# Archiving "Katrina"

## Lessons Learned

**Kris Carpenter Negulescu, Director**

**Gordon Mohr, Chief Technologist**

**The Internet Archive, Web Group**
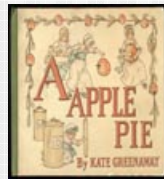
# Agenda

- **Introductions**

- **Project Overview**

- **Lessons Learned**

  - What Worked

  - Challenges & Limitations

- **Technology Review**

- **Future Recommendations**

# What is the Internet Archive?

A digital library of about 3 petabytes of information, including

- Web Pages
- Educational Courseware
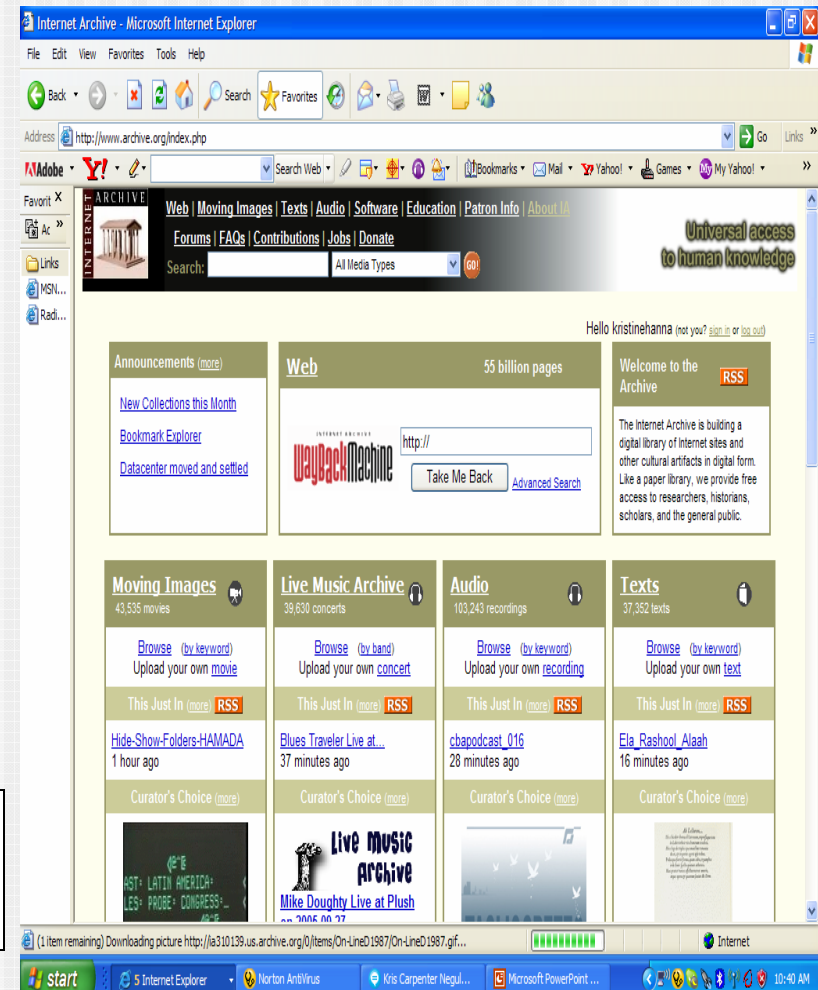- Films & Videos
- Music & Spoken Word
- Books & Texts
- Software

*The archive's combined collections receive 6 million downloads a day*

# How is IA Content Collected?

- **Bi-monthly snapshots of www harvested by Alexa Internet**

- **Web harvests by the Internet Archive**

- **Contributions from the community**

- **Scanning & digitization of public domain texts, stills, moving images, etc.**

**Entire collection accessible for free to the public at www.archive.org**

# Web Archiving Services

# Project Overview

# The 'Katrina' Collection

## The Goal

- To create an historical record of the devastation and the massive relief effort that followed as it was documented on the Web.

## Scope & Timing

- Web content was harvested between September 4 and November 8, 2005
- Collection includes over 61 million unique Web captures, all text searchable, from over 1700 Web sites.

## Key Contributors:

- Library of Congress
- Indiana University, CDL/LSU, University California at Berkeley, and others
- Individuals

**Available at: http://websearch.archive.org/katrina/**

# The Katrina Collection - Tools

- **Heritrix**: a web crawler

  http://crawler.archive.org/

- **Wayback Machine**: an address-based access tool used to locate and view archived web pages

  http://archive-access.sourceforge.net/projects/wayback/

- **NutchWAX**: tools for full text search of archival web content.

  http://archive-access.sourceforge.net/projects/nutch/

NutchWax is an extension of other open source projects:

- **Nutch:** open source web search engine software

  http://lucene.apache.org/nutch/

- **Hadoop:** a framework for running applications on large clusters of commodity hardware.

  http://lucene.apache.org/hadoop/

# Lessons Learned

- Back-Drop
- Defining the Collection
- Harvest & Access

# Back-Drop

- Tight timeline
- No formal submission tools
- No definition of desired collection scope
- No process for recruiting contributors
- No pre-established "permissions to crawl"

# Defining the Collection

# What Worked

**Acquiring Partners & Permissions**

- Able to collect seeds rapidly from a broad range of sources

- Received permission to crawl national and regional press and news sites

# What Worked

## Collection Scope & Quality

Essential contributions came from LoC curators and
subject experts at research universities

- 450+ seeds (~300 from LoC, 100+ from CDL/LSU, and others)
- A subset of the seeds were so unique that they
  would not have been found via a search engine or
  submitted by the general public

# Challenges/Limitations

**Acquiring Partners & Permissions**

- Requests to contribute were distributed ad hoc, artificially limiting participation

- Partners had difficulty convincing experts to drop current work

- More difficult than usual to reach web site operators to acquire permissions

# Challenges/Limitations

## Collection Scope & Quality

- Sans permissions, forced to respect robots.txt

- Curatorial partners and the public submitted blind

- Given time constraints cast "a very big net"
    - Broad expansion of seeds using online directories, blog aggregators, and search engines returned mixed results

# Harvest & Access

# What Worked

**Flexible Action**

- Allocated hardware and engineering resources on short notice to enable timely capture, preservation and access to the collection

- Helped re-establish some Web presences using captures from the IA general archive harvested prior to the event

# What Worked

## Open Source Solutions

- Used Heritrix for rapid, continuous harvests during a six week period

- Automated the tracking and capture of seeds
  - some prior to disappearance of a site; and then when/if it re-appeared

- Nutch/NutchWax enabled full text search
  - 61 million unique captures harvested from 1700 sites

# Challenges/Limitations

## Scalability of Harvesting Tools

- Improved distributed crawl capabilities would have increased the speed, efficiency and effectiveness of crawls

## Prioritized/Parameterized Crawling

- "Smart" content discovery
  - Adaptive revisits and/or selective capture, automated link relationship analysis, geographic clustering, topic based clustering &/or content relationships defined by curatorial resources, etc.

# Challenges/Limitations

## Search & Analysis

- Duplicate management, indexing techniques not yet fully optimized for archived Web content

- Absence of meta data

- Limited options for viewing and analyzing contents

# Technology Review

- Heritrix

- Wayback Machine

- NutchWAX

  with Nutch, Hadoop

# 3 Tasks, 3 Tools

- **Collect**

  Heritrix Web Crawler

- **Redisplay/Browse**

  Wayback Machine Web Archive Viewer

- **Search**

  NutchWAX Web Archive Search Engine

# Heritrix

"Open source, Extensible, Web-scale, Archival
Quality Web Crawling Software"

http://crawler.archive.org

- Collaboratively developed
  - Internet Archive, IIPC, partner libraries, and others
  - First release (0.2) in January 2004; 10 since
- Katrina crawls used version 1.4
  - Latest version: 1.10, released September 2006
- *'Heritrix'* means?   **woman who inherits (heiress)**

# Heritrix Goals

Heritrix was designed to…

- Crawl for completeness/depth
- Be highly configurable – especially in collection scope
- Offer a web-based control interface
- Respect the robots.txt exclusion directives & META robots tags
- Collect material at a measured, adaptive pace, not to disrupt ordinary web site usage

…within the parameters…

- Available under an open source license (LGPL)
- Use Java, leveraging existing open source libraries
- Scale to 100's of millions to a billion documents
- Run well on Linux
  - (but also work wherever Java is available)

# Heritrix

**Current Release: 1.10.1**

**Notable Additions & Changes Since Katrina**

- Ability to distribute a crawl across multiple independent crawlers
- Per-host/domain/queue-grouping collection quotas
- Improved performance and stability in large crawls
- De-duplication add-on
- Expanded crawl configuration options

**Planned Enhancements**

- "Smart" crawling enhancements
- Scaling to ongoing, 1billion+ URL crawls

# Wayback Machine

http://archive-access.sourceforge.net/projects/wayback/

The Wayback Machine is an open source (LGPL) web
archive viewing application in Java designed to
- Display lists of available captures by date
- Allow browsing 'as it was' in a natural manner
- Offer multiple UI modes and index options for different classes of users and deployment sizes

**Current Release: v0.6.0**

**Notable Additions & Changes Since Katrina**
- First open source release in December 2005
- Timeline and proxy modes

# Wayback Machine

# NutchWAX, Nutch, Hadoop

http://archive-access.sourceforge.net/projects/nutch/

NutchWAX includes tools to search Web Archive Collections (WACs). NutchWAX is built on two other open source projects:

- **Nutch**, web-search software based on Lucene Java adding Web specifics such as link-graph analysis, parsers for HTML and other formats, and a batch-oriented crawler.

- **Hadoop**, a framework for running applications (such as the indexing of large scale web content) on clusters of commodity hardware. Hadoop…
  - Implements a computational paradigm named map/reduce
  - Provides a reliable distributed file system on the compute nodes
  - Has been tested on clusters of 600 nodes, but there are reports of support for 900+ nodes for both map/reduce and file system storage.
  - IA's largest cluster for full-text indexing has been 34 nodes.

# NutchWAX

NutchWAX bundles Nutch/Hadoop with extensions for working with archived content, including:

- Adapted Nutch fetcher step to go against archives rather than open net
- Index-time and query-time plug-ins to allow querying of a records' location in a repository
- Awareness of capture-times and multiple captures per URL
- Sending searchers/browsers to an associated Wayback Machine

## Current Release: v0.6.0

### Notable Additions & Changes Since Katrina

- Moved to a map/reduce version of Nutch in May 2006

### Planned Enhancements

- Improved performance and stability for indexing jobs of 100+ million captures

# Recommendations

For Spontaneous, Event-Based, Web Harvests

# Plan Ahead

Establish a broad network of expert contributors & create a communication plan

- Encourage expert participation from research institutions, universities, museums, historical societies, archives, libraries, etc.
  - Secure commitments to participate, where feasible
- Provide incentives for involvement
- Clarify responsibilities for communication & procedures for initiating action

# Plan Ahead

Brainstorm scenarios for event-based, Web harvests

Define criteria that trigger action and standard criteria for collection "quality"

- E.g. Wilma vs. Katrina?, epidemic or isolated threat?, etc.

Create flexible resource pools (hardware, bandwidth, harvesting tools, people, etc.)

# Facilitate
# Contribution & Collaboration

Use Wiki's to organize participation, share information, encourage interaction amongst contributors

Define desired collection "quality" & provide specific guidelines for submission

# Facilitate Contribution & Collaboration

Notify contributors of duplicate seeds as they attempt to submit

Make seed list/s, crawl status, prioritization, scope and frequency transparent to all contributors

Enable experts and curators to edit/expand seed lists, crawl scope, prioritization and frequency

# Thank You!

Kris Carpenter Negulescu, Director

- kcarpenter@archive.org

Gordon Mohr, Chief Technologist

- gojomo@archive.org

The Internet Archive, Web Group