

# Storage Systems

## NPTEL Course

### Jan 2012

(Lecture 14)

K. Gopinath

Indian Institute of Science

# Types of Storage

- “Mobile” Storage
  - Memory Stick, Camera, Smart Phone, Laptop
- Personal Storage
  - PC, “Home” RAID systems, “Home” NFS server
- Dept/Organizational Storage
  - NFS/CIFS server
- Cloud Storage
- Highly Available Storage
- Parallel Storage
- Web-scale storage
- Secure Storage
- Attribute-based Storage (“QoS”)
- Long term storage
  - DNA storage!?

# Parallel Storage

- Access to storage in parallel in HPC appls
- Multiple threads (cores) working on a large distr file across many nodes
- Scalability issues critical
  - Access/Updating of metadata
  - Bottlenecks in file name processing

# Web-scale storage

- Google
  - GFS/BigTable, gmail, ...
- Facebook
- Dropbox, ...
- Cloud storage
  - Amazon, Azure storage services

# Facebook Modular Storage

- Web/Chat
  - High CPU, low memory 16GB, low capacity 250GB disks
- Database
  - Med CPU, high memory 144GB, High IOPs: 3.2TB flash
- Hadoop
  - Med CPU, med memory 48GB, high capacity: 12x3TB SATA
- Haystack (photos, video)
  - Low CPU, low memory 18GB, high capacity: 12x3TB SATA
- Feed (ads, search, multifeed (to feed the “Wall”))
  - High CPU, high memory 144GB, medium 2TB SATA
- “Flash sled”: 500GB-8TB, 600kIOPS
- “Storage sled”: 15 drives, 3kIOPS

# Storage Security

- Appl
  - Can be selective wrt data but need to know what is sensitive
  - Each appl has its own mgmt of security
- FS/OS
  - Can be selective but costly in time in sw
  - Changes in FS/OS?
- Device driver (HBA), netw interface
  - Only for data in flight (SAN)
- Network
  - Heterogenous devices OK (disk/tape)
- Storage Controller
- Device

# Secure Storage: Disk

- Entire drive (incl MBR) encrypted
  - MBR+OS unmodified (transparency); MBR cannot be corrupted
  - Auth occurs before OS/any malicious software loaded
- BIOS attempts MBR read
  - drive redirects to pre-boot area
- Drive loads pre-boot OS
- User enters auth credentials for drive to verify
- If auth successful, drive loads original MBR
- Normal operation commences

# Attribute-based Storage

- Latency, BW, reliability guarantees
  - May be power, longevity, ... also!
- End2end guarantees typ
  - Difficult!
- In research stage still!



# Saving Current Documents for the Next Millennium

Across time (written now 2012 and read then 3012) & space (written here on Earth & read there on Mars!)

Say document written in Postscript/ Wordstar/ Word!

Stored on a SCSI 36GB 15K rpm disk drive on, say, a ext3 filesystem on Linux...

What is the **hard** technical problem?

Drives/device driver/filesystem/kernel/application may become obsolete

Along with document, the associated model of device/software may need to be saved (recursion problem!)

Any current technologies useful?

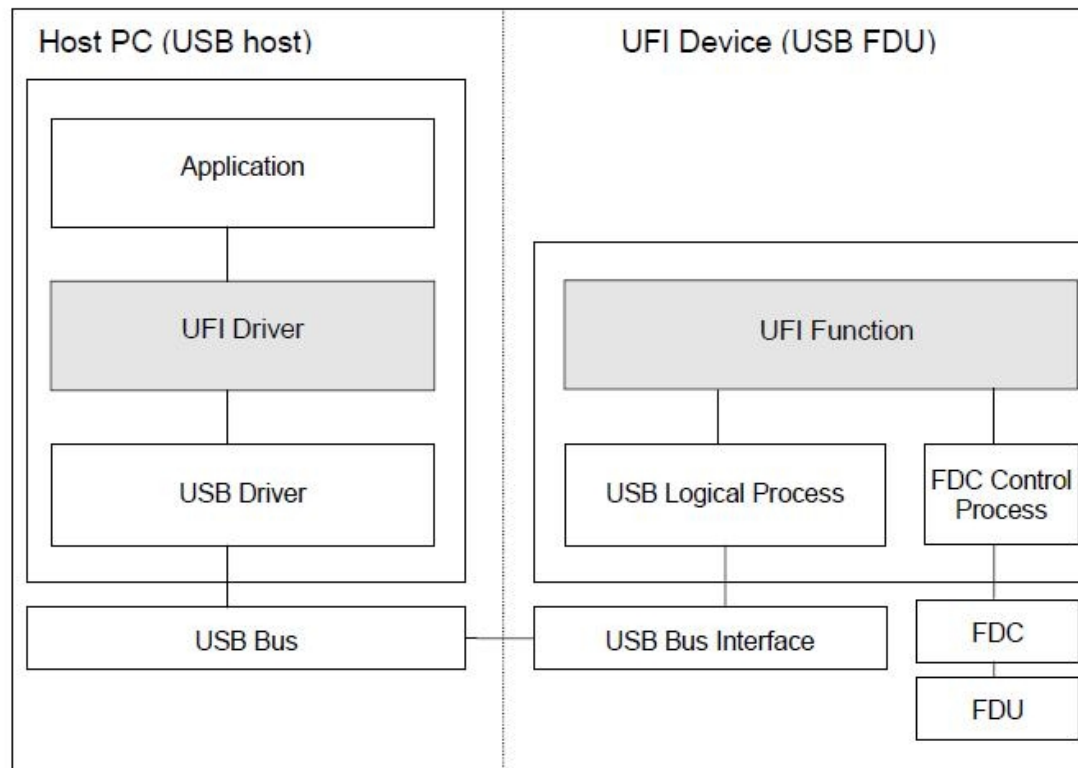
ASN.1? Virtual machines? Virtualization?!

*Till today no effective solution!*

# Evolution

- Design and architecture of subsystems (say, USB) may need to be modified over a period of time.
- Consider a (new) development where we have thin clients with a non-local server storing data and configuration.
- Users would still want the convenience of USB devices; if this is not possible, thin clients will not succeed.
- Though USB devices were expected to be electrically connected to the PC, it is possible to have remote USB functionality with appropriate redesign of the layers, using Internet as the connection.
- We need to “virtualise” some aspects of the previous legacy design.
  - also happened with other devices
    - terminals resulting in pseudo-terminals
    - SCSI resulting in iSCSI designs that allow remote access to disks through Internet.

# Emulation of Floppy Drive (from UFI std)



Note:  indicates software handling UFI commands

UFI: USB Floppy Interface

FDU: Floppy Disk Unit

FDC: Floppy Disk Controller

# Indus Script

- Yet un-deciphered: meaning across time not yet accomplished
- Compare with hieroglyphics (Egyptian Rosetta stone): three scripts side by side
- What is the problem? Not enough contextual info:
  - can see the script (human "readable")
  - no mapping between symbols and phonemes
  - need interpretation of sequences of symbols
  - a problem in archaeology, history, society,...

# Vedas

Transmitted across atleast 3500-5000 years without differing versions

Including exact pronunciation!

“UNESCO proclaimed the tradition of Vedic chant a Masterpiece of the Oral and Intangible Heritage of Humanity on November 7, 2003”

What "technology" used? Redundancy!

Various "pathas" of Samhita text: can recover from a corrupted text due to added redundancy: RAID-like! (Redundant Array of Indep Disks)

Pada-patha: each word in its separate form

Krama-patha: connects a word in pairs

ABCD becomes AB BC CD DE... (“2-mirroring”): 2 copies

Jata-patha: ABBAAB (“3-mirroring”): 3 copies of A, B, ...

Ghana-patha (ABBA ABCCBA ABC BCCB BCDDCB BCD...)(“10x”)

Metrical (similar to checksums!) & Musical

"Information dispersal"

Human Reproduction! (Oral transmission)

Use efficient “virtualizers”!

तम् । भा॒ग॒धे॒ये॒न । वि । मु॒ञ्च॒ति । प्र॒ति॒ष्ठि॒त्यै । यया॑ । रज्ज्वा॑ । उ॒त्त॒मां । गा॒म् । आ॒जे॒त् ।  
 । ता॒म् । भ्रातृ॑व्याय । प्र । हि॒णु॒या॒त् । नि॒र्ऋ॒ति॒म् । ए॒व । अ॒स्मै । प्र । हि॒णो॒ति॒  
 ॥ तै सं २-२-६-५ ॥

तं भा॒ग॒धे॒ये॒न भा॒ग॒धे॒ये॒न तं तं भा॒ग॒धे॒ये॒न वि वि भा॒ग॒धे॒ये॒न तं तं भा॒ग॒धे॒ये॒न वि ॥  
 भा॒ग॒धे॒ये॒न वि वि भा॒ग॒धे॒ये॒न भा॒ग॒धे॒ये॒न वि मु॒ञ्च॒ति मु॒ञ्च॒ति वि भा॒ग॒धे॒ये॒न भा॒ग॒धे॒ये॒न वि  
 मु॒ञ्च॒ति । भा॒ग॒धे॒ये॒ने॒ति॒ भा॒ग॒धे॒ये॒न ॥  
 वि मु॒ञ्च॒ति मु॒ञ्च॒ति वि वि मु॒ञ्च॒ति प्र॒ति॒ष्ठि॒त्यै प्र॒ति॒ष्ठि॒त्यै मु॒ञ्च॒ति वि वि मु॒ञ्च॒ति प्र॒ति॒ष्ठि॒त्यै ॥  
 मु॒ञ्च॒ति प्र॒ति॒ष्ठि॒त्यै प्र॒ति॒ष्ठि॒त्यै मु॒ञ्च॒ति मु॒ञ्च॒ति प्र॒ति॒ष्ठि॒त्यै यया॑ यया॑ प्र॒ति॒ष्ठि॒त्यै मु॒ञ्च॒ति मु॒ञ्च॒ति  
 प्र॒ति॒ष्ठि॒त्यै यया॑ ॥  
 प्र॒ति॒ष्ठि॒त्यै यया॑ यया॑ प्र॒ति॒ष्ठि॒त्यै प्र॒ति॒ष्ठि॒त्यै यया॑ रज्ज्वा॑ रज्ज्वा॑ यया॑ प्र॒ति॒ष्ठि॒त्यै प्र॒ति॒ष्ठि॒त्यै  
 यया॑ रज्ज्वा॑ । प्र॒ति॒ष्ठि॒त्या इति॑ प्र॒ति॒ऽस्ति॒त्यै ॥  
 यया॑ रज्ज्वा॑ रज्ज्वा॑ यया॑ यया॑ रज्ज्वौ॒त्त॒मा॒मु॒त्त॒मां॑ रज्ज्वा॑ यया॑ यया॑ रज्ज्वौ॒त्त॒मा॒म् ॥  
 रज्ज्वौ॒त्त॒मा॒मु॒त्त॒मां॑ रज्ज्वा॑ रज्ज्वौ॒त्त॒मां गां॑ गा॒मु॒त्त॒मां॑ रज्ज्वा॑ रज्ज्वौ॒त्त॒मां गा॒म् ॥  
 उ॒त्त॒मां गां॑ गा॒मु॒त्त॒मा॒मु॒त्त॒मां गा॒मा॒जे॒दा॒जे॒द्वा॒मु॒त्त॒मा॒मु॒त्त॒मां गा॒मा॒जे॒त् । उ॒त्त॒मा॒मि॒त्यु॒त्त॒मा॒म् ॥  
 गा॒मा॒जे॒दा॒जे॒द्वां गा॒मा॒जे॒त्तां ता॒मा॒जे॒द्वां गा॒मा॒जे॒त्ताम् ॥  
 आ॒जे॒त्तां ता॒मा॒जे॒दा॒जे॒त्तां भ्रातृ॑व्याय भ्रातृ॑व्याय ता॒मा॒जे॒दा॒जे॒त्तां भ्रातृ॑व्याय ।  
 आ॒जे॒दि॒त्या॒ऽअ॒जे॒त् ॥

तां भ्रातृ॑व्याय भ्रातृ॑व्याय तां तां भ्रातृ॑व्याय प्र प्र भ्रातृ॑व्याय तां तां भ्रातृ॑व्याय प्र ॥  
 भ्रातृ॑व्याय प प भ्रातृ॑व्याय भ्रातृ॑व्याय प हि॒ण॒या॒द हि॒ण॒या॒त्प भ्रातृ॑व्याय भ्रातृ॑व्याय प

# What is needed?

from *Avoiding Technological Quicksand: Finding a Viable Technical Foundation for Digital Preservation* by Jeff Rothenberg January 1998

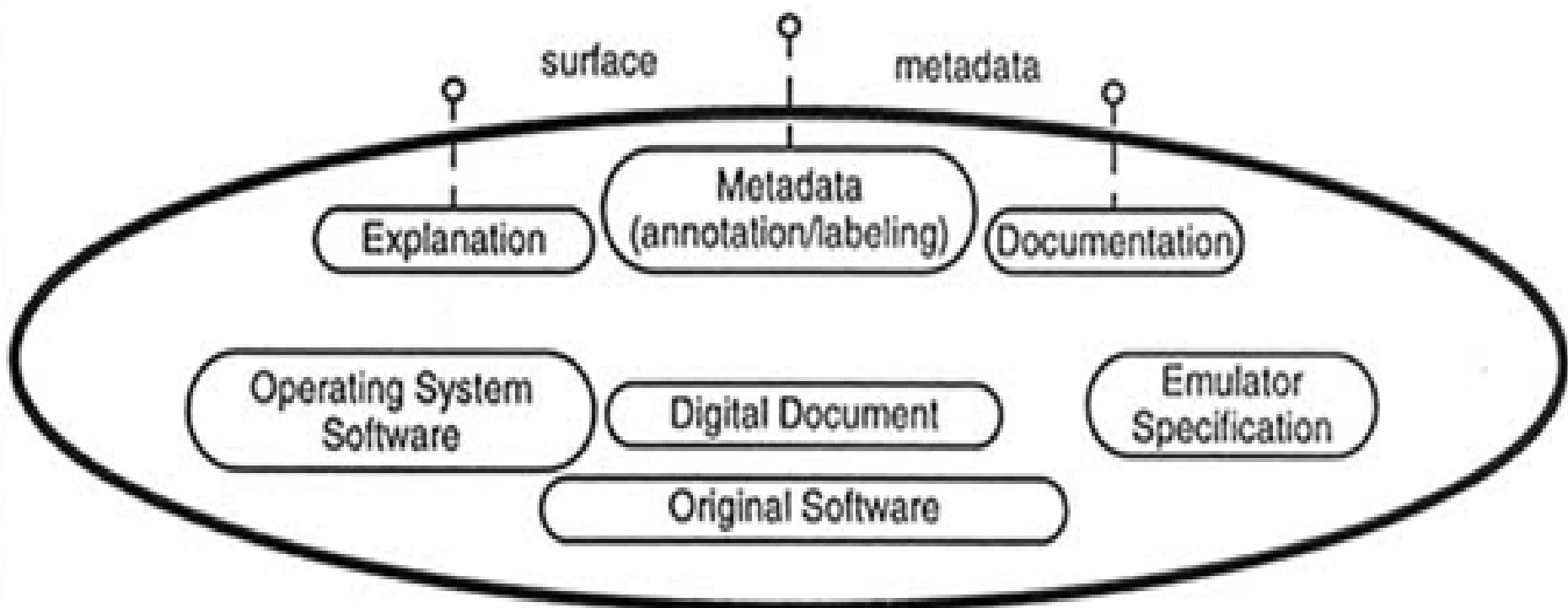


Figure 1: An encapsulated digital document