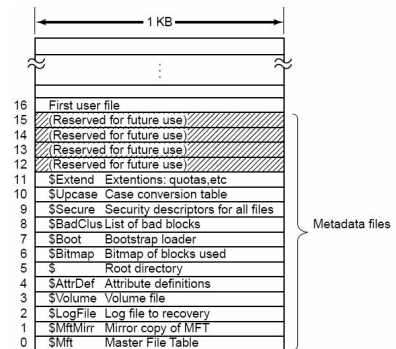
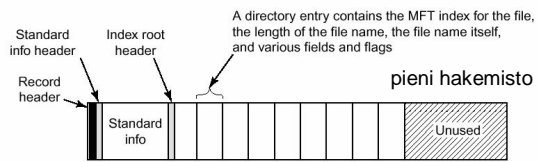


Fig. 11-37. A file that requires three MFT records to store all its runs. [Tane 01]



[Tane 01] Fig. 11-34. The NTFS master file table.

## Hakemiston MFT-tietue



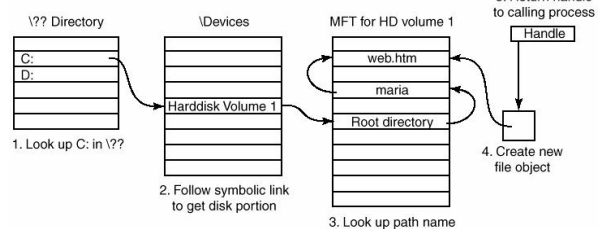
(Fig. 11-38 [Tane01])

- n Pienissä hakemistoissa MFT-tietueet peräkkäisjärjestyksessä
- n Isoissa hakemistoissa MFT-tietueissa B-puun (B-tree) indeksirakenne
  - u nimen etsintä ei ole peräkkäishakua

## Tiedoston käyttö

Fig. 11-39 [Tane 01]

- n `CreateFile("C:\maria\web.htm", ...)`
  - u etsi ensin oikea taltio
  - u juurihakemiston tietue on MFT:ssä, etsi juurihakemistosta alihakemiston tietue
  - u oliomanageri: luo uusi tiedosto-olio ja palauta kahva siihen
  - u prosessi käyttää kahvaa seuraavissa kutsuissa



(Fig. 10-39 [Tane 01])

## NTFS API

- n **NTFS WIN32 API vs Unix API** Fig 11-31 [Tane 01]

- n **Tiedoston kopiointi Win32 APIn avulla** Fig 11-32 [Tane 01]

Fig 6.5 [Tane 01]

Win32 API function	UNIX	Description
CreateFile	open	Create a file or open an existing file; return a handle
DeleteFile	unlink	Destroy an existing file
CloseHandle	close	Close a file
ReadFile	read	Read data from a file
WriteFile	write	Write data to a file
SetFilePointer	lseek	Set the file pointer to a specific place in the file
GetFileAttributes	stat	Return the file properties
LockFile	fcntl	Lock a region of the file to provide mutual exclusion
UnlockFile	fcntl	Unlock a previously locked region of the file

Fig. 11-31. The principal Win32 API functions for file I/O. The second column gives the nearest UNIX equivalent.

[Tane 01]

```

/* Open files for input and output. */
inhandle = CreateFile("data", GENERIC_READ, 0, NULL, OPEN_EXISTING, 0, NULL);
outhandle = CreateFile("newf", GENERIC_WRITE, 0, NULL, CREATE_ALWAYS,
FILE_ATTRIBUTE_NORMAL, NULL);

/* Copy the file. */
do {
    s = ReadFile(inhandle, buffer, BUF_SIZE, &count, NULL);
    if (s && count > 0) WriteFile(outhandle, buffer, count, &count, NULL);
} while (s > 0 && count > 0);

/* Close the files. */
CloseHandle(inhandle);
CloseHandle(outhandle);
    
```

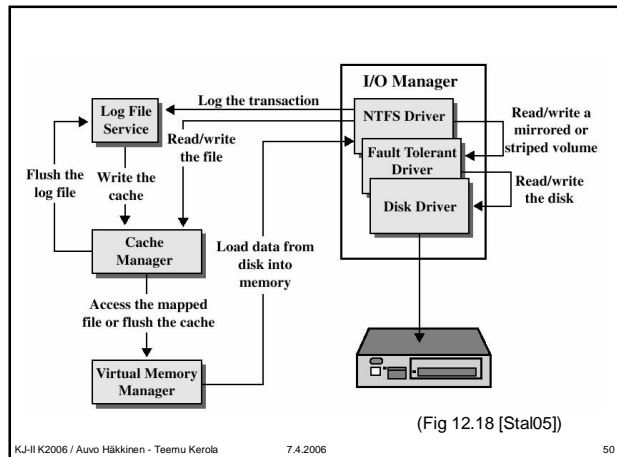
Fig. 11-32. A program fragment for copying a file using the Windows 2000 API functions.

[Tane 01]



## NTFS: Virheistä toipuminen

- n **Log FS** [Fig 12.18 [Stal 05]]
  - u kirjaa lokiin kaikki taltiota muuttavat transaktiot
  - u loki aluksi välimuistissa (vain kirjanpidon rakenteista)
- n **Muuta taltiota**
  - i. talleta lokitapahtuma tiedostovälimuistiin (file cache)
  - ii. tiedostomuutos välimuistiin
  - iii. talleta lokitapahtuma levyille välimuistista } "tapahtuma"
  - iv. talleta muutokset levyille välimuistista } "tapahtuma"
  - v. kommitoidu (**commit**)
- n Jos köllähtää ennenkuin muutokset levyllä, bootti voi palauttaa edeltävän tilanteen lokin avulla (**rollback**)
  - u ei takaa etteikö tiedostojen tietoa katoaisi
  - u järjestelmä säilyy eheänä (koherenttina)



## NTFS virheistä toipuminen

- n **Huono levysektori?**
  - u kirjoittamassa? kirjoita muualle OK
  - u lukemassa? *too bad*, data menetetty
  - u lukemassa master boot recordia tai boot sector'ia
    - F *really too bad*, taltio ehkä menetetty ...
    - F ... ellei ehjää kopiota löydy
  - u NTFS Log FS pitää tiedostojärjestelmän muuten eheänä
- n **Yleinen lääke: Cluster Remapping**
  - u sektori otetaan pois käytöstä ja kyseinen looginen sektori mapataan muualle levyille
    - F voi hidastaa peräkkäiskäyttöä jatkossa

## Käyttöjärjestelmät II

### Linux ext3fs

## Linux ext3fs [Tweedie talk, 20.7.2000] [click](#)

(<http://olstrans.sourceforge.net/release/OLS2000-ext3/OLS2000-ext3.html>)

- n **ext3fs = ext2fs with journaling**
  - u journal in special file, or in special device
- n **Problem: time spend in recovering file system after a crash**
  - u fsck (e2fsck) takes too long – hours for big systems
  - u must have better availability
- n **Complete compatibility with ext2fs (ext2fs ↔ ext3fs)**
  - u clean, unmounted ext3fs has no journal, can mount as ext2fs

TKTL: ext3fs kaikissa tiedostopalvelimissa

## Linux ext3fs

- n **Extra layer on top of ext2fs: JFS (journaling FS)**
  - u independent on actual file system (ext2fs)
    - F ext2fs does not know about journaling!
  - u arbitrary modifications in buffer cache
  - u transactional semantics
    - F "do all these 5 updates, or none of them"
  - u API to add transactions onto a block device
- n **File update**
  - u no data to disk until transaction commit
    - F no quarantine when it will be written (**write behind**)
    - F written to disk from JFS cache of updates
  - u **write to log first, commit log, then do file update**
    - F disks can guarantee one sector write with power failure
    - F use special sector updates as commit blocks for log

## Kertauskysymyksiä

- n Mihin tarvitaan VFS:ää?
- n Kuinka ext2 poikkeaa "iciwanhasta" UNIXin tdstojärjestelmästä?
- n Mistä ext2:n tehokkuus / luotettavuus?
- n Mitä tietoja superlohkossa/indeksisolmussa?
- n Miksi i-solmuja? Miksi ei attribuutit ja nimi samassa paikassa?
- n Mihin tarvitaan NFS-protokollaa?
- n Mitä hyötyä on NFS:n tilattomuudesta?
- n Mitä on "journaling" ja "logging"? Windows vs. Linux?