

# Storage Systems

## NPTEL Course

### Jan 2012

(Lecture 15)

## K. Gopinath

## Indian Institute of Science

# 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 । हि॒णो॒ति  
 ॥ 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 । भा॒ग॒धे॒ये॒ने॒ति 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 हि॒णु॒या॒त् हि॒णु॒या॒त् भ्रातृ॑व्याय भ्रातृ॑व्याय प

# What is needed?

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

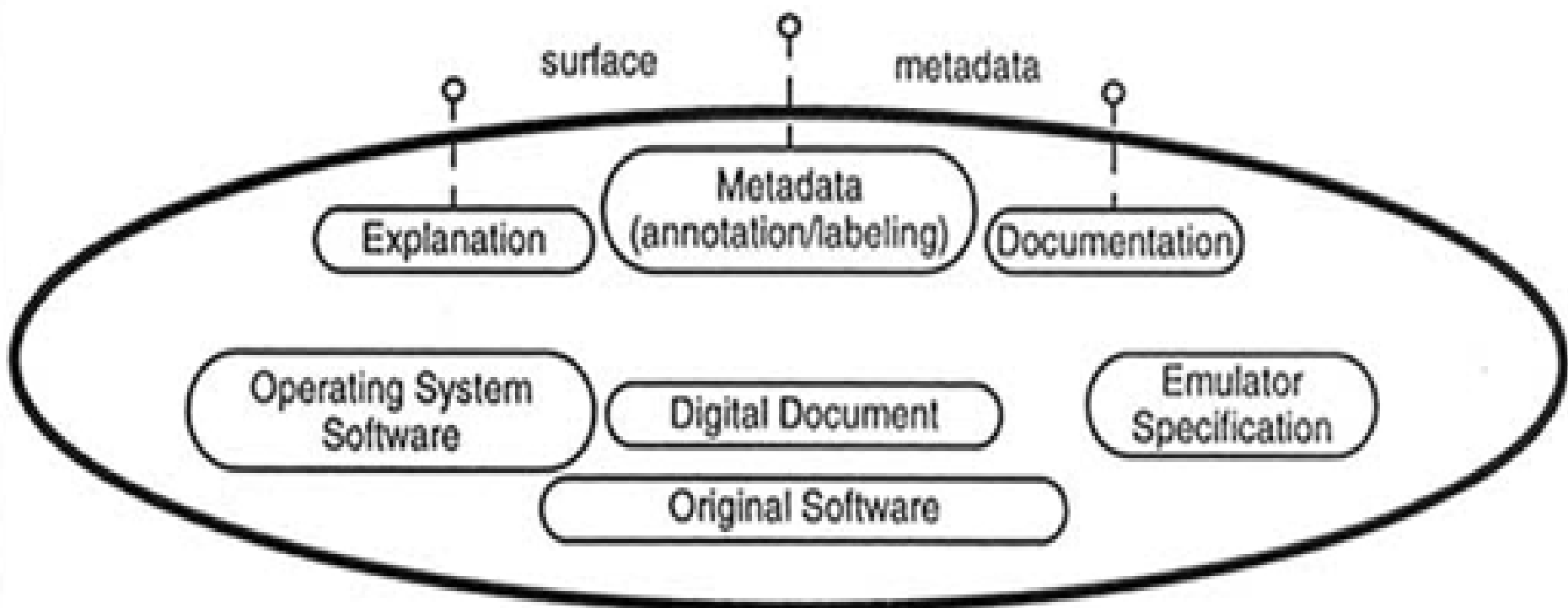


Figure 1: An encapsulated digital document

# What is needed? (contd)

*from Avoiding Technological Quicksand: Finding a Viable*

*Technical Foundation for Digital Preservation by Jeff Rothenberg*

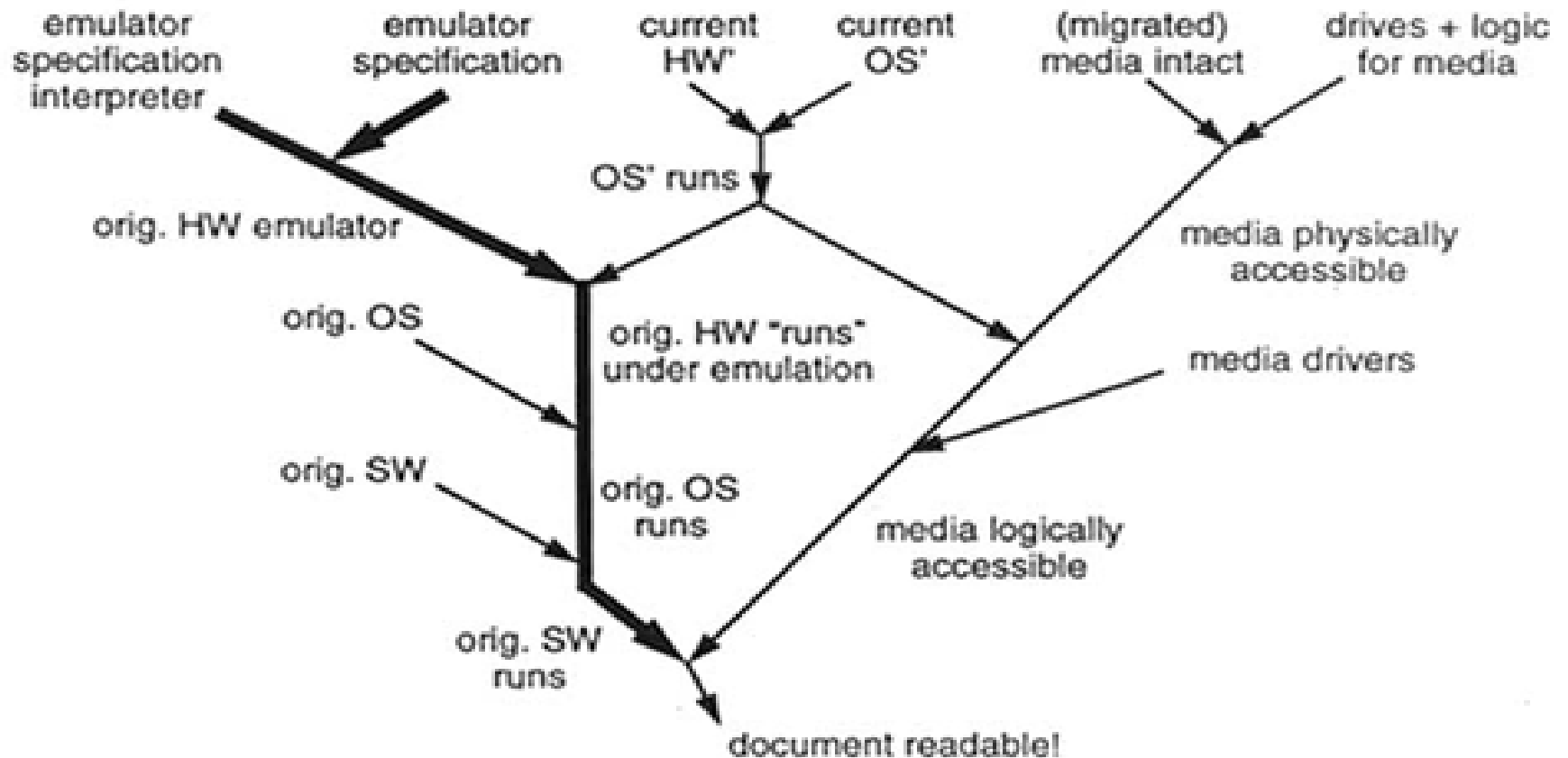


Figure 2: Using emulation to read an obsolete digital document

# Film Archival

- Data preservation requires active energy to move it from one format to the next new generation format
  - needs to be done every few years
  - cost is so high that many movies shot digitally stored in analog form as a fallback in case the migration unsuccessful
- Storing a digital master record of a movie costs about \$12,514 a year, versus \$1,059 to archive a conventional film master in a salt/limestone mine
  - US Academy of Motion Picture Arts & Sciences after a yearlong study of digital archiving in the movie business (2007)

# DNA storage (2013)

- Longevity of Data problem:
  - Hardware changes all the time!
  - But not DNA structure!
  - Very compact too
- 2013: Use ternary encoding (A, T, C with G for breaking long seqs of A/T/C as they get missequenced often)
  - For addl redundancy: 100-bases-long DNA inserts staggered by 25 bases so that consecutive fragments have a 75-base overlap
  - storage density about 2.2 Petabytes per gram
    - enough DNA to recover the data about ten additional times

Towards practical, high-capacity, low-maintenance information storage in synthesized DNA (Nick Goldman et al Nature'13)



# Insurable Storage

- Digital Documents as insurable property
  - Provide economic incentives for storage service producers and consumers to jointly create a marketplace for a diversity of differentially-priced services
- Insurable Storage Services

# 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!?

# Conclusion

- Wide variety of storage designs
  - Each requires different combinations of sw/hw
- Many are still in research stage...

# Interfaces to Storage

- Note that Unix tries to treat everything as a file!
  - Hence some “file” ops might not have anything with storage ops
- Libc
  - `<stdio.h>` in C/Unix but may be arbitrary (for eg, S3 or in parallel user FS)
  - 3 types of buffering:
    - unbuffered, block buffered, and line buffered (default block)
  - FILE structure: a file descriptor, current stream position, eof/error indicators, a pointer to the stream's buffer, if applicable
  - May be able to redefine some functions or load a different shared obj
  - User level fs vs kernel
- System call level: POSIX.1-2008 now
  - Inode (metadata about file)
  - Filesystems may be loadable at runtime (except that contains root image)
- Device Driver
  - Typically loadable at runtime
- Device
- Implications for new models of security/perf...
  - Encryption
  - Mandatory access control/Info flow models
  - Compression

# libc

- File access
- *fopen* opens a file
- *freopen* opens a different file with an existing stream
- *fflush* synchronizes an output stream with the actual file
- *fclose* closes a file
- *setvbuf* sets the buffer and its size for a file stream
  - *setbuf* sets the buffer for a file stream
- *fwide* switches a file stream between wide character I/O and narrow

## Direct input/output

- *fread* reads from a file
- *fwrite* writes to a file

# libc

- Unformatted input/output
- *fgetc getc fgetwc getwc* reads a byte/wchar\_t from a file stream
- *fgets fgetws* reads a byte/wchar\_t line from a file stream
- *fputc putc fputwc putwc* writes a byte/wchar\_t to a file stream
- *fputs fputws* writes a byte/wchar\_t string to a file stream
- *getchar getwchar* reads a byte/wchar\_t from stdin
- *gets* reads a byte string from stdin (deprecated in C99, obsoleted in C11)
- *putchar putwchar* writes a byte/wchar\_t to stdout
- *puts* writes a byte string to stdout
- *ungetc ungetwc* puts a byte/wchar\_t back into a file stream

# libc

- Formatted input/output
- *scanf fscanf sscanf wscanf fwscanf swscanf* reads formatted byte/wchar\_t input from stdin, a file stream or a buffer
- *vscanf vscanf vsscanf vwscanf vfwscanf vswscanf* reads formatted input byte/wchar\_t from stdin, a file stream or a buffer using variable argument list
- *printf fprintf sprintf snprintf wprintf fwprintf swprintf* prints formatted byte/wchar\_t output to stdout, a file stream or a buffer
- *vprintf vfprintf vsprintf vsnprintf vwprintf vfwprintf vswprintf* prints formatted byte/wchar\_t output to stdout, a file stream, or a buffer using variable argument list
- *perror* writes a description of the current error to stderr

# libc

- File positioning
- *ftell* returns the current file position indicator
- *fgetpos* gets the file position indicator
- *fseek* moves file position indicator to a specific location in a file
- *fsetpos* moves file position indicator to a specific location in a file
- *rewind* moves the file position indicator to the beginning in a file

## Error handling

- *clearerr* clears errors
- *feof* checks for the end-of-file
- *ferror* checks for a file error

## Operations on files

- *remove* erases a file
- *rename* renames a file
- *tmpfile* returns a pointer to a temporary file
- *tmpnam* returns a unique filename



# A Common System Interface: Posix 1

- *access* Tests for file accessibility
- *chdir* Changes current working directory
- *chmod* Changes file mode
- *chown* Changes owner and/or group of a file
- *close* Closes a file
- *closedir* Ends directory read operation
- *creat* Creates a new file or rewrites existing one
- *dup* Duplicates an open file descriptor
- *dup2* Duplicates an open file descriptor
- *execl* Executes a file
- *execle* Executes a file
- *execvp* Executes a file
- *execv* Executes a file
- *execve* Executes a file
- *execvp* Executes a file
- *\_exit* Terminates a process
- *fcntl* Manipulates an open file descriptor
- *fdopen* Opens a stream on a file descriptor
- *fork* Creates a process
- *fpathconf* Gets config variable for an open file
- *fstat* Gets file status
- *getcwd* Gets current working directory
- *link* Creates a link to a file
- *lseek* Repositions read/write file offset
- *mkdir* Makes a directory
- *mkfifo* Makes a FIFO special file
- *open* Opens a file
- *opendir* Opens a directory
- *pathconf* Gets config variables for a path
- *pipe* Creates an interprocess channel
- *read* Reads from a file
- *readdir* Reads a directory
- *rename* Renames a file
- *rewinddir* Resets the readdir() pointer
- *rmdir* Removes a directory
- *stat* Gets information about a file
- *umask* Sets the file creation mask
- *unlink* Removes a directory entry
- *utime* Sets file access & modification times
- *write* Writes to a file