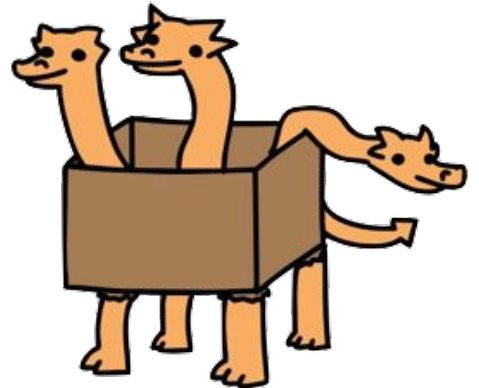# Hydra in a Box:
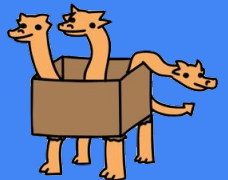## Building A Next-Generation Platform for Digital Collections

Hannah Frost, Stanford University
Gretchen Gueguen, DPLA
Mark A. Matienzo, DPLA
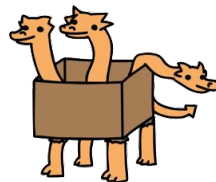DPLAFest 2016 — April 14, 2016

# Project Overview

- A Time for Change
- The Vision
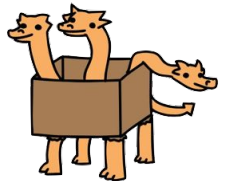- Project Partners
- Project Goals
- Timeline

# A Time for Change

- Conversations between Stanford University, DPLA, and DuraSpace informed project design

- Current digital collections platforms originate in an earlier phase of the web, which explain current limitations

- Infrastructure needs in the DPLA Hub network

  - Legacy systems unable to leverage modern affordances of the web

  - Lack of scalable and sustainable aggregation workflows

  - Lack of support for linked data and metadata enrichment

  - Perceived lack of "obvious choices" for replacement systems for digital collections
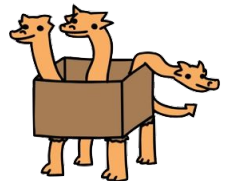
# The Vision

- A product and service that is easy to use, easy to integrate, and that

- Reduce barriers (including cost) to DPLA contribution

- Allow digital collections to be not just **on** the web, but **of** the web

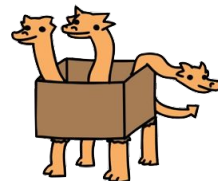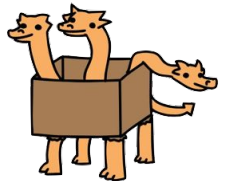- Expand and diversify both the DPLA and Hydra communities

# Project Partners

# Project Goals

- Development of turnkey, Hydra-based application that leverages and improves on core components

- Development/integration of metadata aggregation & enrichment tools

- Connect components with DPLA hubs, current Hydra partners, and prospective Hydra adopters

- Work toward a hosted service

# Timeline

- May 2015-November 2017 (30 months)

- Design process: May 2015-March 2016

- Development: March 2016-November 2017

- Service development and community engagement: throughout project

# Design Phase

Discovery Phase (Fall 2015)

- Literature review and product/service analysis
- Surveys, interviews, and focus groups
- Community outreach

Information Architecture (Winter 2016)

- User requirements and personas
- Requirements - functional & technical
- Models and wireframes

Visual Design (Spring 2016)

# Design Phase



Hydra-in-a-Box Design Process: Tasks & Timeline

Updated 2016-04-07

### Discovery

Phase 1 → Phase 2

Stakeholder goals
Community & landscape surveys

Interviews & focus groups

### Information Architecture

Phase 1 → Phase 2

User personas
Requirements prioritization

Conceptual sitemaps
Wireframes

### Visual Design

Phase 1

Visual design mockups
Visual design style guide

### Infrastructure & Technical Exploration

Phase 1 → Phase 2 → Phase 3

Community engagement
Collaborative development on core gems

Technical prototyping

Gap analysis
Community feedback

### Development

Phase 1 →

Formal development & user testing

| Discovery: Phase 1 | Discovery: Phase 2 | Info Arch: Phase 1 | Info Arch: Phase 2 |
|---|---|---|---|
| | Tech Exploration: Phase 1 | Tech Exploration: Phase 2 | Tech Exploration: Phase 3 |
| Summer | Fall | Winter | Spring |

2015                    2016

# Key Areas of Progress

Design, Requirements and Specifications team:

● Community survey insights

● Analysis of user interviews, focus groups

● Content types requirements for data modeling

# Community Survey

256 complete responses

311 repositories

Mostly small, US academic libraries



Public college/university library
34%

Independent research
library/archives
14%

Museum
10%

Over half of respondents represent
either a public or private college or
university library.

Regional consortium
8%

Historical
society
4%

Private college/university library
20%

Public library
6%

Government

# Survey Insights

- Expectations of our project

- Satisfaction levels
  - Users of hosted services tend to be more satisfied than users of local deployments

- Strengths and weaknesses of existing repository options

- 53% plan to migrate to another system
  - Most to a Fedora-based solution
  - Rest are "not sure" what's next

# User Interviews

- Completed 21 individual or small-group interviews and 4 focus groups
  - 55 individuals in total
  - 46 institutions in the US and Canada
  - 29 hours of recorded content

- Interviews held either in-person or through videoconference; focus groups held in-person

- Coded and analyzed process to further identify potential requirements

# Content Analysis Visualizations

# Interviewee's Notable Quote

"... How many of these different systems do you need? You can have your digital collections with images and documents, you can have your IR, you can have your digital preservation system, and you can add Omeka on top of that to do exhibits. It's just too much to have four or five different systems."
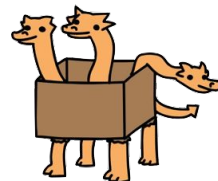
# Content Type Analysis

| | IMAGES | | | | | | | | | TEXT | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | digital photo | digital photo album | scanned poster | scanned postcard | scanned looseleaf manuscript | scanned folder of archival docs | scanned photo album (pages) | scanned scrapbook (pages and clippings) | scanned newspaper | scanned monograph with cleaned OCR | scanned mono with r |
| **OBJECT STRUCTURE** | | | | | | | | | | | |
| single file | yes | no | yes | no | no | no | no | no | no | no | no |
| multi-file, sequenced, flat | no | yes | no | yes | yes | yes | yes | yes | yes | yes | yes |
| multi-file, unsequenced, flat | no | no | no | no | no | no | no | no | no | no | no |
| multi-file, hierarchical | no | no | no | no | no | no | no | yes | yes | no | no |
| | | | | | | | | | | | |
| **DESCRIPTIVE METADATA** | | | | | | | | | | | |
| basic ("core") | yes | yes | yes | yes | yes | yes | yes | yes | yes | yes | yes |
| geo | yes | yes | no | yes | no | no | no | no | no | no | no |
| serial (volume, issue, etc.) | no | no | no | no | no | no | no | no | yes | no | no |
| archival order/hierarchy | yes | yes | yes | yes | yes | yes | yes | yes | no | no | no |
| | | | | | | | | | | | |
| **DELIVER DERIVATIVES** | maybe | no | yes | yes | yes | yes | yes | yes | yes | yes | yes |
| | | | | | | | | | | | |
| **APPLICATION BEHAVIOR (ITEM LEVEL)** | | | | | | | | | | | |
| zoom | yes | yes | yes | yes | yes | yes | yes | yes | yes | yes | yes |
| front/back | no | no | no | yes | yes | yes | no | yes | no | no | no |
| page turn | | | | | | | | | yes | | |

# Early Technical Exploration

- Deploying to the Cloud
  - Leverage services for institutions without local infrastructure

- Simplifying installation and configuration
  - Users should not need to be technical to set up and maintain an instance

- Determining a starting point for application development
  - Build on existing community-based work if possible
  - Sufia 7.0 - actively under development

# Repository Development

- Assembled an all-star technical team
  - 10 Engineers: software development, data modeling, development operations
  - Contributions from other institutions (Penn State, maybe others)
  - Led by Michael Giarlo

- First work cycle: March - June 2016
  - Series of one-week sprints
  - Recorded demos of iterative progress, available to the public

- First milestone: Deploy our application based on Sufia 7 to the cloud
  - Priority content types
  - Configuration UI
  - Administrative dashboard
  - Batch import

# Follow our progress

# Aggregator Needs

- More flexible mapping standard than XSLT

- Ability to harvest from multiple sources

- Reconciliation services that utilize linked data

- Enhanced quality control tools

- Ability to normalize and create consistencies in data values

- Easily get data in and out

- Robust enough to handle multiple feeds and multiple sources

- Processes to move data from one repository to another resemble aggregation workflows

# DPLA's Aggregation System, Heiðrún

Three Main Functions

- Harvest
  - Source agnostic
- Map
  - Mapping DSL expressed in Ruby
  - Maps to RDF triples
- Enrich
  - Modular enrichments written to normalize and enhance data

# Roadmap

- Completing requirements now
- April - July
  - Design remaining infrastructure
  - Develop user interface requirements further
- August - November
  - Develop dashboard tools
  - Analyze convergence points with repository
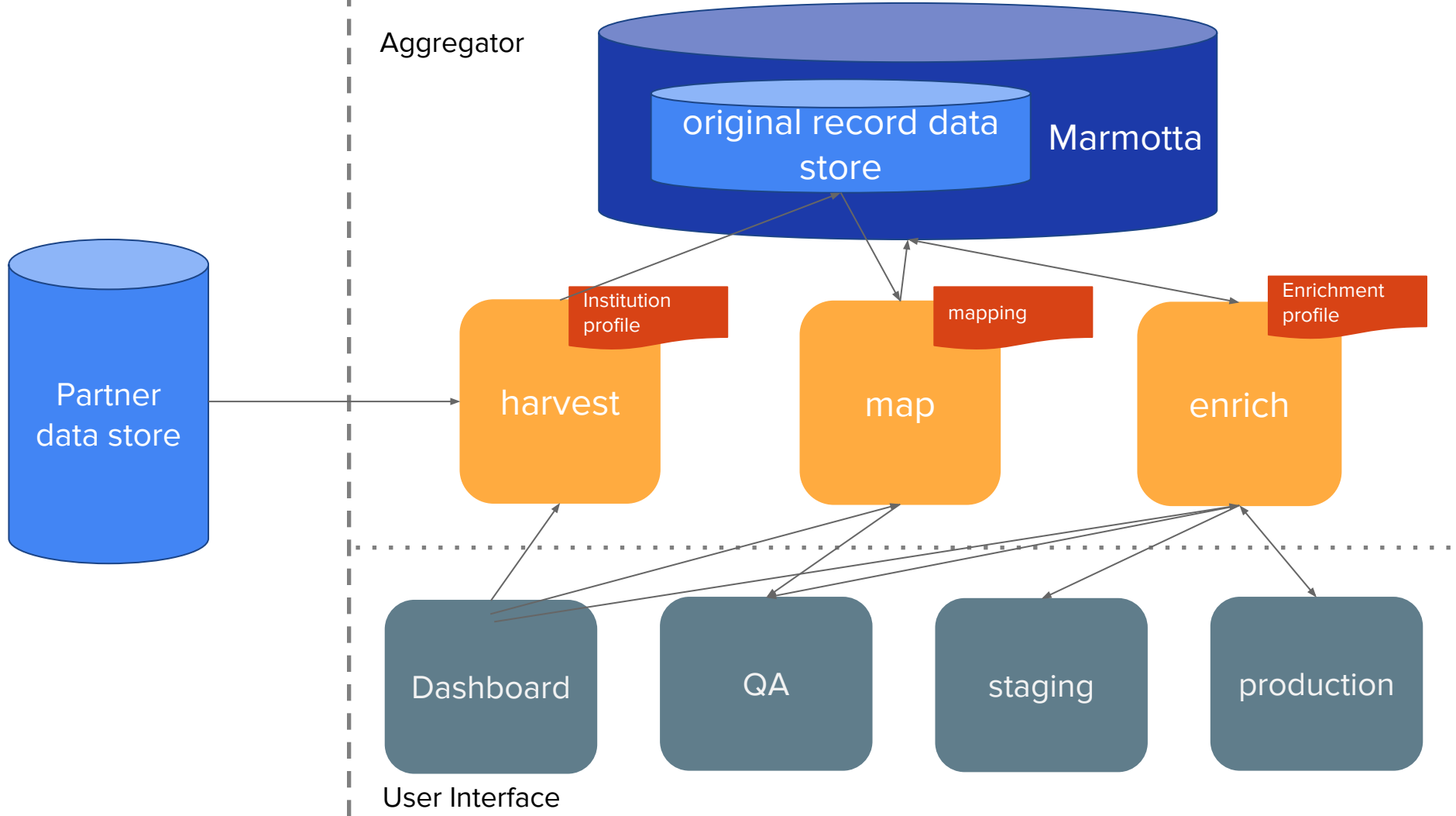  - Plan for improvements to QA interface
  - Begin User Testing
- November - March 2017
  - Develop QA improvements
  - Refine interfaces and infrastructure
  - Implement job scheduling

# Developing a Hosted Service

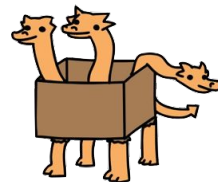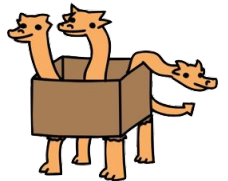- Project partners collaborating to develop requirements for a cloud-hosted service based on the repository product under development

- Market research underway, starting with analysis of information discovered during the design phase

- Evaluating tiered service models depending on needs of potential adopters

- Significant technical work to focus on develop a shared, maintainable, and scalable service

# More Information

Visit our website and blog: **http://hydrainabox.org/**

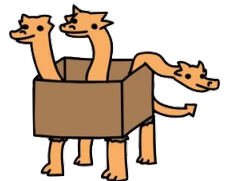Follow us on Twitter: @HydraInABox

Public information list

hybox-info@googlegroups.com

Contact us

hybox-contact@googlegroups.com

**Thank You!**

Hannah Frost
hfrost@stanford.edu
@feefifofannah

Gretchen Gueguen
gretchen@dp.la
@G_AmSpinnrade

Mark A. Matienzo
mark@dp.la
@anarchivist

___

http://bit.ly/dplafest-hybox