"**Diffused Knowledge Immortalizes Itself**"

**Sir James Mackintosh 1765-1832**

# MOTIVATION

http://lockss.stanford.edu

LOTS OF COPIES KEEP STUFF SAFE

# Paper Library System

Libraries act for their institution to

– Acquire copies of important "stuff"

– Keep copies on shelves

– Give access to local readers

Libraries cooperate to

– Supply copies to other libraries

- a reader can easily to find *a* copy

- a "bad guy" has trouble finding and destroying all copies

# Paper Library System

Libraries ensure content persists simply
by supporting their local communities

*A cooperative, affordable,
decentralized, 'archive system'
with LOTS OF COPIES*

# LOCKSS "Library System"

Libraries act for their institution to

– Acquire copies of important "stuff"

– Keep copies in transparent web caches

– Give access to local readers

Libraries cooperate to

– Detect and repair damage

- a reader can easily find *a* copy
- a "bad guy" has trouble finding and destroying all copies

# LOCKSS "Library System"

Libraries ensure content persists simply by supporting their local communities

*A cooperative, affordable, decentralized, 'archive system' with LOTS OF COPIES*

Long Lived: *slow, determined, indestructible*

# LOCKSS

- Open source
- Peer to peer
- Persistent access preservation system
- Web delivered information

Production:  Released April 2004
Support:  Mellon, NSF, Stanford Libraries
Software:   www.sourceforge.net
Teams: Production and Research

# Research Team

Stanford, Harvard, HP Labs & Intel Labs

Award winning research: ACM 2004

- Best Paper SOSP
- Grand Finals 2$^{nd}$ place all student research
- 7 other key research papers accepted so far

Investigating LOCKSS communication

- Scaling, attack resistance

Production Team:

deploys findings, builds system

LOCKSS

**Partners**

*50+ Publishers*

*90+ Libraries*

# LIBRARIES and PUBLISHERS

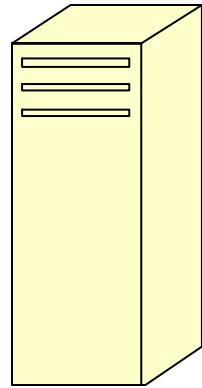™ "...let us save what remains: not by vaults and locks which fence them from the public eye and use in
nsigning them to the waste of time, but by such a multiplication of copies, as shall place them beyond the reach of accident." Jefferson, Thomas.
1] 1984. Thomas Jefferson to Ebenezer Hazard, Philadelphia, February 18, 1791. In Thomas Jefferson: Writings: Autobiography, Notes on the State of Virginia, Public and Private Papers, Addresses,
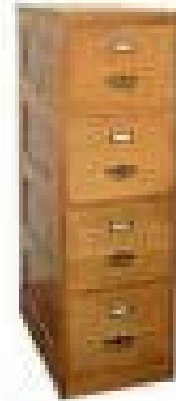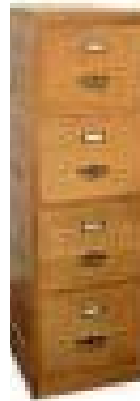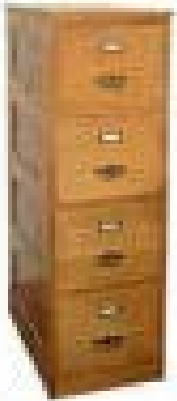rs, edited by Merrill D. Peterson. New York: Library of America

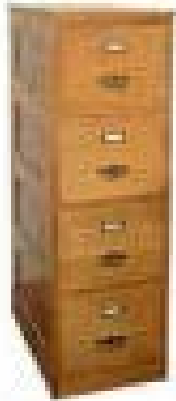| escriptions | Status | Technical Specifications |
|---|---|---|
| ef | Libraries and Endorsing Publishers | Overview |
| blished Papers | | Security |
| ss Releases & News Articles | **April 5, 2004 - Production software released** | Network Integration |
| | | The Plug-in |
| | **ask us to send you the software and instructions** | OAIS |
| | | Research - FAQ |
| ublishers | Librarians | LOCKSS Alliance |
| blisher Actions | Collection Development | Description |
| | Humanities Project | |
| | Title Registry | |
| | User Interface Demo | |
| elated Work | Frequently Asked Questions | Government Documents |
| er Work | Frequently Asked Questions | LOCKSS-DOCS |
| | Software | Credits |
| | License | Funders & Technical Wizards |

**LOCKSS software turns a PC**
into a persistent web cache
into a preservation tool

1 PC holds
~2,500 e-j
years

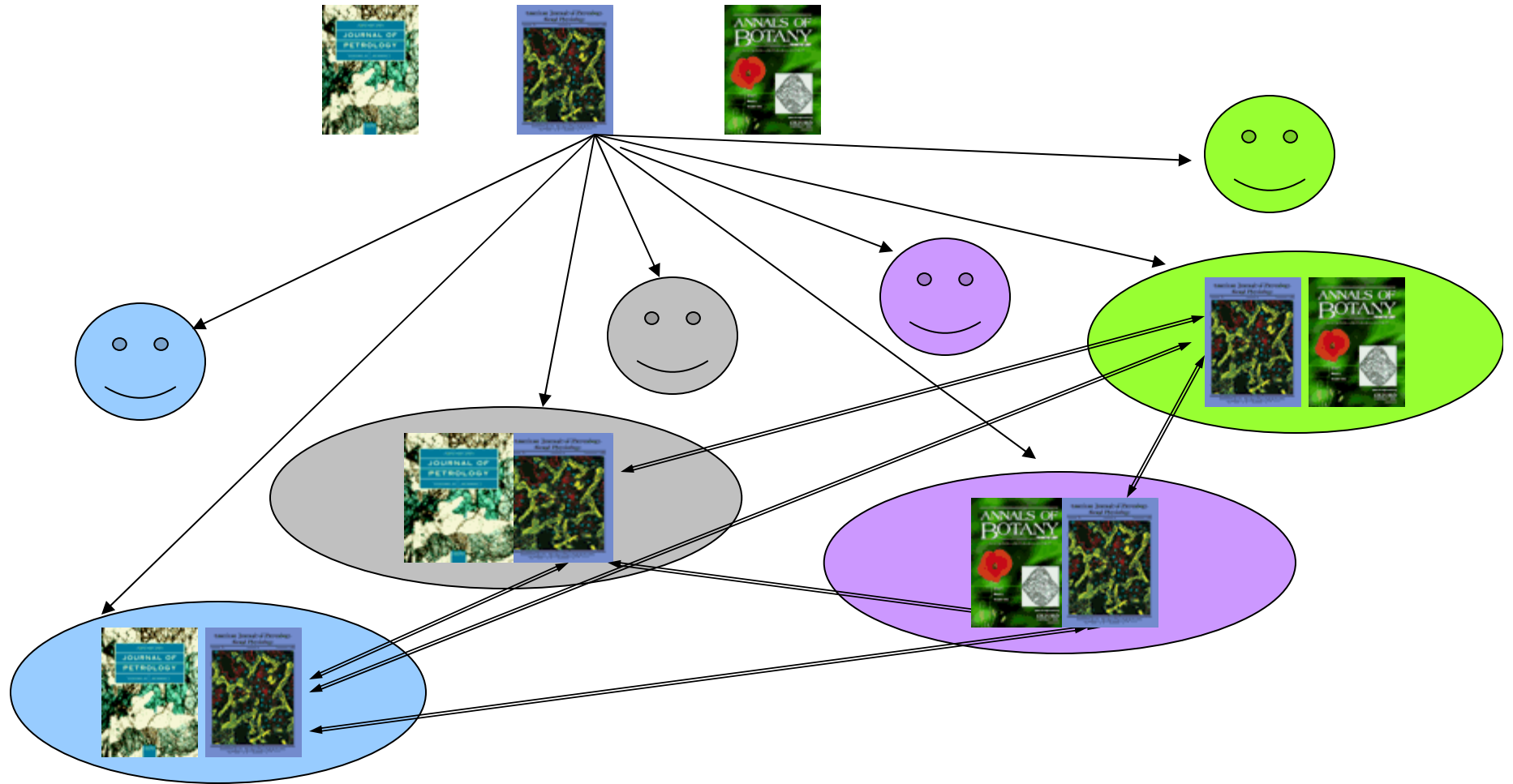600MHz-128MB RAM-Bootable CD drive-Floppy disk drive

LOCKSS

# LOCKSS Caches

- Crawls and collects HTTP content
    - All formats (PDF, HTML, JPEG, TIF, Audio, Video)
- Preserves content integrity
    - Independent collection
    - Cooperate to audit and repair damage
- Provides access
    - Via web browser
    - Content is never "dark"
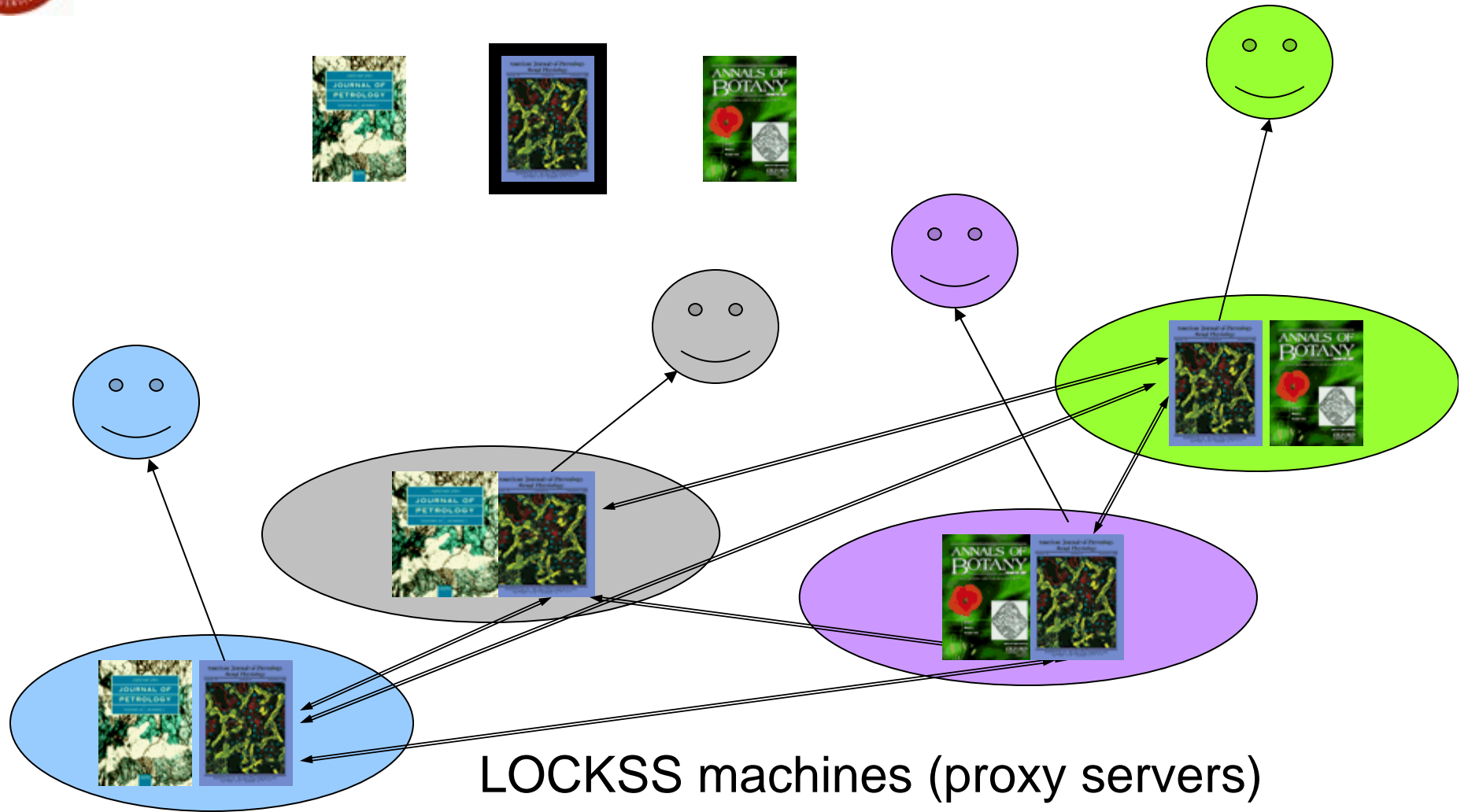
# Approximate Data Flows

## LOCKSS machines

# Approximate Data Flows



LOCKSS machines (proxy servers)

*Prevent the publisher from revoking access rights to back content*

# You're Crazy

A research library's serial collection on a PC
?

# Hardware Costs
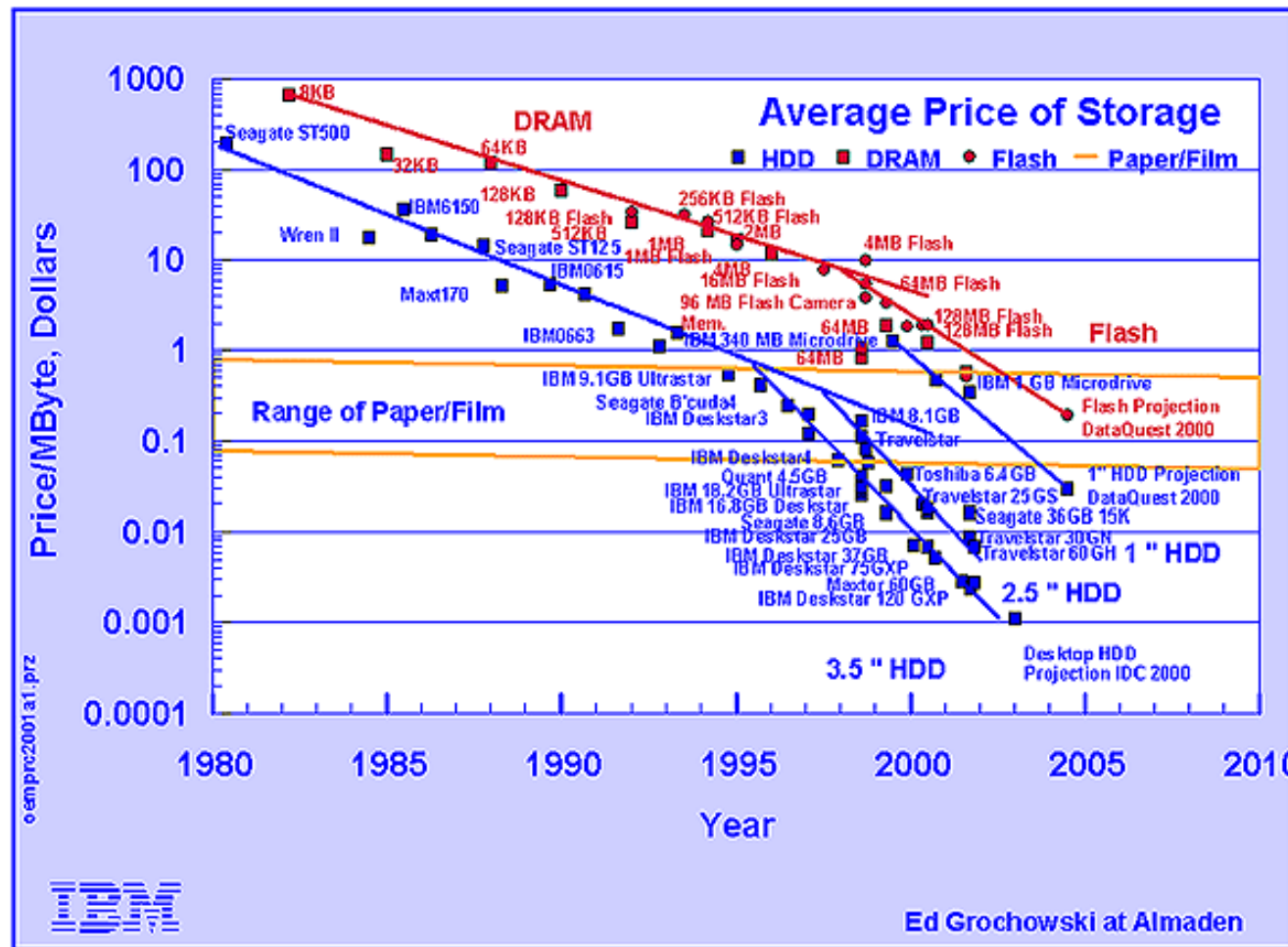
*HDD prices decline by 50% a year*



http://www.almaden.ibm.com/sst/html/leadership/g05.htm

# Terabytes of E-Journals

Median e-journal size is less then  0.5 GB/ year

1 Terabyte (1000 GB) = 2000 journal years

|      | J-yr storage | TB/PC | J-yrs/PC |
|------|--------------|-------|----------|
| 2004 | $0.35        | 1.44  | 2,880    |
| 2005 | $0.28        | 2.88  | 5,760    |
| 2006 | $0.14        | 5.76  | 11,520   |
| 2007 | $0.07        | 11.52 | 23,000   |

1 terabyte for $1,199.00

# Look and Feel to Readers

Configure LOCKSS as a web proxy

Example:

– PNAS table of contents page

- from web (9/11/02)
- from LOCKSS cache

NAS Online

ARCHAEA   **An International Microbiological Journal**

ME | HELP | FEEDBACK | SUBSCRIPTIONS | ARCHIVE | SEARCH | TABLE OF CONTENTS

itution: **STANFORD UNIV MED CENTER** || Sign In as Member / Individual

**t to be notified by email when new content goes on-line?** [Sign up for eTOCs]

'ther Issues: ⬅ ➡

.ble of Contents: **Jan 2 2001; 98 (1)** [Index by Author] [Cover]

**COMMENTARIES**

**PERSPECTIVES**

| **Physical Sciences:** | **Biological Sciences:** | Immunology |
|---|---|---|
| Mathematics | Biochemistry | Medical Sciences |
| Statistics | Biophysics | Microbiology |
| | Cell Biology | Neurobiology |
| **Social Sciences:** | Developmental Biology | Physiology |
| Anthropology | Ecology | Plant Biology |
| | Evolution | |
| | Genetics | |

**CORRECTIONS**

**Find articles in this issue containing these words:**

[          ]  Enter  [Browse & Search All Issues]

ee an article, click its [Full Text] link. **To review many abstracts**, check the boxes to the left of the titles you want, and click the 'Get All Checked Abstract(s)' button. **To see one** ract at a time, click its [Abstract] link.

ear |    Get All Checked Abstract(s)

# NAS Online

ME | HELP | FEEDBACK | SUBSCRIPTIONS | ARCHIVE | SEARCH | TABLE OF CONTENTS

**nt to be notified by email when new content goes on-line?** [Sign up for eTOCs]

Other Issues: ⬅ ➡

# ble of Contents: Jan 2 2001; 98 (1) [Index by Author] [Cover]

**Find articles in this issue containing these words:**

[                    ] Enter | [Browse & Search All Issues]

# What to Collect and Preserve?

- E-Journals
  - Titles you've paid for and are leasing
  - Freely available titles
- Other genres
  - Newspapers, Gov Docs

*http delivered  - serial  - stable URLs*
*– authoritative version*

# Easy for publishers to participate

Publisher give permission (copyright materials) to:

- Libraries
- LOCKSS crawler

*Blanket license permissions*

*no individual library negotiations*

# Publisher License

## Permit libraries

- Collect materials as published for preservation

- Use material consistent with original license terms

- Provide copies for audit and repair to other caches only if they've had copy in the past

**JOURNAL OF**
**Histochemistry & Cytochemistry**

HOME  HELP  FEEDBACK  SUBSCRIPTIONS  ARCHIVE  SEARCH

## Archive of 2003 Online Issues:

### ← 2003 →

| January | February | March |
|---|---|---|
| **January**; 51 (1): 1 - 132 | **February**; 51 (2): 133 - 267 | **March**; 51 (3): 271 - 404 |
| April | May | June |
| **April**; 51 (4): 407 - 554 | **May**; 51 (5): 555 - 696 | **June**; 51 (6): 697 - 852 |
| July | August | September |
| **July**; 51 (7): 853 - 980 | **August**; 51 (8): 981 - 1112 | **September**; 51 (9): 1113 - 1248 |
| October | November | December |
| **October**; 51 (10): 1249 - 1391 | **November**; 51 (11): 1393 - 1574 | **December**; 51 (12): 1575 - 1712 |

# Distributed Repository Model Technology

## Uses many "unreliable repositories" (PCs)

- Robustness through redundancy

- Inexpensive consumer hardware

- Low sys admin overhead (less 1 hour/mo)

## Leverages web technology

- HTTP delivered and displayed content, all formats

- No need to replicate publisher's system

- Automated content ingestion over time

## No single point of failure

# Distributed Repository Model Business

## Costs shared widely

- Total system is never a line item
- Low management overhead
- Low capital cost

## IP issues simplified

- Straight forward blanket license terms
- No "negotiated" access
- Locally owned collections

## No single point of failure

*Budget cuts = key threat to long term access*

# LOCKSS and "Central Repositories"

## Benefits

- System stability improves with some reliable peers
- Diversity improves reliability and attack-resistance

## Requirements

- Implement LOCKSS repository interface
- Run system on mega-servers
- More metadata may be needed for access

# LOCKSS Alliance

## Publishers and libraries work together

- Define policies and best practice
- Develop and share technology
- Share core team costs
  - For limited time, to give model a chance
  - Contributions not required to participate, but
  - Critical amount of support required
  - Suggested contributions on web site

# Taking Action

LOCKSS Program

- is in a nascent stage of development
- needs the community's support to go forward
- shows great promise

There are few actions librarians can take now to preserve digital information for future generations.

The risks of going forward are few. The risks of doing nothing are extremely high.

http://lockss.stanford.edu

LOTS OF COPIES KEEP STUFF SAFE

# Frequent Questions

## OAIS

*Formal statement of Conformance to ISO 14721:2003* May 2004

## Format Migration

# Format Migration

Replacing web format takes a long time

– Both servers and browsers to be updated

– Society pays conversion for popular formats

During this long time we can

– Update cache software with converter

– Preserve content in original format

– Convert on output from old to new format

– Rewrite intra-journal links on output

*.jpg to .png test conversion mid 2004*

# Metadata

## Format metadata

- Collected from HTTP headers and the HTML
- Sufficient for browsers (now & near term)
- Demonstrate format migration based on this metadata
- Incorporate Harvard's JHOVE

## Bibliographic metadata

- For Ingest  OAI  metadata crawler.
- For Export  OAI metadata export capability
- Exploring automatically extracting OAI bibliographic metadata from the text

# When HTTP
# is no longer supported as a protocol?

- Servers will export content using old and new transport protocol.
- LOCKSS caches can be upgraded to support both old and new transport protocols

No "flag day on web"

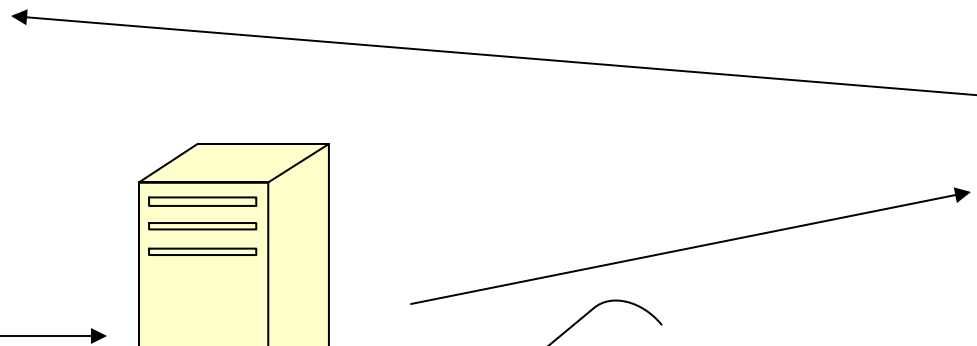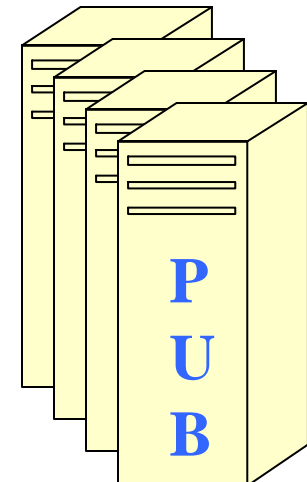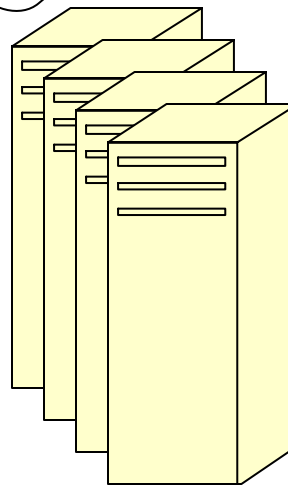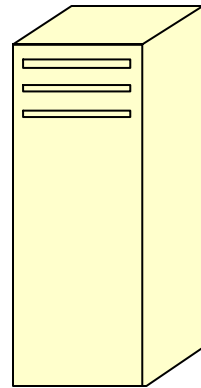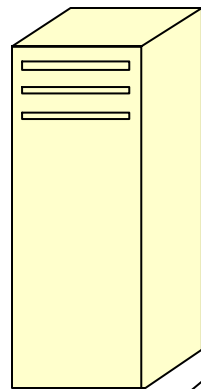Long period of format overlap for common formats

# Collection Access
# LOCKSS and Local Networks
# *publisher is available*

**PAC File**
or **Proxy**

**LOCKSS**

**P U B**

# Collection Access
# LOCKSS and Local Networks
## *publisher is unavailable*

**PAC File**
or **Proxy**

**LOCKSS**

P
U
B