

LOTS OF COPIES KEEP STUFF SAFE

# Distributed Digital Preservation with LOCKSS

Nicholas Taylor ([@nullhandle](#))

Program Manager, [LOCKSS](#) and [Web Archiving](#)

[Stanford Libraries](#)

[Massive Storage Systems and Technology](#)

14 May 2018

# overview

- LOCKSS background
- preservation principles
- distributed preservation
- what's next for LOCKSS?



"LAX on take off" by Doug under [CC BY-NC-ND 2.0](#)





# LOCKSS Background



# (digital) libraries



# lots of copies (were already) keeping stuff safe

- print journal holdings **incidentally resilient:**
  - distributed
  - decentralized
  - irrevocable
  - tamper-evident
  - publisher-independent



“serials” by [Timothy Vollmer](#) under [CC BY 2.0](#)

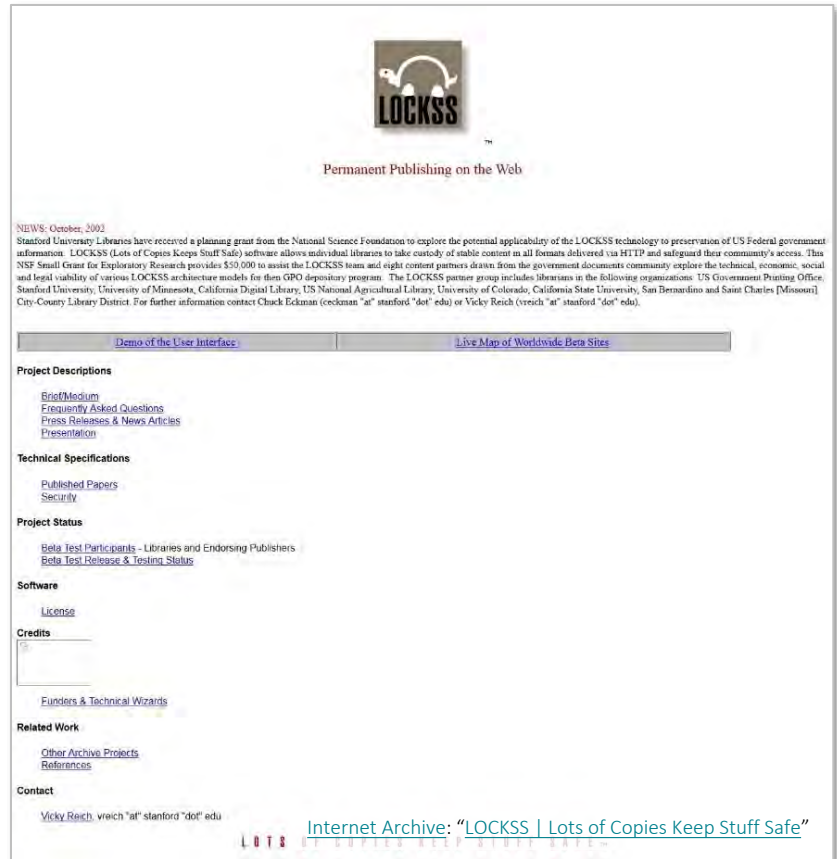






# LOCKSS but for e-journals

- p2p software for e-journal preservation
- restore preservation features of print journal holdings for digital
- re-empower libraries, individually + communally
- improve durability of digital scholarly record



The screenshot shows the LOCKSS website homepage. At the top center is the LOCKSS logo, which consists of a stylized white bird-like figure with its wings spread, positioned above the word "LOCKSS" in a bold, black, sans-serif font. Below the logo is the tagline "Permanent Publishing on the Web".

Below the tagline is a news section dated "NEWS: October, 2002". The text of the news item states: "Stanford University Libraries have received a planning grant from the National Science Foundation to explore the potential applicability of the LOCKSS technology to preservation of US Federal government information. LOCKSS (Lots of Copies Keeps Stuff Safe) software allows individual libraries to take custody of stable content in all formats delivered via HTTP and safeguard their community's access. This NSF Small Grant for Exploratory Research provides \$50,000 to assist the LOCKSS team and eight content partners drawn from the government documents community explore the technical, economic, social and legal viability of various LOCKSS architecture models for their OJO depository program. The LOCKSS partner group includes librarians in the following organizations: US Government Printing Office, Stanford University, University of Minnesota, California Digital Library, US National Agricultural Library, University of Colorado, California State University, San Bernardino and Saint Charles [Minnesota] City-County Library District. For further information contact Chuck Eckman (eckman "at" stanford "dot" edu) or Vicky Reich (reich "at" stanford "dot" edu)." Below the news text are two links: "Demo of the User Interface" and "Live Map of Worldwide Beta Sites".

The main content area is organized into several sections with sub-headers:

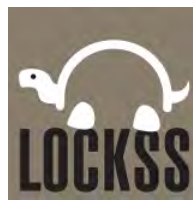
- Project Descriptions**: Includes links for "Brief/Medium", "Frequently Asked Questions", "Press Releases & News Articles", and "Presentation".
- Technical Specifications**: Includes links for "Published Papers" and "Security".
- Project Status**: Includes links for "Beta Test Participants - Libraries and Endorsing Publishers" and "Beta Test Release & Testing Status".
- Software**: Includes a link for "License".
- Credits**: Includes a link for "Fundors & Technical Wizards".
- Related Work**: Includes links for "Other Archive Projects" and "References".
- Contact**: Includes a link for "Vicky Reich, reich "at" stanford "dot" edu".

At the bottom of the page, there is a footer with the text "Internet Archive: 'LOCKSS | Lots of Copies Keep Stuff Safe'" and a small logo that says "LOTS OF COPIES KEEP STUFF SAFE" with a stylized bird icon.



# LOCKSS for more than e-journals

- set out to build **e-journal preservation system**
- ended up building **generic digital preservation core**
- growing number of communities use to **preserve other digital materials**





# community + digital preservation

- communities **complement** LOCKSS:
  - **resilience** against organizational failure
  - native **heterogeneity**
- preservation is an **active** community effort
- lots of **communities** keep stuff safe



["Redwood Canopy"](#) by [Floyd Stewart](#) under [CC BY-NC-SA 2.0](#)





# Preservation Principles

# lots of copies

- intuitive **best practice**
- LOCKSS typically operates w/ **4+ copies**
- **enlist copies to attest** to expected integrity value
- lots of copies **enables**:
  - majority **votes w/ minority** of participating copies
  - **higher-confidence attestations** via landslide agreement



"Untitled" by [Craig Donachy](#) under [CC BY-NC-ND 2.0](#)





# routine audit + repair

- ensuring **long-term bit integrity**
  - must **read data** to know it's good
  - easier to **repair data** sooner
- network nodes **conduct polls** to validate integrity of distributed copies
- more nodes = **more security**
  - more nodes can be **down**
  - more copies can be **corrupted**
  - ...and polls will **still conclude**
- nonces **force re-hashing**
- peers are **mutually-distrusting**



# fail slowly

- fast-operating, tightly-coupled systems **fail quickly**
- LOCKSS is **conservative + sophisticated** about repairs
- polls run slow to enable **detection + mitigation** of cause of damage



“Galápagos giant tortoise on Santa Cruz Island”  
by [John Solaro \(soolaro\)](#) under [CC BY-ND 2.0](#)

# threat model

- familiar threats:
  - media + hardware obsolescence
  - software obsolescence
  - natural disaster
- more typical threats:
  - economic failure
  - organizational failure
  - operator error
- security threats:
  - internal attack
  - external attack





# distributed + decentralized

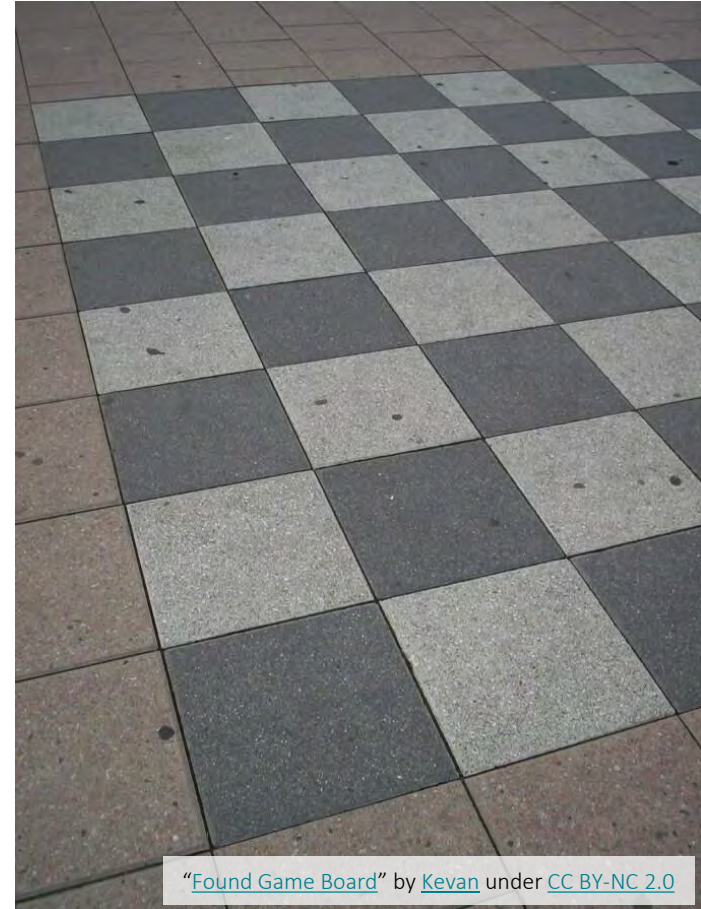
- no monopoly on copy-making
- more copies doesn't mitigate **correlated risk**
- independent, **de-correlated copies**
- minimize central points of failure or vulnerability



"Domino's" by [david pacey](#) under [CC BY 2.0](#)

# no centralized fixity store

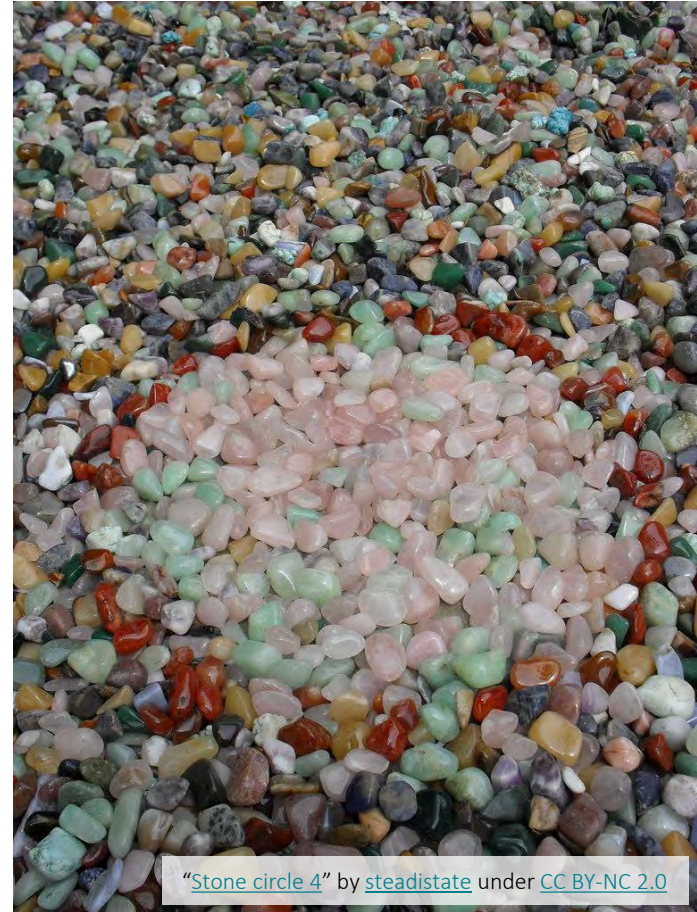
- fixity data **subject to same threats** as data whose integrity it assures
- fixity data is **more vulnerable**, in fact, since more valuable + more centralized
- LOCKSS uses fixity data in **limited ways**



"Found Game Board" by [Kevan](#) under [CC BY-NC 2.0](#)

# local custody

- if preserving data is **core to mission**, LOCKSS helps maintain that competency + commitment in-house
- **unencumbered access** for use by designated community
- **conserving autonomy** + leverage w/ content + service providers
- **jurisdictional transparency** + control



"Stone circle 4" by [steadistate](#) under [CC BY-NC 2.0](#)





# Distributed Preservation

"Catho longtime [explored]" by [Bill Collison](#) under [CC BY-NC 2.0](#)

# where does distributed preservation fit?

- may be **integrated into own infrastructure** (e.g., offsite replication)
- as a **wholly hosted service**:
  - for some, may be **main preservation solution**
  - for others, may **supplement local preservation**



*"Stone stacks"* by [Jack Malvern](#) under [CC BY-NC 2.0](#)



# use cases

- scholarly record
- government documents
- web archives
- collaborative collections
- any types of content **valued in common** by a community





# distributed preservation providers

- **hosted services** w/ varied architectures, service tiers, levels of assurance, replication factors
- replication nodes include **memory orgs + cloud**
- none (including LOCKSS) require **local preservation infrastructure**
- LOCKSS provides **opportunity for co-preservation**



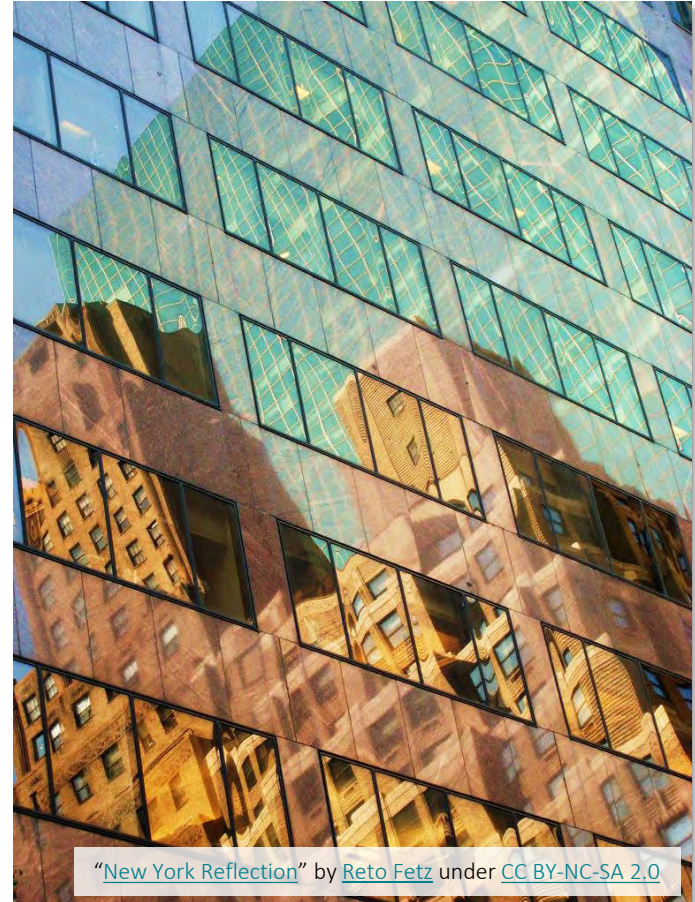


What's Next?

"The two bridges" by [Frank Schulenburg](#) under [BY-SA 2.0](#)

# re-architecture rationale

- de-silo + enable **external integrations**
- foster **developer community**
- capitalize on **work of broader communities**
- create space for **system enhancements**
- evolve w/ web + **digital preservation ecosystem**



["New York Reflection"](#) by [Reto Fetz](#) under [CC BY-NC-SA 2.0](#)



# anticipated outcomes

- functional parity + **backward compatibility**
- **components providing value** outside of end-to-end system
- **better integration** + data hand-offs w/ other apps
- increased use to **preserve repository content**
- increased use to preserve content managed by **non-memory institutions**





# Questions

“Any Questions?” by [Matthias Ripp](#) under [CC BY 2.0](#)