

# File System (II)

# Roadmap for Lecture

- NTFS Recovery Support
  - Log File Service Operation
  - NTFS Recovery Procedures
  - Fault-Tolerance Support
  - Volume Management -  
Striped and Spanned Volumes
- Encryption File System Security
  - Encrypting File System (EFS) Terminology
  - EFS Operation
  - Data Encryption and Decryption
  - Windows EFS Architecture
  - Encryption Process Details

# NTFS Recovery Support

- Transaction-based logging scheme
- Fast, even for large disks
- Recovery is limited to file system data
  - Use transaction processing like SQL server for user data
  - Tradeoff: performance versus fully fault-tolerant file system
- Design options for file I/O & caching:
  - **Careful write:** VAX/VMS fs, other proprietary OS fs
  - **Lazy write:** most UNIX fs, OS/2 HPFS

# Careful Write File Systems

- OS crash/power loss may corrupt file system
- Careful write file system orders write operations:
  - System crash will produce predictable, non-critical inconsistencies
- Update to disk is broken in sub operations:
  - Sub operations are written serially
  - Allocating disk space: first write bits in bitmap indicating usage; then allocate space on disk
- I/O requests are serialized:
  - Allocation of disk space by one process has to be completed before another process may create a file
  - No interleaving sub operations of the two I/O requests
- Crash: volume stays usable; no need to run repair utility

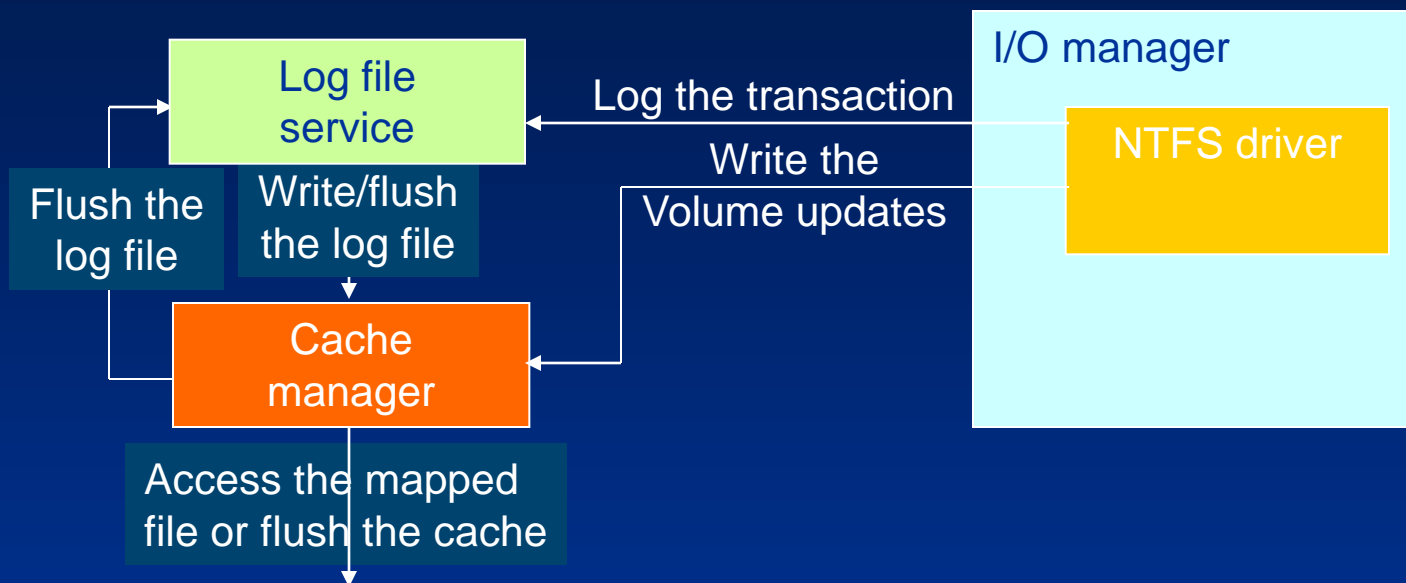
# Lazy Write File Systems

- Careful file system write sacrifices speed for safety
- Lazy write improves performance by write back caching
  - Modifications are written to the cache;
  - Cache flush is an optimized background activity
- Less disk writes; buffer can be modified multiple times before being written to disk
- File system can return to caller before op. is completed
- Inconsistent intermediate states on volume are ignored
- Greater risk / user inconvenience if system fails

# Recoverable File System (Journaling File System)

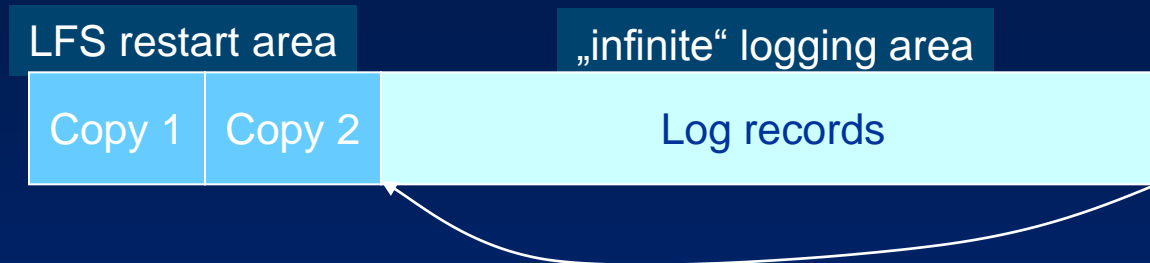
- Safety of careful write fs / performance of lazy write fs
- Log file + fast recovery procedure
  - Log file imposes some overhead
  - Optimization over lazy write: distance between cache flushes increased
- NTFS supports *cache write-through* and *cache flushing* triggered by applications
  - No extra disk I/O to update fs data structures necessary: all changes to fs structure are recorded in log file which can be written in a single operation
  - In the future, NTFS may support logging for user files (hooks in place)

# Log File Service (LFS)



- LFS is designed to provide logging to multiple kernel components (clients)
- Currently used only by NTFS

# Log File Regions



- NTFS calls LFS to read/write restart area
  - Context info: location of logging area to be used for recovery
  - LFS maintains 2nd copy of restart area
  - Logging area: circularly reused
  - LFS uses logical sequence numbers (LSNs) to identify log records
- NTFS never reads/writes transactions to log file directly
- During recovery:
  - NTFS calls LFS to read forward; recorded transactions are redone
  - NTFS calls LFS to read backward; undo all incompletely logged transactions



# Operation of the LFS/NTFS

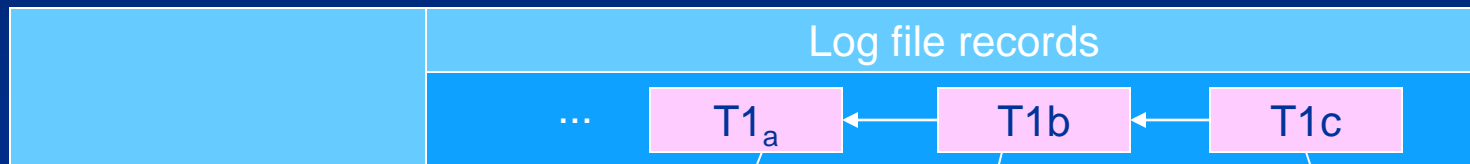
1. NTFS calls LFS to record in (cached) log file any transactions that will modify volume structure
2. NTFS modifies the volume (also in the cache)
3. Cache manager calls LFS to flush log file to disk (LFS implements flushing by calling cache manager back, telling which page to flush)
4. After cache manager flushes log file, it flushes volume changes

Transactions of unsuccessful modifications can be retrieved from log file and un-/redone

Recovery begins automatically the first time a volume is used after system is rebooted.

# Log Record Types

- Update records (series of ...)
  - Most common; each record contains:
  - **Redo information:** how to reapply on subop. of a committed trans.
  - **Undo information:** how to reverse a partially logged sub operation
- Last record commits the transaction (not shown here)



**Redo:** Allocate/initialize an MFT file record  
**Undo:** Deallocate the file record

**Redo:** Set bits 3-9 in the bitmap  
**Undo:** Clear bits 3-9 in the bitmap

**Redo:** Add the filename to the index  
**Undo:** Remove the filename from the index

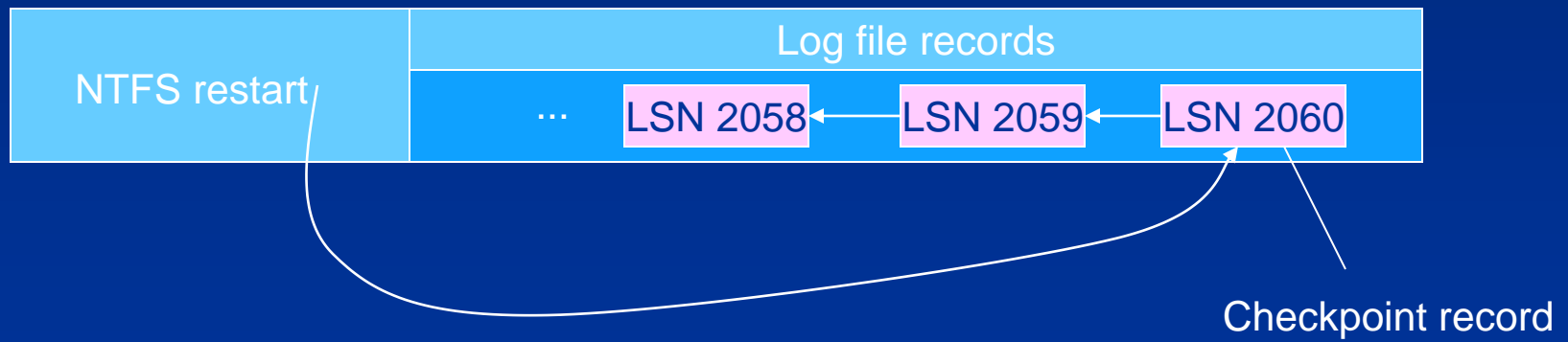
Recovery: redo committed/undo incompletely logged transact.

# Log Records (contd.)

- Physical vs. logical description of redo/undo actions:
  - Delete byte range on disk vs. Delete file “a.dat”
  - NTFS writes update records with physical descriptions
- NTFS writes update records (usually several) for:
  - Creating a file
  - Deleting a file
  - Extending a file
  - Truncating a file
  - Setting file information
  - Renaming a file
  - Changing security applied to a file
- Redo/undo ops. must be idempotent (can be applied multiple times)

# Checkpoint Records

- NTFS periodically writes a checkpoint record
  - Describes, what processing would be necessary to recover a volume if a crash would occur immediately
  - How far back in the log file must NTFS go to begin recovery
  - LSN of checkpoint record is stored in restart area



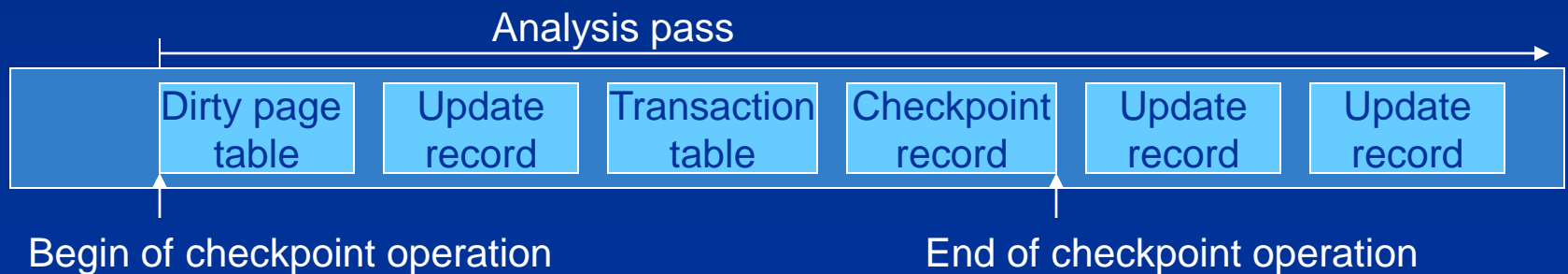
# Log File Full

LFS presents log file to NTFS as if it were infinitely large

- Writing checkpoint records usually frees up space
- LFS tracks several numbers:
  - Available log space
  - Amount of space needed to write an incoming log record and to undo the write
  - Amount of space needed to roll back all active (no committed) transactions, should that be necessary
- Insufficient space: “Log file full” error & NTFS exception
  - NTFS prevents further transactions on files (block creation/deletion)
  - Active transactions are completed or receive “log file full” exception
  - NTFS calls cache manager to flush unwritten data
  - If data is written, NTFS marks log file “empty”; resets beginning of log file
- No effect on executing programs (except short I/O pause)

# Recovery - Principles

- NTFS performs automatic recovery
- Recovery depends on two NTFS in-memory tables:
  - *Transaction table*: keeps track of active transactions (not completed)  
(sub operations of these transactions must be removed from disk)
  - *Dirty page table*: records which pages in cache contain modifications to file system structure that have not yet been written to disk
- NTFS writes checkpoint every 5 sec.
  - Includes copy of transaction table and dirty page table
  - Checkpoint includes LSNs of the log records containing the tables



# Recovery - Passes

## 1. Analysis pass

- NTFS scans forward in log file from beginning of last checkpoint
- Updates transaction/dirty page tables it copied in memory
- NTFS scans tables for oldest update record of a non-committed trans.

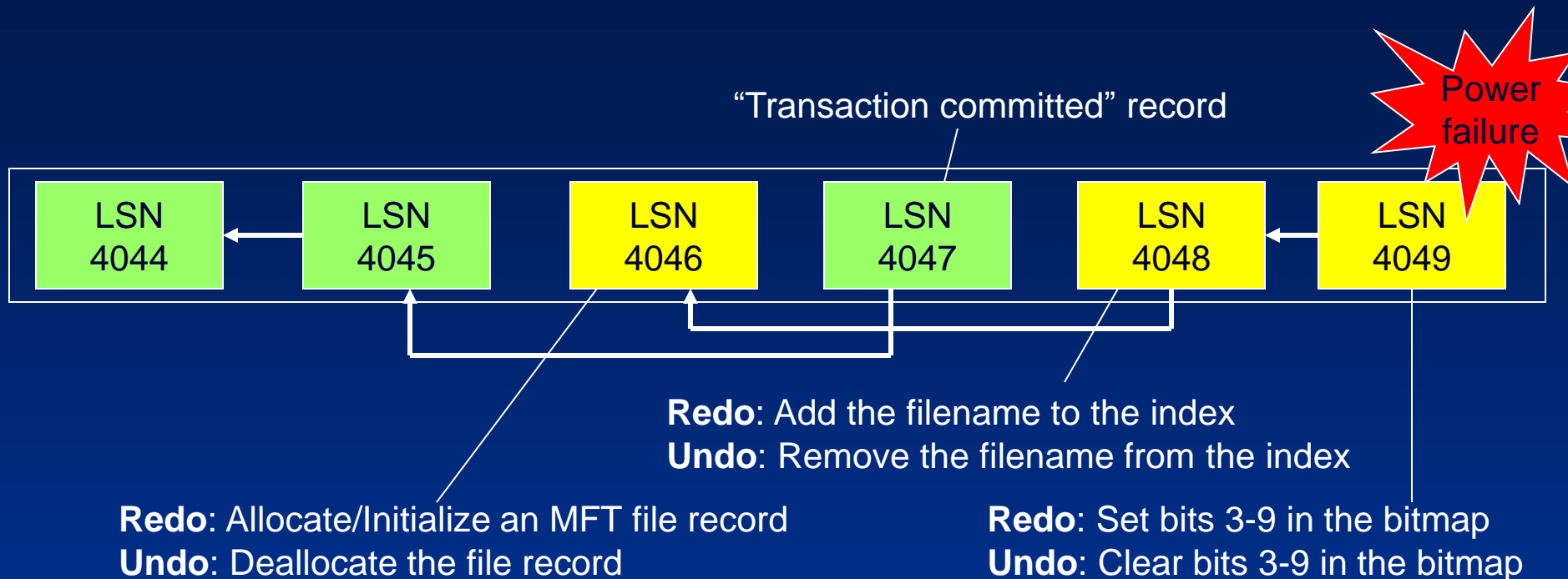
## 2. Redo pass

- NTFS looks for “page update” records which contain volume modification that might not have been flushed to disk
- NTFS redoes these updates in the cache until it reaches end of log file
- Cache manager “lazy writer thread” begins to flush cache to disk

## 3. Undo pass

- Roll back any transactions that weren't committed when system failed
- After undo pass – volume is at consistent state
- Write empty LFS restart area; no recovery is needed if system fails now

# Undo Pass - Example



- Transaction 1 was committed before power failure
- Transaction 2 was still active
- NTFS must log undo operations in log file!
  - Power might fail again during recovery;
  - NTFS would have to redo its undo operations



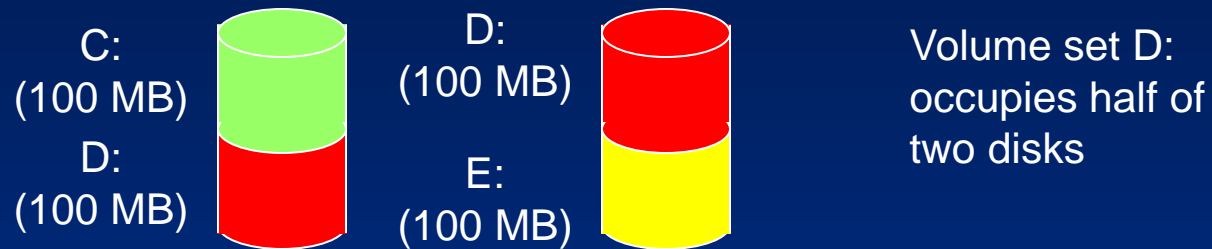
# NTFS Recovery - Conclusions

- Recovery will return volume to some pre-existing consistent state (not necessarily state before crash)
- Lazy commit algorithm: log file is not immediately flushed when a “transaction committed” record is written
  - LFS batches records;
  - Flush when cache manager calls or check pointing record is written (once every 5 sec)
  - Several parallel transactions might have been active before crash
- NTFS uses log file mechanisms for error handling
- Most I/O errors are not file system errors
  - NTFS might create MFT record and detect that disk is full when allocating space for a file in the bitmap
  - NTFS uses log info to undo changes and returns “disk full” error to caller

# Fault Tolerance Support

- NTFS' capabilities are enhanced by the fault-tolerant volume managers FtDisk/DMIO
  - Lies above hard disk drivers in the I/O system's layered driver scheme
  - FtDisk – for basic disks
  - DMIO – for dynamic disks
- Volume management capabilities:
  - Redundant data storage
  - Dynamic data recovery from bad sectors on SCSI disks
- NTFS itself implements bad-sector recovery for non-SCSI disks

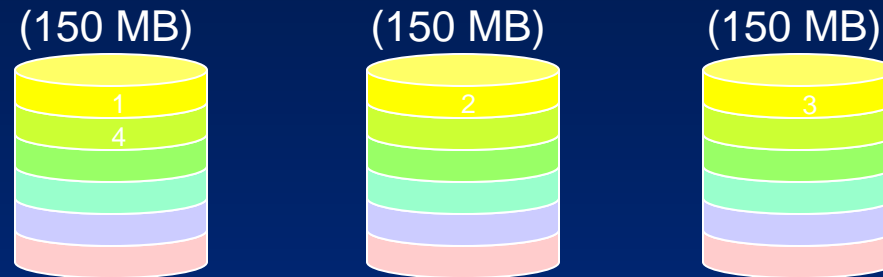
# Volume Management Features – Spanned Volumes



## Spanned Volumes:

- single logical volume composed of a maximum of 32 areas of free space on one or more disks
- NTFS volume sets can be dynamically increased in size (only bitmap file which stores allocation status needs to be extended)
- FtDisk/DMIO hide physical configuration of disks from file system
- Tool: Windows Disk Management MMC snap-in
- Spanned volumes were called volume sets in Windows NT 4.0

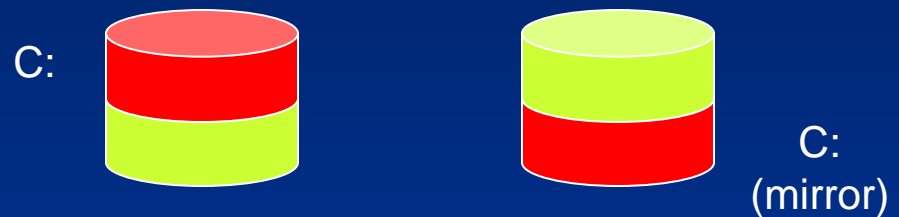
# Striped Volumes



- Series of (up to 32) partitions, one partition per disk (of same size)
- Combined into a single logical volume
- FtDisk/DMIO optimize data storage and retrieval times
  - Stripes are relatively narrow: 64KB
  - Data tends to be distributed evenly among disks
  - Multiple pending read/write ops. will operate on different disks
  - Latency for disk I/O is often reduced (parallel seek operations)

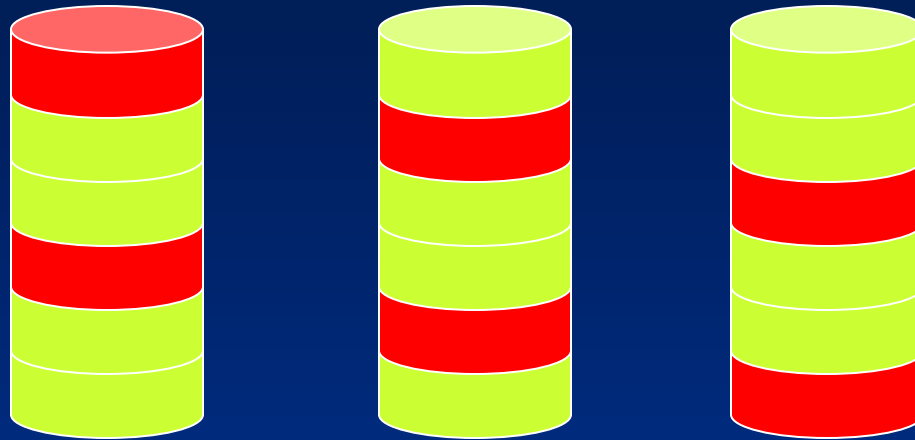
# Fault Tolerant Volumes

- FtDisk/DMIO implement redundant storage schemes
  - Mirror sets
  - Stripe sets with parity
  - Sector sparing
- Tools: Windows Disk Management MMC snap-in



- **Mirrored Volumes (a.k.a. RAID-1):**
  - Contents of a partition on one disk are duplicated on another disk
  - FtDisk/DMIO write same data to both locations
  - Read operations are done simultaneously on both disks (load balancing)

# RAID-5 Volumes

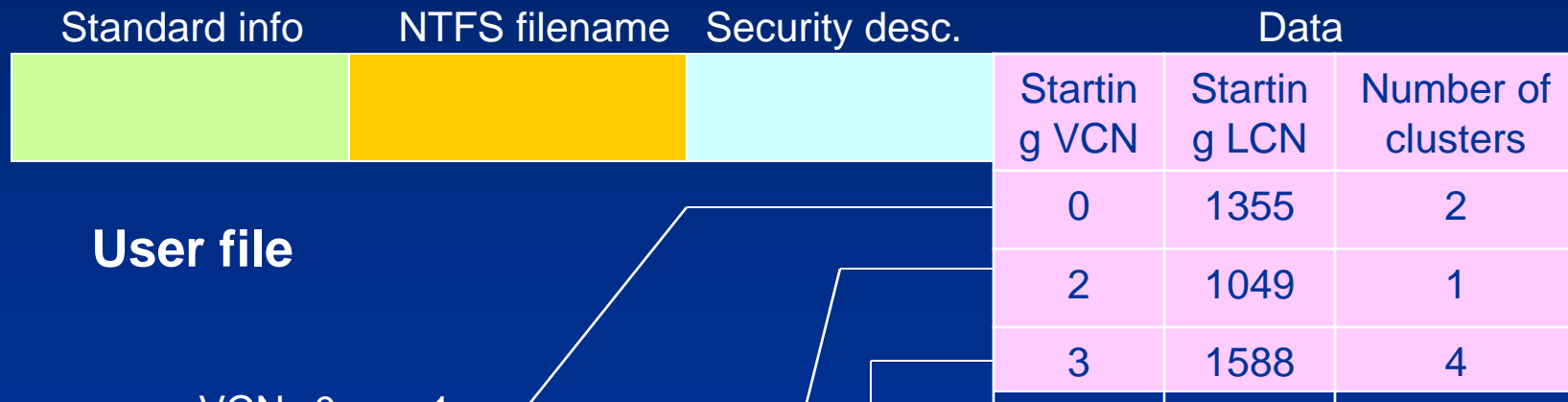
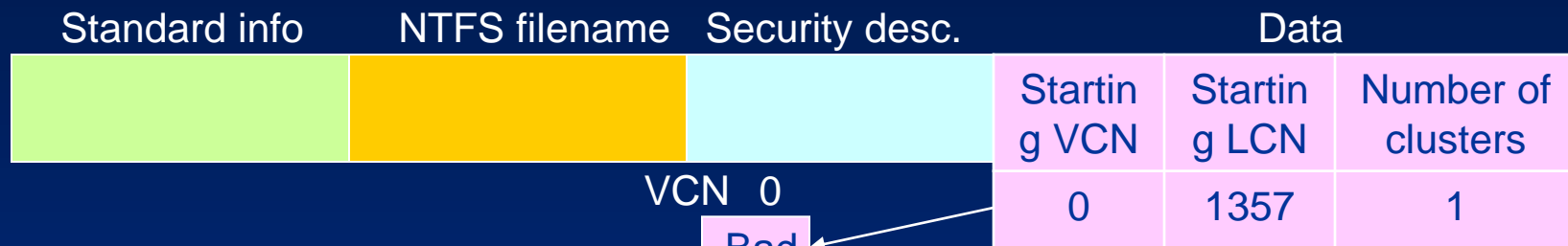


- Fault tolerant version of a regular stripe set
- Parity: logical sum (XOR)
- Parity info is distributed evenly over available disks
- FtDisk/DMIO reconstruct missing data by using XOR op.
- However I/O performance degraded until the failed disk is replaced

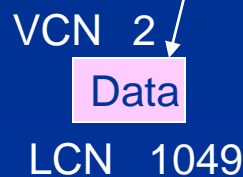
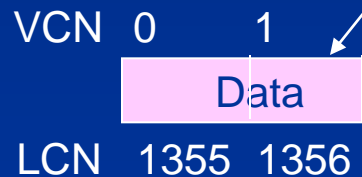
# Bad Cluster Recovery

- Sector sparing is supported by FtDisk/DMIO
  - Dynamic copying of recovered data to spare sectors
  - Without intervention from file system / user
  - Works for certain SCSI disks
  - FtDisk/DMIO return bad sector warning to NTFS
- Sector re-mapping is supported by NTFS
  - NTFS will not reuse bad clusters
  - NTFS copies data recovered by FtDisk/DMIO into a new cluster
- NTFS cannot recover data from bad sector without help from FtDisk/DMIO
  - NTFS will never write to bad sector (re-map before write)

# Bad-cluster re-mapping



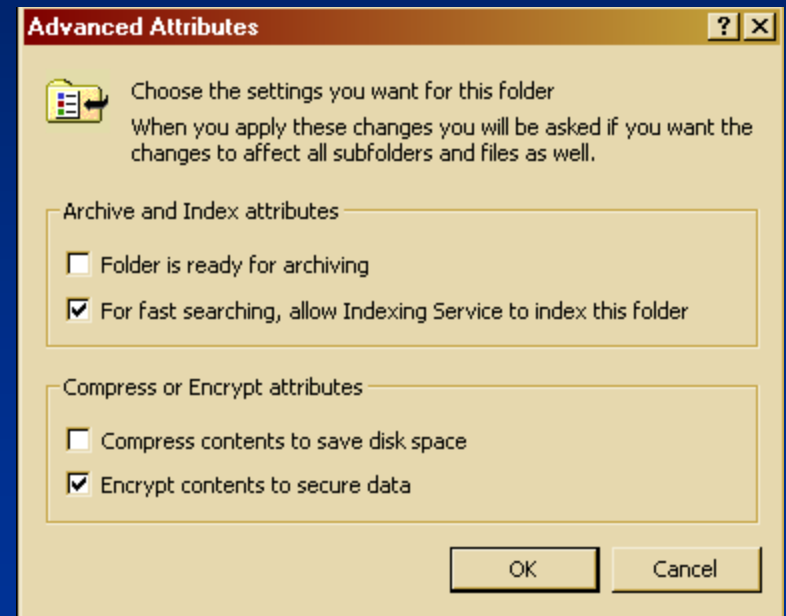
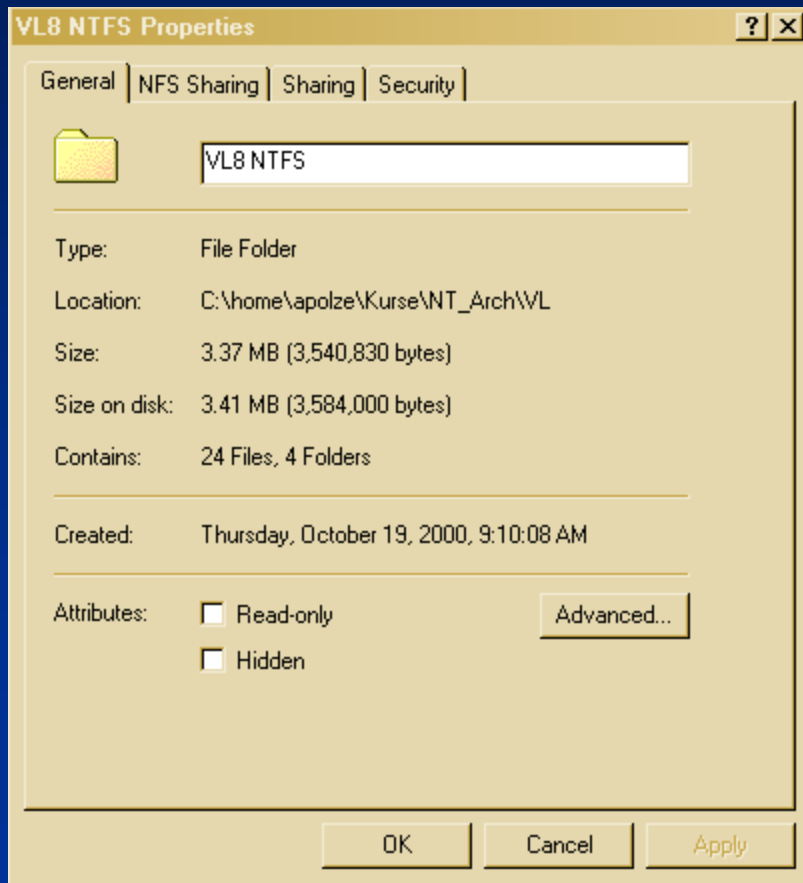
## User file





# Encrypting File System Security

- EFS relies on Windows cryptography support
  - Transparent encryption through Windows Explorer or cipher-utility



# EFS operation

- When a file is encrypted...
  - EFS generates random File Encryption Key (FEK) to encrypt file content
  - Stronger variant of Data Encryption Standard (3DES or AES) to encrypt file content (fast, shared secret)
  - File's FEK is stored with file and encrypted using the file creator's RSA public key (slow)
- File can be decrypted...
  - only with the user's private RSA key
  - What about lost keys?
- FEK can be stored in multiple encryptions...
  - Users can share an encrypted file
  - Can store a recovery key to allow recovery agents access to files
- Secure public/private key pairs are essential
  - Stored on \Users\xxx\AppData\Roaming\Microsoft\Crypto\RSA

# Basic Terminology

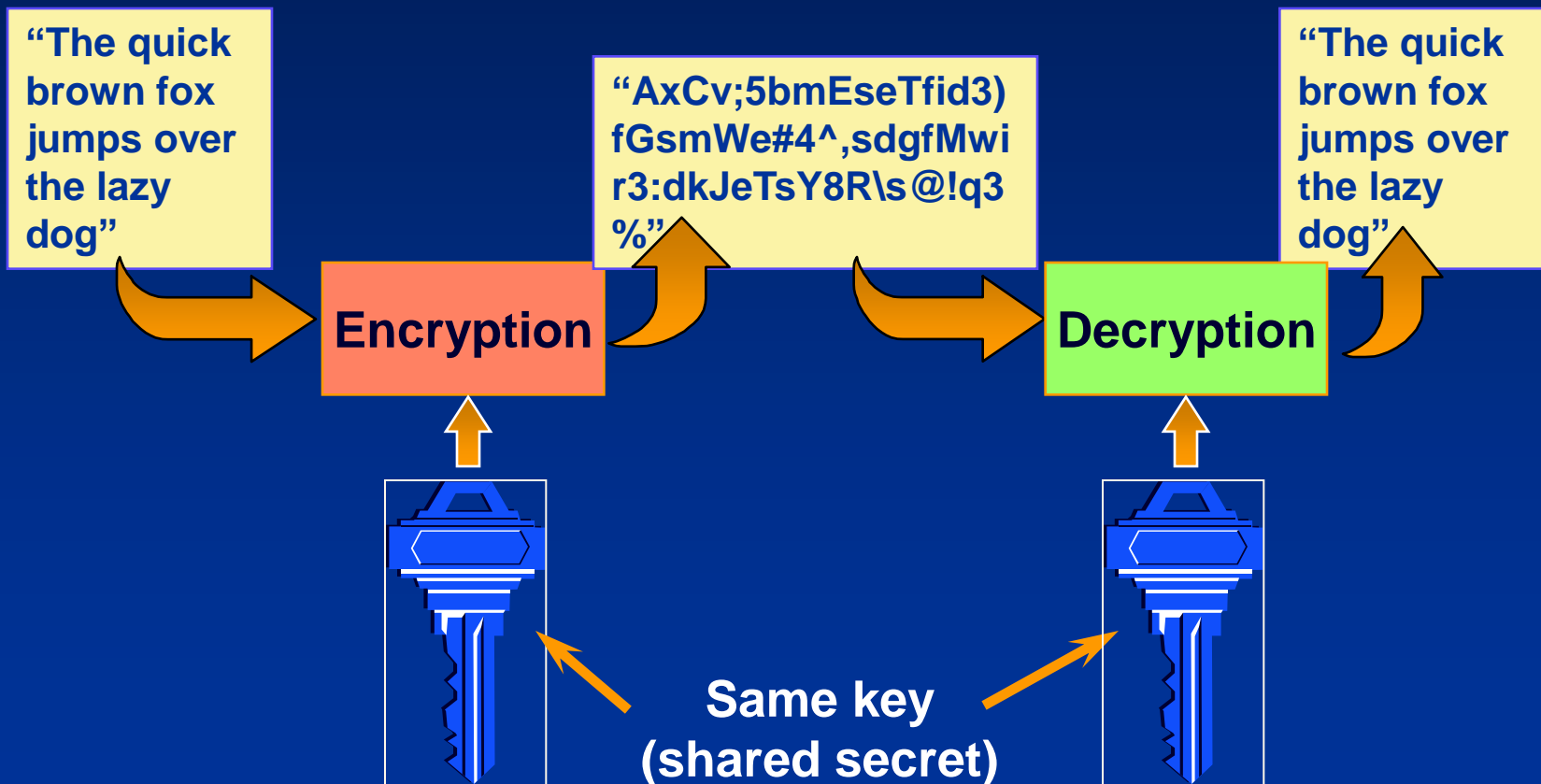
- Plaintext
  - The stuff you want to secure, typically readable by humans (email) or computers (software, order)
- Ciphertext
  - Unreadable, secure data that must be decrypted before it can be used
- Key
  - You must have it to encrypt or decrypt (or do both)
- Cryptoanalysis
  - Hacking it by using science
- Complexity Theory
  - How hard is it and how long will it take to run a program

# Symmetric Key Cryptography

Plain-text input

Cipher-text

Plain-text output



# Symmetric Pros and Cons

- Weakness:

- Agree the key beforehand
- Securely pass the key to the other party

- Strength:

- Simple and really very fast (order of 1000 to 10000 faster than asymmetric mechanisms)
  - Super-fast if done in hardware (DES)
  - Hardware is more secure than software, so DES makes it really hard to be done in software, as a prevention

# Public Key Cryptography

- Knowledge of the encryption key doesn't give you knowledge of the decryption key
- Receiver of information generates a pair of keys
  - Publish the public key in directory
- Then anyone can send him messages that only she can read

# Public Key Encryption

Clear-text Input

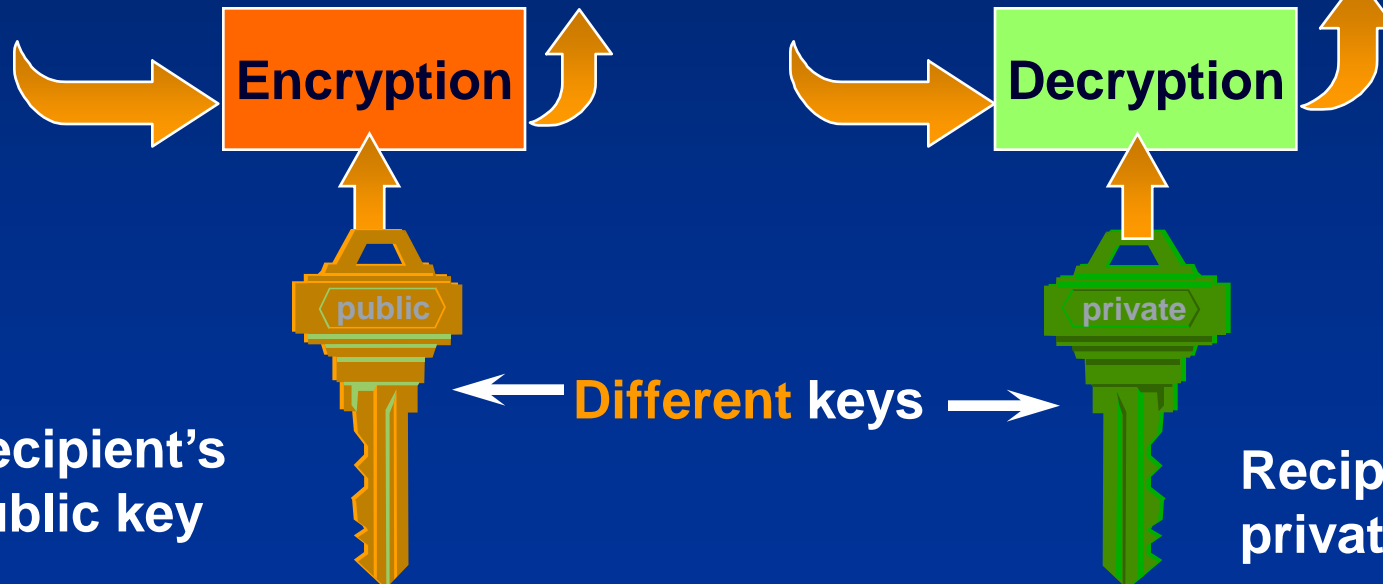
“The quick  
brown fox  
jumps over  
the lazy  
dog”

Cipher-text

“Py75c%bn&\*)9|fDe^  
bDFaq#xzfFr@g5=&n  
mdFg\$5knvMd’rkveg  
Ms”

Clear-text Output

“The quick  
brown fox  
jumps over  
the lazy  
dog”



Recipient's  
public key

Recipient's  
private key

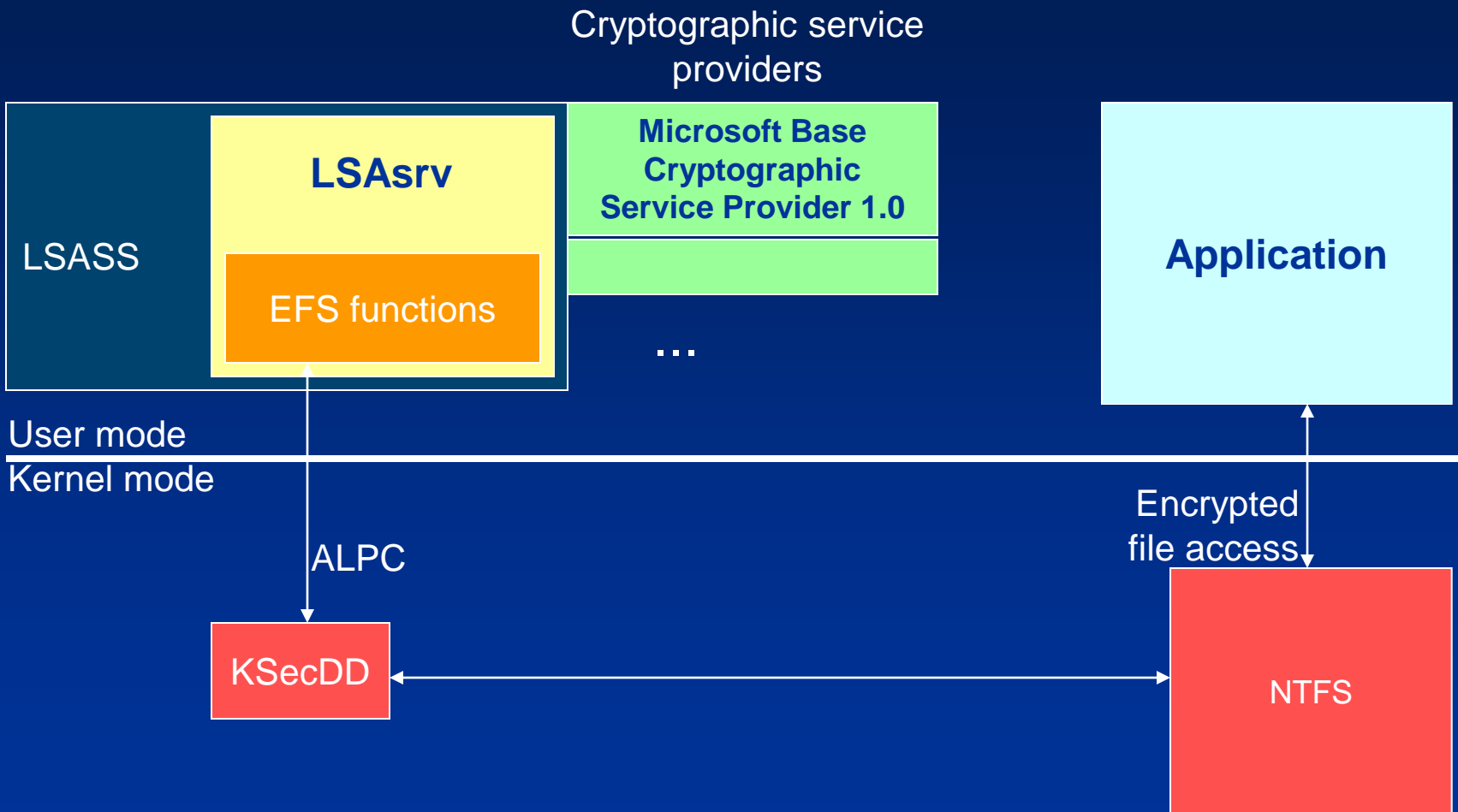
Different keys

# Problem of Key Recovery

- What if you lose the private key? 😊
- Data recovery by authorized agents
  - Integrated key management
- Windows:
  - Flexible recovery policy
    - Enterprise, domain, or per machine
  - Encrypted backup and restore
    - Integrated with Windows backup
- Potential weakness but you can opt not to use it!



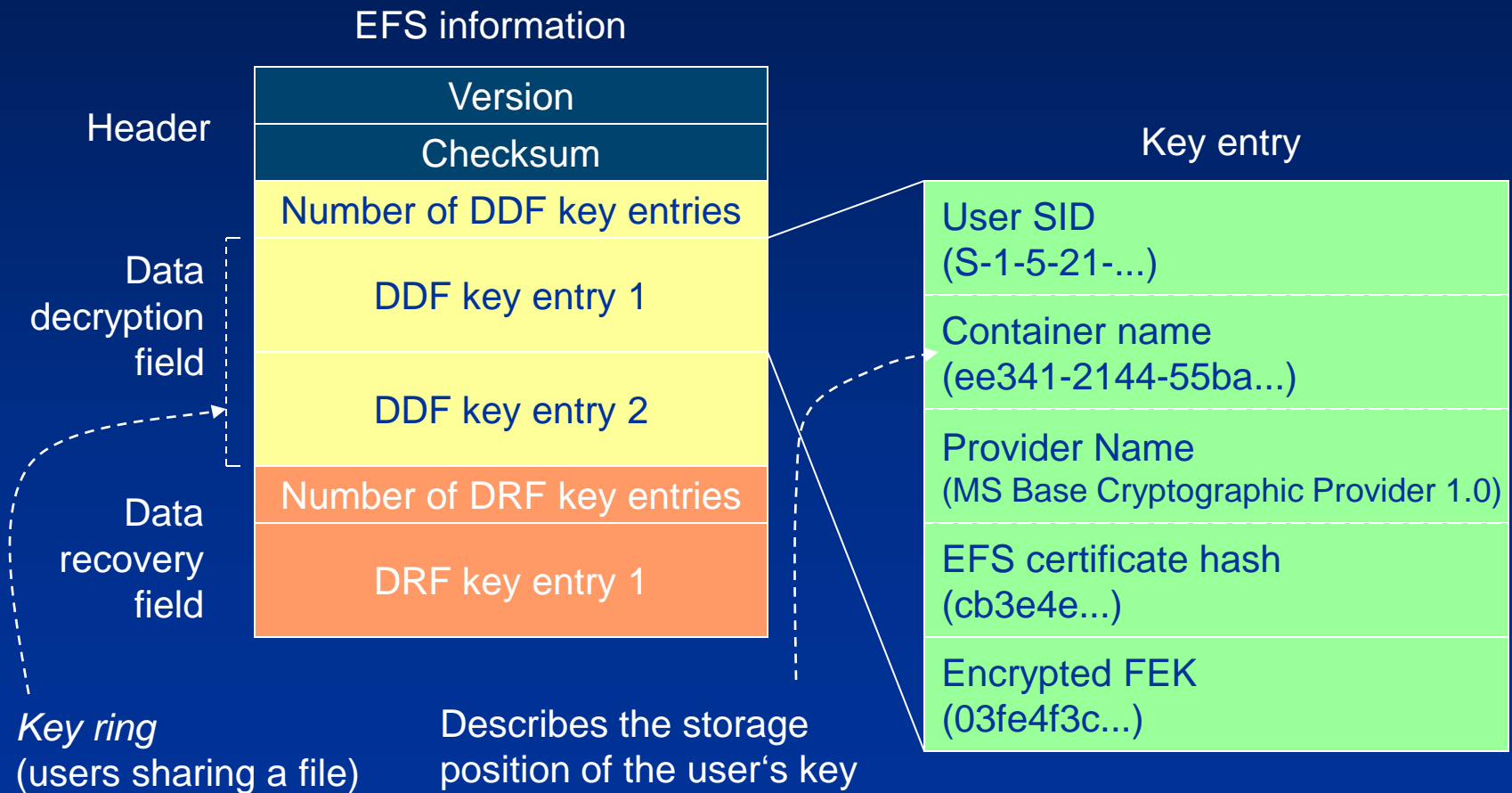
# Windows EFS Architecture



# EFS Components

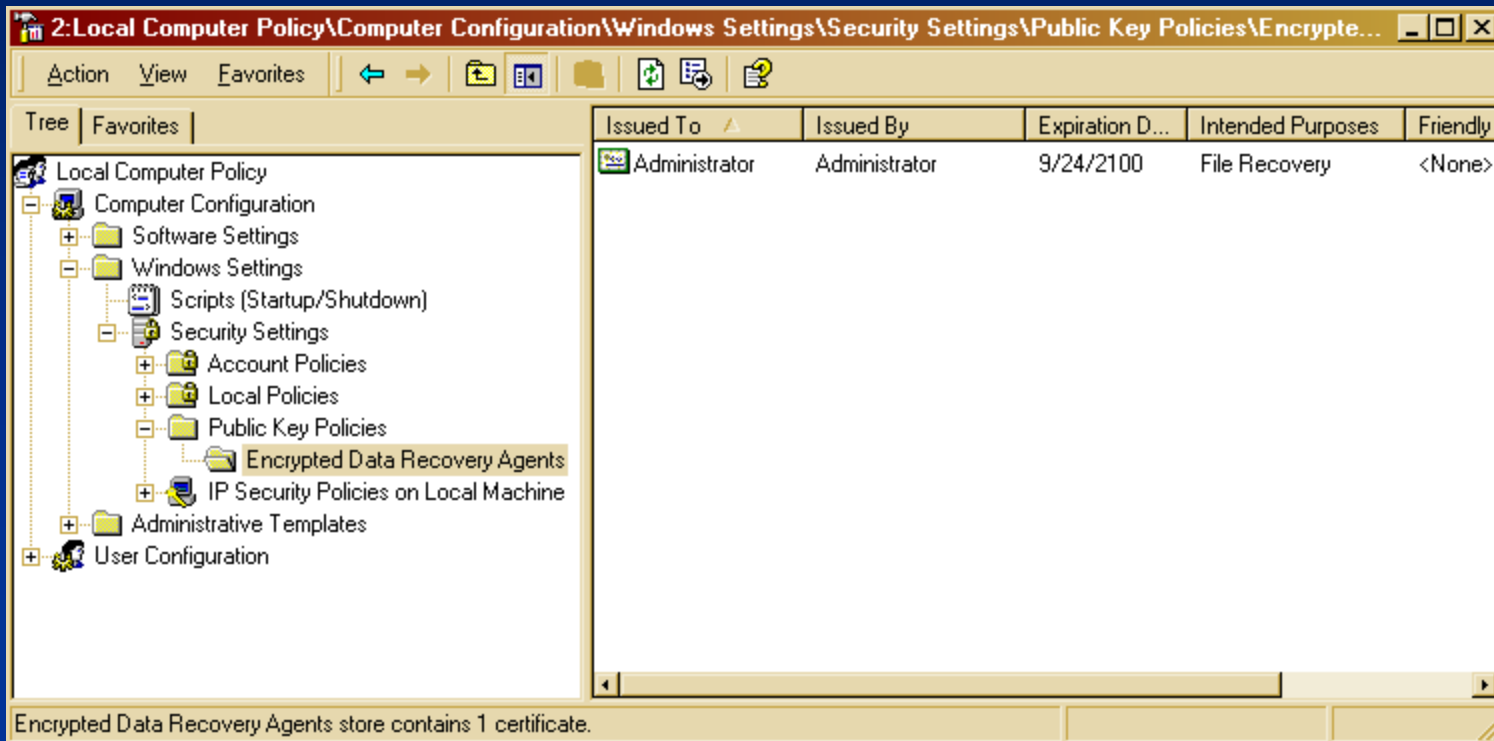
- Local Security Authority Subsystem
  - LSASS (\Winnt\System32\lsass.exe) manages logon sessions
  - EFS obtains FEKs from LSASS
- KSecDD device driver implements comm. with LSASS
- LSAsrv listens for ALPC comm.
  - Passes requests to EFS functions
  - Uses functions in MS CryptoAPI (CAPI) to decrypt FEK for EFS
- Crypto API ...
  - is implemented by Cryptographic Service Provider (CSP) DLLs
  - Details of encryption/key protection are abstracted away

# Format of EFS information and key entries for a file

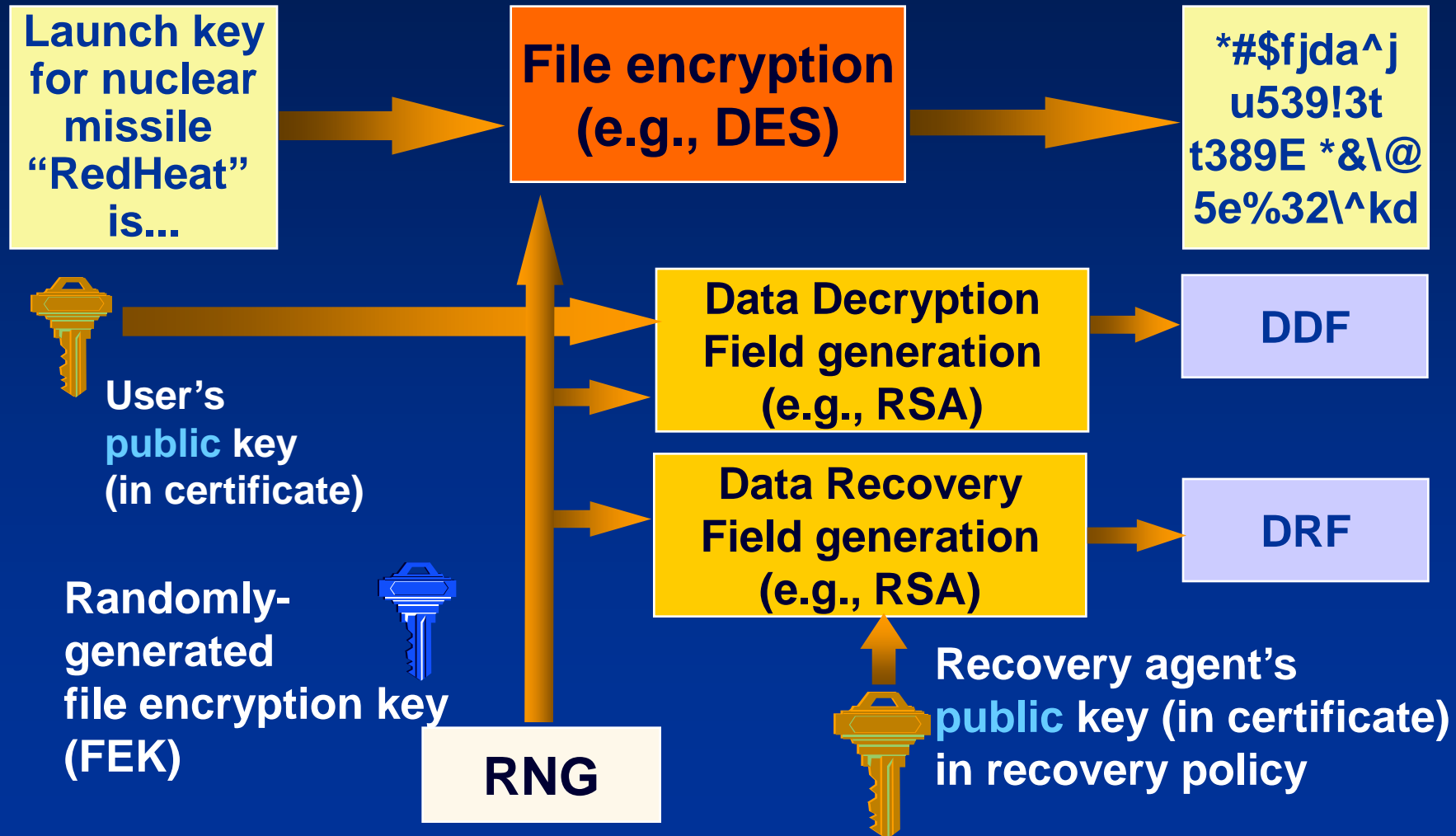


# Encrypted Data Recovery Agents group policy

- Use Group Policy MMC snap-in to configure recovery agents (...list may be empty)



# Data Encryption Process



# Encryption Process Details

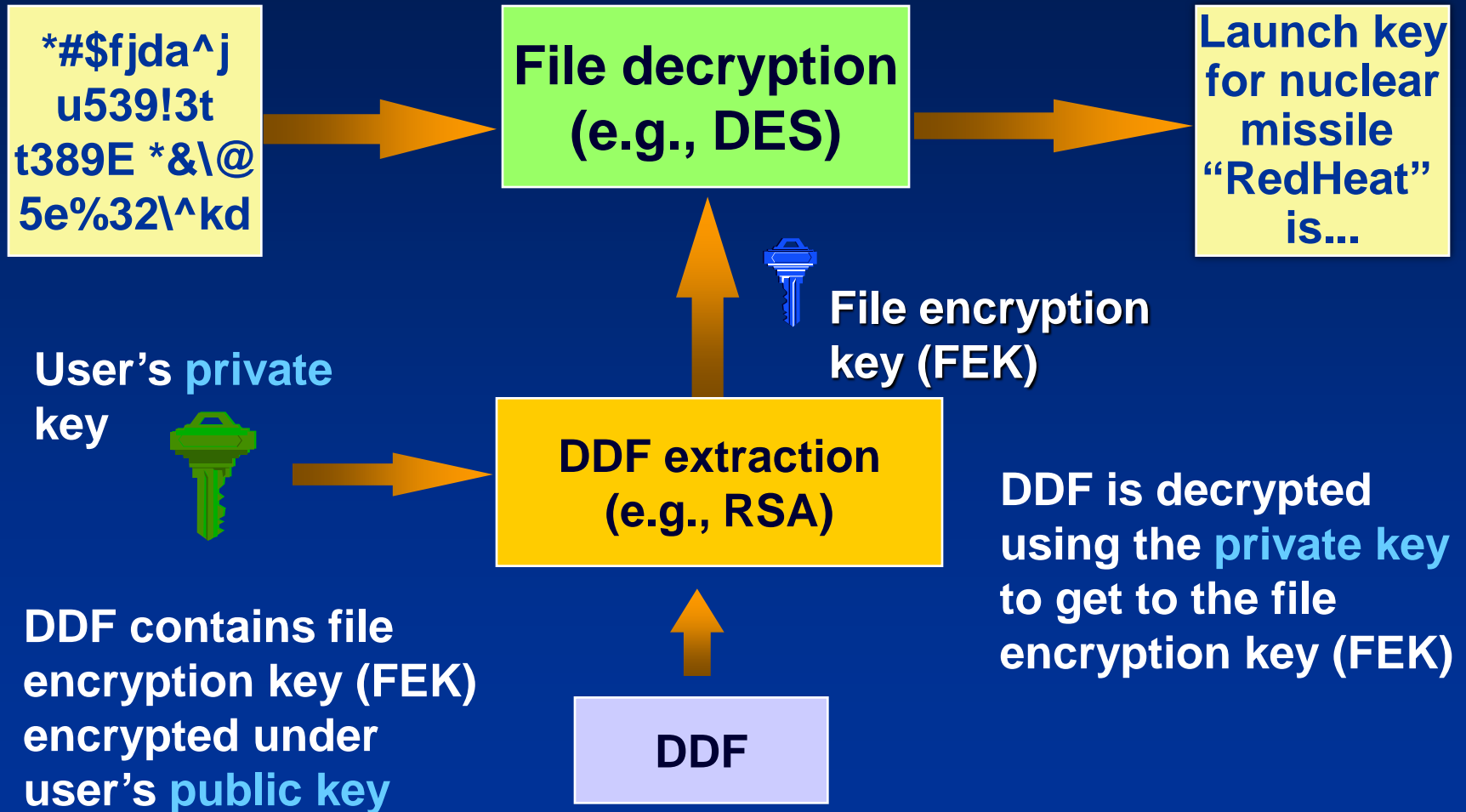
1. User profile is loaded if necessary
2. A log file Efsx.log is created
  - In system volume info dir; x is unique number
3. Base Cryptographic Provider 1.0 generates random 128-bit FEK
4. User EFS private/public key pair is generated or obtained
  - HKEY\_CURRENT\_USER\Software\Microsoft\Windows NT\CurrentVersion\EFS\CurrentKeys\CertificateHash identifies the user's key pairs
5. A DDF key ring is created for the file with an entry for the user
  - Entry contains copy of FEK encrypted with user's public key
6. A DRF key ring is created for the file
  - Has an entry for each recovery agent on the system
  - Entries contain copies of FEK encrypted with agents' public keys

# Encryption Process Details (contd.)

7. A backup file is created (Efs0.tmp)
  - Same directory as original file
8. DDF and DRF rings are added to a header
  - EFS attributes - \$LOGGED\_UTILITY\_STREAM
9. Backup file is marked encrypted, original file is copied to backup
10. Original file's contents are destroyed
  - Backup is copied to original
  - This results in encrypting the file contents
11. The backup file is deleted
12. The log file is deleted
13. The user profile is unloaded (if it was loaded in step 1)

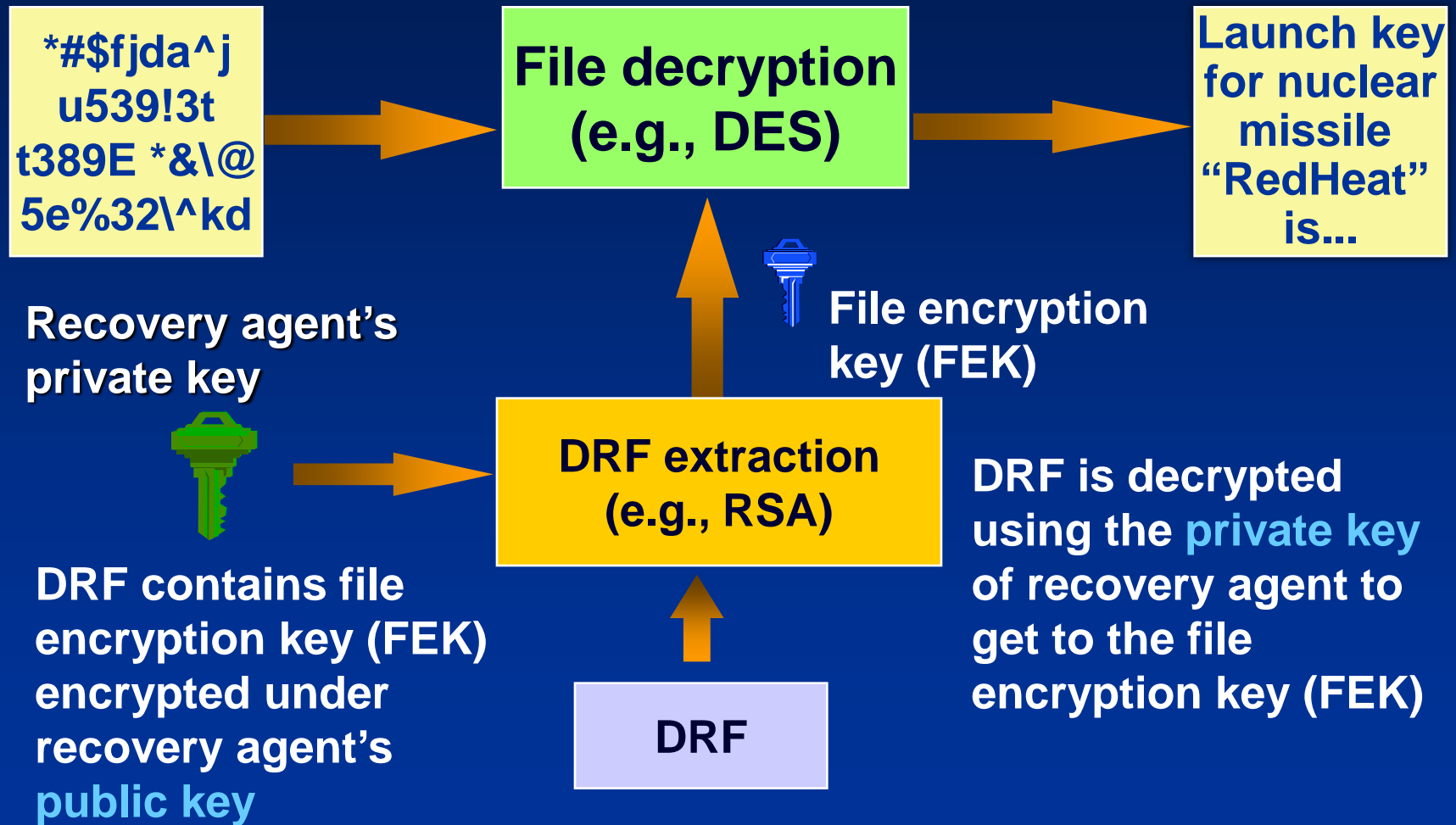
In case of system crash, either original file or backup contain valid copy of the file content.

# Data Decryption Process

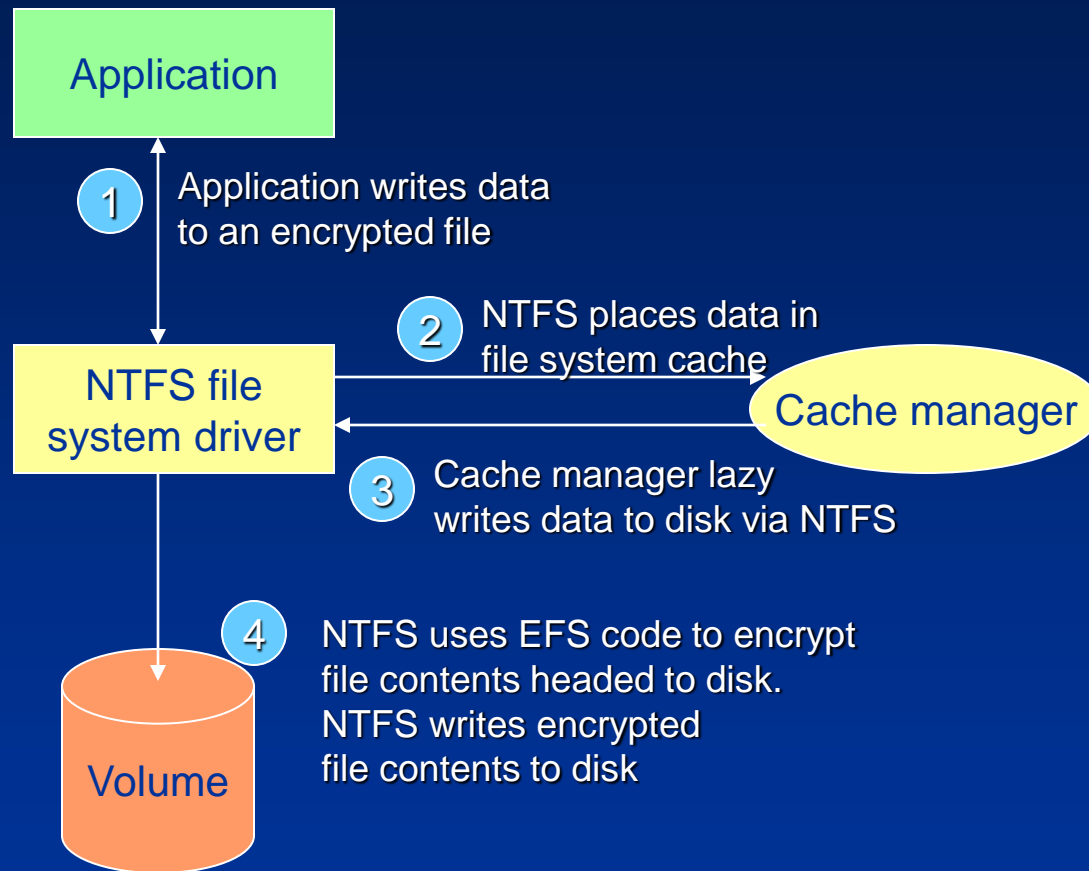




# Data Recovery Process



# Flow of EFS



# Backing Up Encrypted Files

- Data is never available in unencrypted form
  - Except to applications that access file via encryption facility
- EFS provides a facility for backup programs:
  - No need to encrypt or decrypt files when backup
  - EFS API: *OpenEncryptedFileRaw()*, *ReadEncryptedFileRaw()*, *WriteEncryptedFileRaw()*, *CloseEncryptedFileRaw()*
  - Implemented in Advapi32.dll, use ALPC to invoke function in LSAsrv
  - LSAsrv calls *EfsReadFileRaw()* to obtain file's EFS attribute and the encrypted contents from NTFS driver
  - Similarly, *EfsWriteFileRaw()* is invoked to restore file's contents

# Further Reading

- Mark E. Russinovich *et al.* Windows Internals, 5th Edition, Microsoft Press, 2009.
  - Chapter 11 - File Systems  
NTFS Recovery Support (from pp. 974)
  - Chapter 8 - Storage Management  
Multipartition Volume Management (from pp. 661)
  - Encrypting File System Security (from pp. 990)
- *Applied Cryptography*, B. Schneier, John Wiley & Sons, ISBN 0-471-12845-7
- *Handbook of Applied Cryptography*, A.J. Menezes, CRC Press, ISBN 0-8493-8523-7

# Source Code References

- Windows Research Kernel sources do not include NTFS
- A raw file system driver is included in `\base\ntos\raw`
- Also see `\base\ntos\fstrl` (File System Run-Time Library)