

#### From Hydras to TACOs

Harlow, Christina; Fahy, Erin 2018 Dostupný z http://www.nusl.cz/ntk/nusl-380920

Dílo je chráněno podle autorského zákona č. 121/2000 Sb. Licence Creative Commons Uveďte původ-Neužívejte komerčně-Nezpracovávejte 4.0

Tento dokument byl stažen z Národního úložiště šedé literatury (NUŠL).

Datum stažení: 23.05.2024

Další dokumenty můžete najít prostřednictvím vyhledávacího rozhraní nusl.cz .

# From Hydras to TACOs: **Evolving the Stanford** Digital Repository

#### ELAG 2018 Christina Harlow, Erin Fahy



STANFORD UNIVERSITY LIBRARIES

#### **Goals of this Talk**

- 1. Introduce the Stanford Digital Repository
- 2. Discuss our Approach(es) to re-architecting our system
- 3. Introduce SDR3, TACO, & our redesign so far
- 4. What's next?

#### **Goals of this Talk**

- 1. Introduce the Stanford Digital Repository
- 2. Discuss our Approach(es) to re-architecting our system
- 3. Introduce SDR3, TACO, & our redesign so far
- 4. What's next?

#### We'd really love to hear your feedback on this work!

And special thanks to the Bootcamp group that went through a fast-paced deep dive of some of this work on Monday.

## **1. Some Context on SDR**

### **Stanford Digital Repository (SDR)**

Currently in it's second iteration, architecturally (i.e. "SDR2")

Been working for over ten years

Guided by a '3 Spheres Topology'



## **Stanford Digital Repository**

Variety of digital resources & assets:

- Bulk ingest of digitization labs work,
- Institutional repository self-deposit,
- Electronic dissertations & theses self-deposit,
- Geo-datasets,
- Web archiving,
- Electronic resources cataloged & preserved,

...

#### **Stanford Digital Repository Metrics**

Manages roughly **1.6 million** distinct resources currently

Has about half a petabyte (455 TB) of digital assets in our preservation layer

~426 TB of digital assets in our repository staging systems

455 TB of digital assets & 59.1 GB of metadata in our access system(s)



# 2. Our Approach(es) to re-architecting our system



#### **SDR2 'Retrospective'**

- Lack of full system comprehension
- Lots of unmaintained codebases & workflows
- Over-engineered components
- Pain points on adding new features or processes
- Mismatch of design(s) & implementation(s)

"There are a lot of interaction points between layers of the technology stack and you often need to know a lot about all of these interactions even if you are only currently concerned with one part of the stack."



## SDR3 Design Cycle

- 3 months of daily 1 hour meetings with architect, engineers, product owners, administrators, & others
- Produced requirements independent of system expectations
- Built shared understanding of our current needs & conceptual architecture
- In tandem: did a 'current state' deep dive on our existing code
- Generated a high-level conceptual design & plan



#### **Hyrax Analysis: SDR Options**

- 1. Do not use Hyrax at all for SDR3. Non-starter.
- 2. Use Hyrax for SDR3 entirely. However...
  - a. ~38% of our core, reviewed requirements are not covered by Hyrax.
  - b. ~24% of those are 'Maybe', i.e. require config, model changes, or coding.
  - c. Most alignment with UI / Self-Deposit, direction of analytics, web dev.
  - d. Least alignment in overall architecture, bulk processing, back-end needs.
- 3. Integrate Hyrax & SDR3 via components & 'seams'.

#### Hyrax Analysis: Possible 'Seams'

- <u>Valkyrie</u>'s "internal air gap" approach for flexible data stores
- Actor Stack, Sipity, or Delayed Jobs:
  - Write <u>Hyrax MiddlewareStack</u> as seam to our Management API & asynchronous processing.
- Rely on both internal air gaps as well as crisp boundaries via ReST APIs.
  - Independent scalability.
  - Migration 'hinge' for components that don't or shouldn't fit into Hyrax.
  - Keeps separate areas of our work most aligned with the Samvera community:
    - self-deposit & access/discovery currently
    - analytics and administration dashboards in the future

## Fedora 4 / Fedora API Analysis

- Incompleteness & uncertainty of specification work
- Graph store limitations
  - Keep Linked Data out of our back-end system
- Complexity & Comprehensibility
- Performance & Extensibility
- Data & Resource Handling
- System Expectations
- Re-approach Fedora overlap with our data publication (Access) systems





#### **SDR3 Evolutionary Plan**

SDR3 Design Kick-Off (x3) & Hydrox Analysis Phase

(10/03/2017-01/12/2018)

TACO Skeleton Prototype Phase

- TACO Prototype Work Cycle (01/12/2018-04/26/2018)
  - SDR3 Design Iteration (4 month)
- ETDs ⇔ TACO Prototype, Bulk smoke test (3 months)
  - SDR3 Design Iteration (1 month)

TACO Prototype Integration Phase

ETDs 🗢 TACO "go live" & data migration

ETDs 🗢 Hydrus "go live" & data migration

. . .

#### SDR3 High Level Conceptual Design (so far)



#### **TACO Prototype's Work Cycle 1 Goals**

Functional Goals	Technological Goals	Process Goals
Deposit resources (binaries & metadata) into repository via API.	Drive forward Department API specifications, implementations, & practices.	Work towards new core with something visible to limited stakeholders to make it real-er.
Retrieve deposited resources from repository via API.	Test implementation options for our current SDR3 design.	Get feedback on SDR3 design, & check for roadblocks.
Persist resources.	Vet our data models & metadata profiles.	Keep to high-level, extensible functional blocks.
Perform skeletal resource processing (i.e. workflows).	<ul> <li>Test feasibility of possible technologies:</li> <li>Hyrax integration points.</li> <li>Test throughput / scalability.</li> <li>SDR2 &amp; SDR3 analytics.</li> <li>Inform cloud practices.</li> <li><u>Cloud first but Cloud neutral.</u></li> </ul>	Showcase internal / lower stack rewrites needed for moving middle and end-user codebases forward.



#### Go & Docker for TACO Codebase

- Ability to be modular, with APIs as clean boundaries & work in Cloud (AWS).
- Decision to use compiled language coupled Docker for deployment.
- Efficient Docker container deployment with small, executable binaries (as opposed to platforms that require an operating system and server).
- Focusing on compilable language for small, efficient services led us to Go language.



#### Go & Swagger Prototype Codebase

Additional TACO Prototype goals included:

- rapid development and delivery;
- SWAGGER specification support for consistent API to Code translations & share-ability of APIs across languages;
- support for continuous deployment & cloud solutions;
- parallelization fit for horizontal scalability.



#### **AWS Selections for TACO**

- Docker containers for sending off the codebase binary.
- AWS ECS (elastic container service) for running this image.
- CircleCI for Continuous Integration with AWS ECS & Docker due to its use by industry for similar set-ups.
- <u>Terraform</u> for building out AWS infrastructure
- AWS DynamoDb for metadata persistence for the prototype.
  - Very likely to use RDS in production.
- AWS S3 for binaries for the prototype.

#### **Cloud-first but Cloud-neutral**

Our considered & kept-in-mind graceful degradation paths:

- Docker => Docker is reusable.
- AWS ECS => Any system or VM that can run Docker. Docker swarm?
- Swagger 2.0 => Specification Built for Translateability
- Go + go-swagger => Just use Ruby.
- AWS dynamodb => CouchDB or Postgres.
- AWS s3 => File system.
- AWS kinesis => Kafka or Spark Streaming when ready.

#### Kafka / Kinesis?

- Early design had event driven system for managing resource state & asych, DAG-based processing
- Put too much intelligence into TACO
- Kinesis deemed not suitable
- Re-designing to use Kafka-inspired event system within our Provenance & State Service
- Our asychronous, DAG-driven processing inspired by <u>Airflow</u> becomes parallel to SDR3



#### **Special Note: Fedora 4 API vs TACO API**

- TACO API aims to be much simpler than Fedora API.
- Decoupled from Linked Data Platform at this level of our stack.
  - We are supporting JSON / JSON-LD, which allows LD higher up.
- Reduced API calls, leading to increased performance.
  - Up to %50 less if we include ACLs, FileSets & ORE proxy ordering.





Files (Binaries)

{

\$schema: "http://json-schema.org/draft-06/schema#",

#### title: "Digital Repository Object",

description: "Domain-defined abstraction of a 'work'. Digital Repository Objects' abstraction is describable for our domain's purposes, i.e. for management needs within our system.", type: "object",

#### - required: [

"@context",

"@type",

"externalIdentifier",

"label",

"tacoIdentifier",

- "version",
- "administrative",

"access",

"identification",

"structural"

#### 1,

- properties: {

```
- @context: {
```

```
description: "URI for the JSON-LD context definitions.",
type: "string"
```

```
},
```

},

```
- @type: {
```

description: "The content type of the DRO. Selected from an established set of values.",
type: "string",

- enum: [

"http://sdr.sul.stanford.edu/models/sdr3-object.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-3d.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-agreement.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-document.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-document.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-age.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-age.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-age.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-photograph.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-manuscript.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-manuscript.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-manuscript.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-map.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-media.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-media.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-webarchive-binary.jsonld", "http://sdr.sul.stanford.edu/models/sdr3-webarchive-seed.jsonld", TACO Metadata Application Profiles (JSON Schema)

## 4. What's Next?

#### **Current Design Work**

- File system analysis for options & costs
- Asynchronous & batch processing system design work going on
  - Heavily influenced by Apache Airflow
- Metadata efforts have free range approach
  - Starting with a metadata use cases analysis before jumping into schemas / ontologies
  - JSON[-LD] & JSON Schema used for flexibility, separation of external semantics & internal data shapes
- Preparing for next work cycle to revise & connect TACO ultimately to a self-deposit system & a bulk load job

## **Keeping Community Connections**

- Samvera architecture & front-end work re-approach
- Interest in architectural overlaps with FOLIO
- Code4Lib Spark in the Dark overlaps
- Using PCDM, MOAB => OCFL, revisiting other places to share our data specifications
- Blacklight, IIIF, & related Access systems community work untouched
- Looking outside of cultural heritage for community partners & ideas
  - $\circ$  Airflow
  - AWS
- Asking our community friends & experts like ELAG participants for feedback

# Questions or Feedback?

#### cmharlow@stanford.edu @cm\_harlow

https://github.com/sul-dlss-labs/taco/



STANFORD UNIVERSITY LIBRARIES

