# Costs of Curation

INFX 551
Data Curation

- Next week: Summary Lecture + Presentations

- Project Updates

- Lecture & Discussion

- Activity (s)

- Gordon Kennedy - WSDoT

# Policy background for open data in the USA…

WE THE PEOPLE ASK THE FEDERAL GOVERNMENT TO TAKE OR EXPLAIN A POSITION ON AN ISSUE OR POLICY:

# Require free access over the Internet to scientific journal articles arising from taxpayer-funded research.

**Created by J.W. on May 20, 2012**

## Signature Count

65,704 SIGNED                    25,000 GOAL

Read the memorandum here: https://obamawhitehouse.archives.gov/sites/default/files/microsites/ostp/ostp_public_access_memo_2013.pdf

EXECUTIVE OFFICE OF THE PRESIDENT
OFFICE OF MANAGEMENT AND BUDGET
WASHINGTON, D.C. 20503

THE DIRECTOR

May 9, 2013

M-13-13

MEMORANDUM FOR THE HEADS OF EXECUTIVE DEPARTMENTS AND AGENCIES

FROM:      Sylvia M. Burwell
           Director

           Steven VanRoekel
           Federal Chief Information Officer

           Todd Park
           U.S. Chief Technology Officer

           Dominic J. Mancini
           Acting Administrator, Office of Information and Regulatory Affairs

SUBJECT:   Open Data Policy—Managing Information as an Asset

"…this Memorandum requires agencies to collect or create information in a way that supports downstream information processing and dissemination activities."

Impact is twofold:

- Unfunded set of mandates to federal agencies  - both memoranda offer choices and broad time frame for compliance (and don't spell out sanctions).

- Spurs innovation for compliance - the government often says it doesn't want to pick winners or losers - but it does want the game to be played. This is how the open data game started.

SCIENCE PRIORITIES

# Who Will Pay for Public Access to Research Data?

When economic models and infrastructure are not in place to ensure access and preservation, federally funded research data are "at risk."

**Francine Berman[1] and Vint Cerf[2]**

On 22 February, the U.S. Office of Science and Technology Policy (OSTP) released a memo calling for public access for publications and data resulting from federally sponsored research grants (*1*). The memo directed federal agencies with more than $100 million R&D expenditures to "develop a plan to support increased public access to the results of research funded by the Federal Government." Perhaps even more succinctly, a subsequent *New York Times* opinion page sported the headline "We Paid for the Research, So Let's See It" (*2*). So who pays for data infrastructure?

What happens to valuable data when project funding ends? Consider, for example, a 3-year research project in which valuable sensor data are collected from an environmentally sensitive area. Those data may be useful not just for the duration of the project but for the next decade or more to collaborators and a broader community of researchers. For the first 3 years, the costs of stewardship (including development of a database that supports analysis, access to the
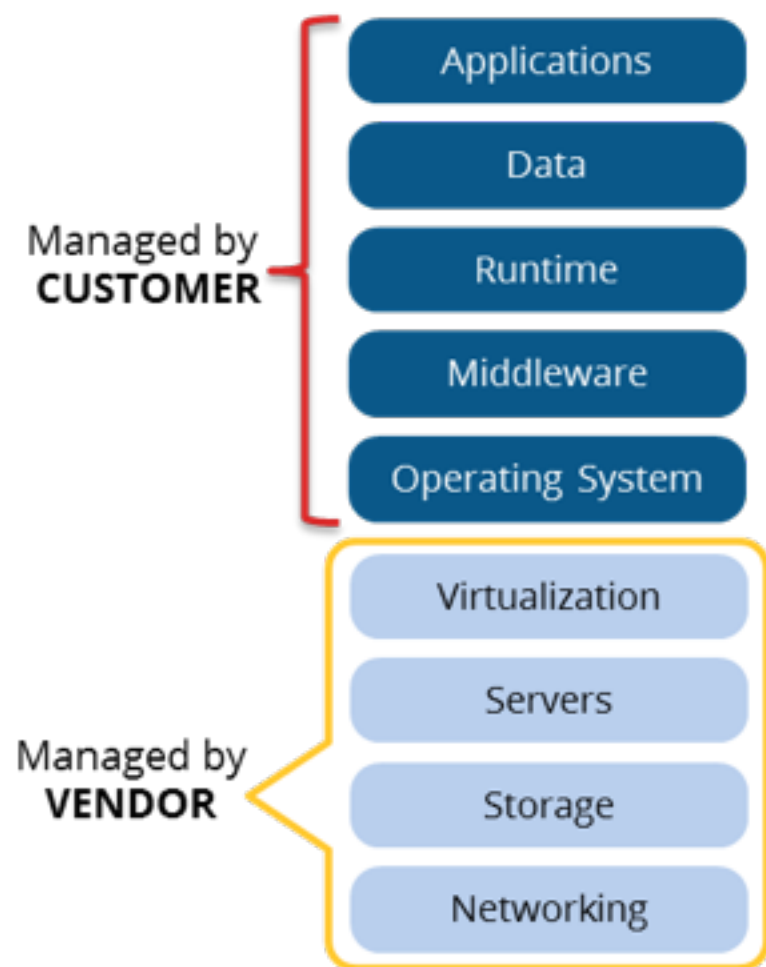
Before we decide who is paying, lets figure out the bill.

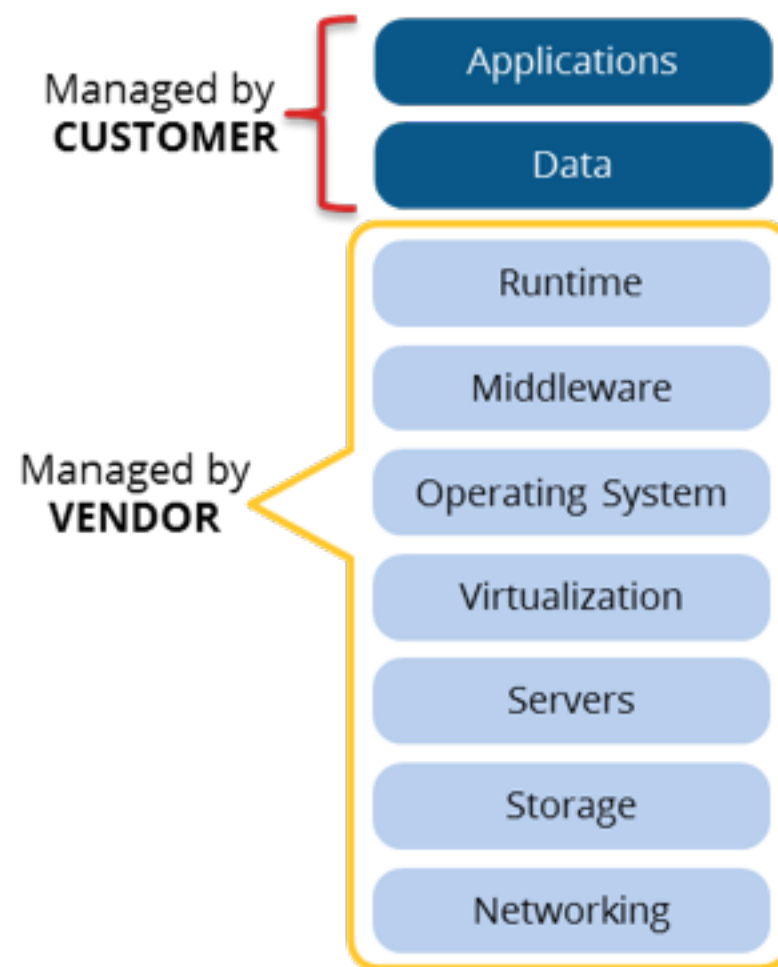What are the costs of running a open data, and curation programs?

- Infrastructure (hardware, software & implementation)

- Governance

- Staff & Skills Development
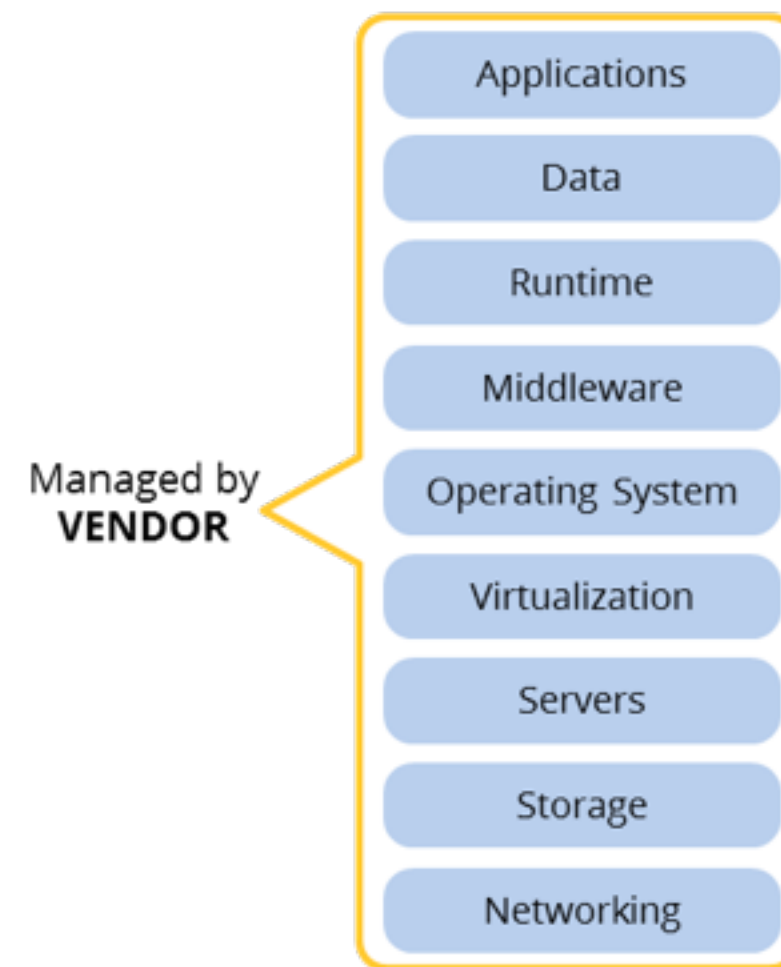
- Maintenance

# Infrastructure

- Open data portal development costs

    - (Open source is free like a pony)

- Cost of storage - including cloud storage, tape backups, and local storage

- Cost of publishing datasets, creating APIs, and automating preservation actions.

- Hardware (rent or roll your own).

    - X_as_a_Service

    - SaaS, Paas, IaaS

## Infrastructure
(as a Service)

Managed by **CUSTOMER**:
- Applications
- Data
- Runtime
- Middleware
- Operating System

Managed by **VENDOR**:
- Virtualization
- Servers
- Storage
- Networking

## Platform
(as a Service)

Managed by **CUSTOMER**:
- Applications
- Data

Managed by **VENDOR**:
- Runtime
- Middleware
- Operating System
- Virtualization
- Servers
- Storage
- Networking

## Software
(as a Service)

Managed by **VENDOR**:
- Applications
- Data
- Runtime
- Middleware
- Operating System
- Virtualization
- Servers
- Storage
- Networking

amazon web services

heroku

Google Apps

# QDR<sup>BETA</sup>

QUALITATIVE DATA
REPOSITORY

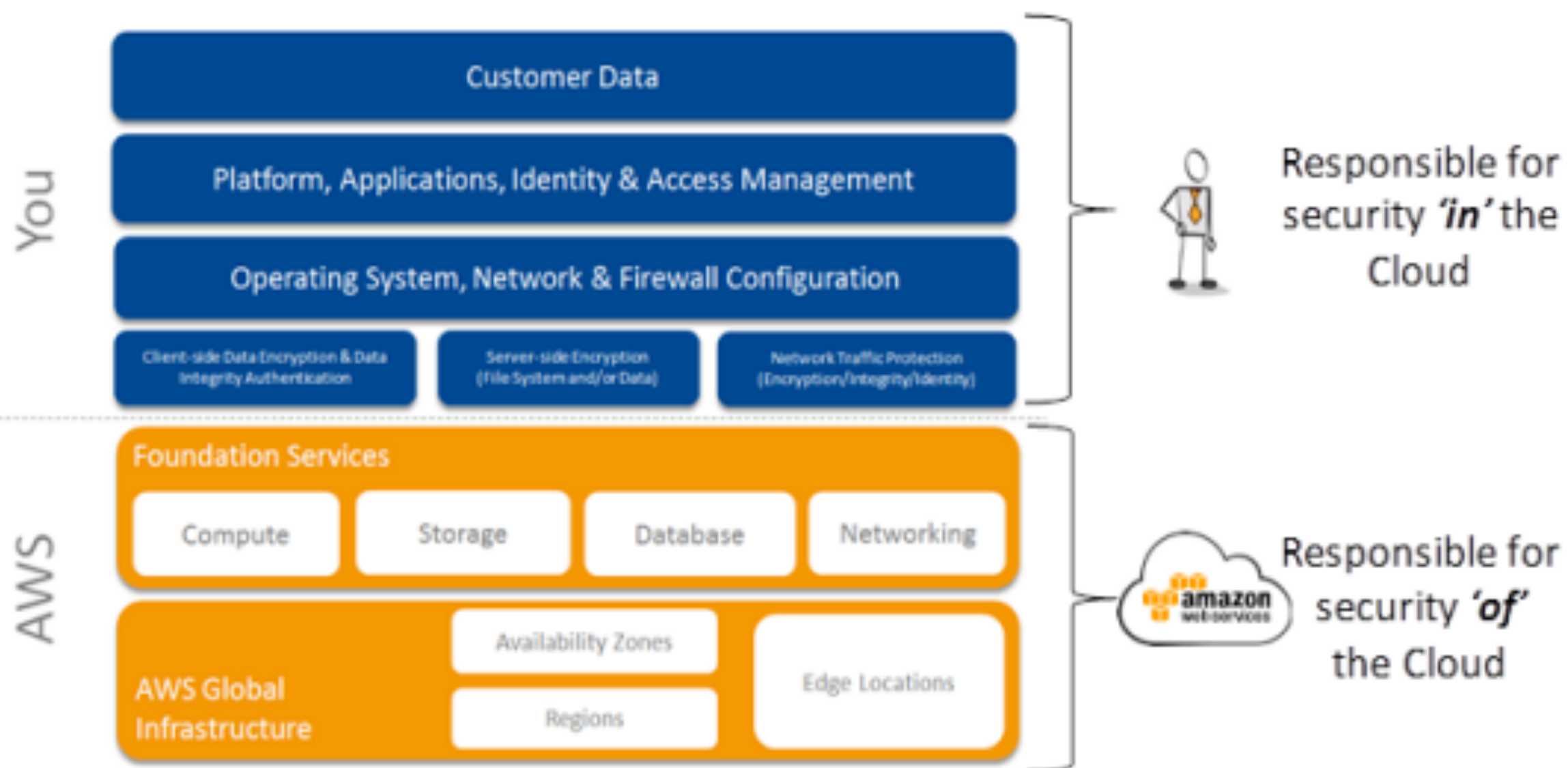| Amazon EC2 Service (US-East) | | | $ | 2256.56 |
|---|---|---|---|---|
| Compute: | $ | 0.00 | | |
| Intra-Region Data Transfer: | $ | 0.15 | | |
| EBS Volumes: | $ | 7.00 | | |
| Reserved Instances (One-time Fee): | $ | 2227.00 | | |
| Elastic IPs: | $ | 3.66 | | |
| Elastic LBs: | $ | 18.30 | | |
| Data Processed by Elastic LBs: | $ | 0.40 | | |
| VPC Peering Data Transfer: | $ | 0.05 | | |
| AWS Data Transfer In | | | $ | 0.00 |
| US-East / US Standard (Virginia) Region: | $ | 0.00 | | |
| AWS Data Transfer Out | | | $ | 6.66 |
| AWS Support (Business) | | | $ | 224.21 |
| AWS Support Plan Minimum: | $ | 100.00 | | |
| Support for Reserved Instances (One-time Fee): | $ | 124.21 | | |
| Free Tier Discount: | | | $ | -21.18 |
| Total One-Time Payment: | | | $ | 2351.21 |
| Total Monthly Payment: | | | $ | 115.04 |

Total = $338.00 Month average
Total infrastructure costs are ~ $10k per year.

Governance

- Getting legislation or institutional policy written and approved

    - Maintaining and enforcing the policy or legal requirements

    - Legal costs to comply with open data legislation

- Managing requests or questions about datasets

- Protecting privacy / providing security

# Sustainability

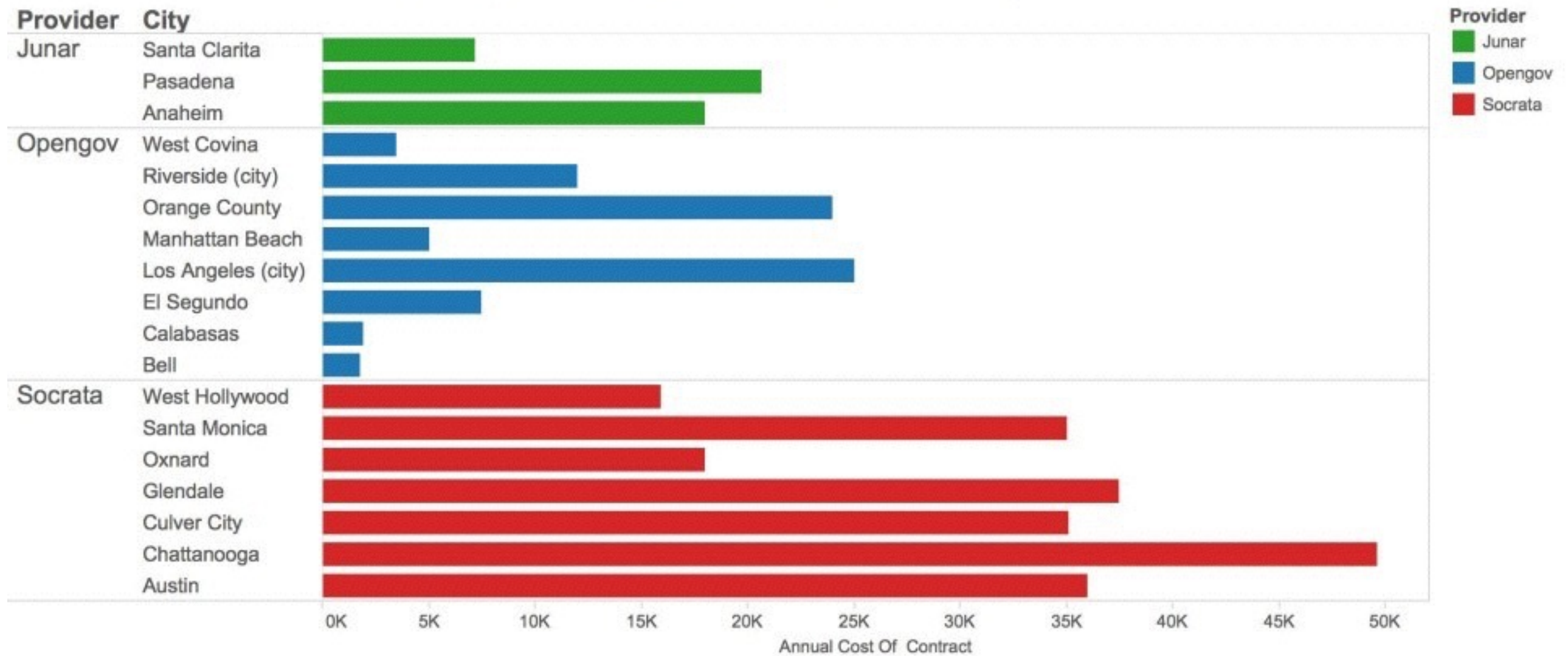- Hardware and software maintenance cost

- Data inventorying, publishing, optimizing discovery, updating relevant standards, preservation checks / actions

- Labor for doing both of the above tasks

- Liability costs in case of publication of nonpublic information (i.e. have a lawyer on retainer)

- Impact and analytics

- Vendor contracts (EzID; Socrate; perma.io; etc)

## Annual Cost of Contract OD Portal for Cities

| Provider | City | Annual Cost Of Contract |
|----------|------|------|
| Junar | Santa Clarita | ~7K |
| | Pasadena | ~21K |
| | Anaheim | ~18K |
| Opengov | West Covina | ~3.5K |
| | Riverside (city) | ~12K |
| | Orange County | ~24K |
| | Manhattan Beach | ~5K |
| | Los Angeles (city) | ~25K |
| | El Segundo | ~7.5K |
| | Calabasas | ~2K |
| | Bell | ~1.5K |
| Socrata | West Hollywood | ~16K |
| | Santa Monica | ~35K |
| | Oxnard | ~18K |
| | Glendale | ~37K |
| | Culver City | ~35K |
| | Chattanooga | ~49K |
| | Austin | ~36K |

**Provider**
- Junar (green)
- Opengov (blue)
- Socrata (red)

https://govex.jhu.edu/open-data-how-much-does-it-cost/

# Staff & Skills Development

- Developers (front and back end), and curators

- Awareness raising activities to promote use

- Capacity building for the use of data within government, discipline, or school

- Ongoing training and capacity building

# QDR

- Directors

  - Content - full time (salary)

  - Technical - .05 time (salary)

- Curators - 3 part-time (GA student salary)

- Developers

  - Front-end - Contract - $75/ hr

  - System Admin / Development / Operations - $70/hr

Labor costs > $110,000 / annual.

UKDA
- employs 64.5 people.
- total budget UK Data Archive (2010-11) ~ £3.43 million
    - £2.43 million - Staff
    - £1 million - Infrastructure


Dryad
- 4-6 FTE
- Total budget $350,000 per year
    - $300,000 - Staff
    - $5,000-$10,000 - Infrastructure

# Cost Recovery
## (academic, discipline or research repositories)

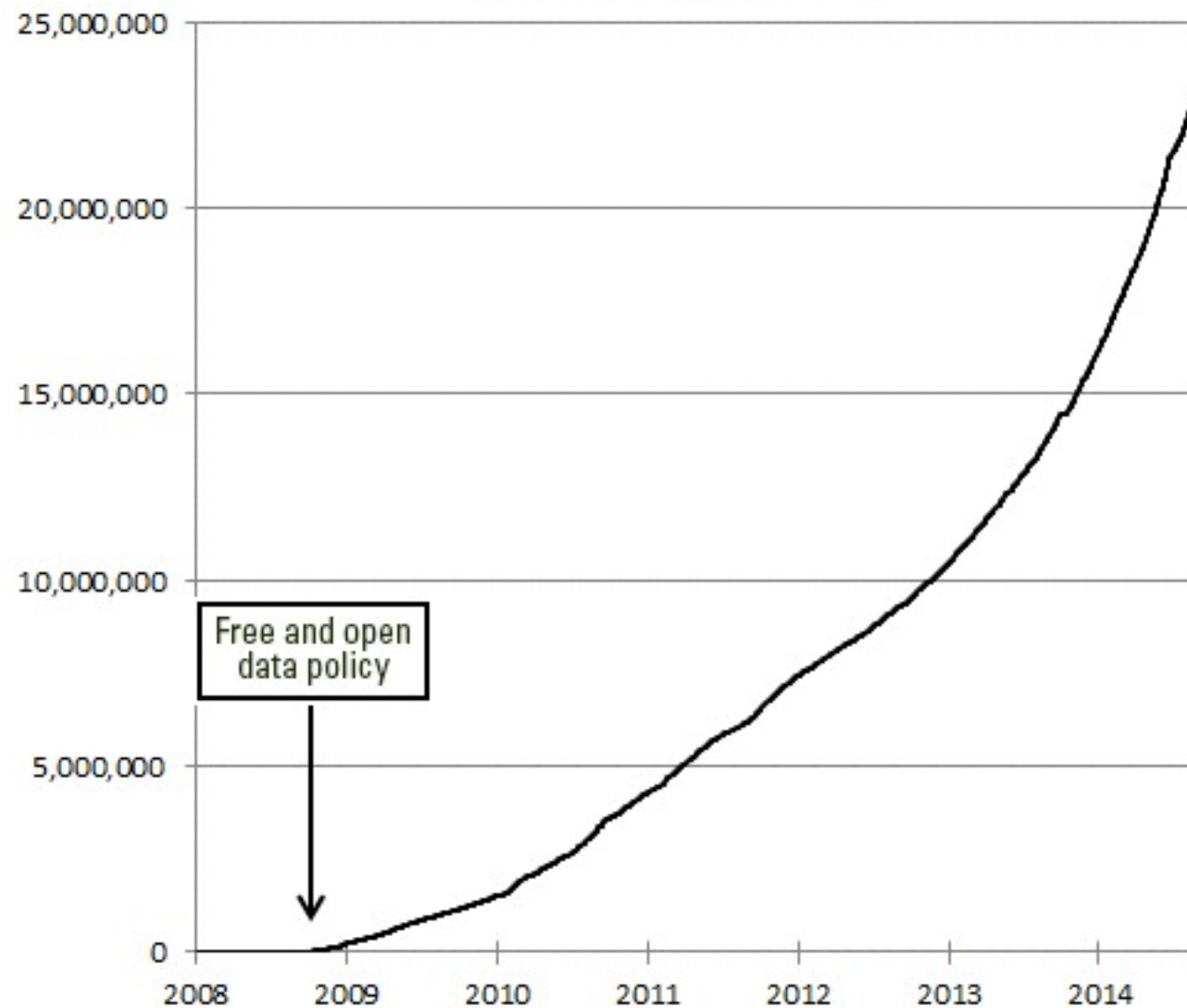1. **Membership**
2. **Submission fees**

Both of these create a barrier to entry, they differ on assumes the cost.

3. **Institutional Support**
4. **Federal funding for Special Projects**

| Funding Models | Potential for Economic Stability Needed for Long-Term Sustainability | Potential for Open Access to Research Data | Potential for Equity for Deposits by Individual Researchers | Potential for Equity for Universities/ Institutions |
|---|---|---|---|---|
| **Membership Dues** | Moderate; subject to institutional budgets and priorities | Low | Moderate | Low |
| **Submission Fees** | Low to Moderate; subject to policies of funding agencies and publications; | High | Low; costs transferred from end users to data producers | Low |
| **Institutional support** | Moderate; subject to institutional budgets and priorities | High | Low | Low |
| **Federally-sponsored Special Projects** | Low; subject to changes in national research priori- | High | Limited to designated research | High |

http://datacommunity.icpsr.umich.edu/sites/default/files/WhitePaper_ICPSR_SDRDD_121113.pdf

## Landsat Scenes Downloaded from USGS EROS Center (Cumulative)



Free and open data policy

## Estimated Productivity Savings from Uses of Landsat

**Landsat Applications** and their **Estimated Annual Efficiency Savings**
1. USDA Risk Management Agency → over $100 million
2. U.S. Government Mapping → over $100 million
3. Monitoring Consumptive Agricultural Water Use → $20–$80 million
4. Monitoring Global Security → $70 million
5. Landsat Support for Fire Management → $28–$30 million
6. Forest Fragmentation Detection → over $5 million
7. Forest Change Detection → over $5 million
8. World Agriculture Supply and Demand Estimates → over $3–$5 million
9. Vineyard Management and Water Conservation → $3-$5 million/year
10. Flood Mitigation Mapping → over $4.5 million
11. National Agricultural Commodities Mapping → $1.9 million/year
12. Waterfowl Habitat Mapping and Monitoring → $1.9 million/year
13. Coastal Change Analysis Program → $1.5 million
14. Forest Health Monitoring → $1.9 million/year
15. NGA Global Shoreline → over $90 million (one time)
16. Wildfire Risk Assessment → $25-$50 million (one time)

# $436 million per year

https://landsat.gsfc.nasa.gov/landsat-seen-as-stunning-return-on-public-investment/

## 6. Payment

### 6.1 Charges at Acceptance

Dryad does not charge any fees for Submissions that are not Accepted by the Repository. Dryad only Accepts Submissions that meet the Content criteria described in Section 2.1, including being in press or previously published by a scholarly publisher.

### 6.2 Hierarchy of Payments

In determining the party responsible for the Data Publishing Charges, Dryad will apply the following rules in sequence:

- If a Submitter qualifies for a Waiver, no Data Publishing Charge will be incurred.
- If Submission is covered by a Subscription Plan, the Data Publishing Charge will be covered by the Subscription Plan.
- If Submission is covered by a Voucher Account or a Deferred Payment Plan, the Voucher Account or Deferred Payment Plan will be charged.
- If Submitter has a Single Use Voucher, the Voucher will be redeemed.
- If none of the above applies, Submitter will be charged.
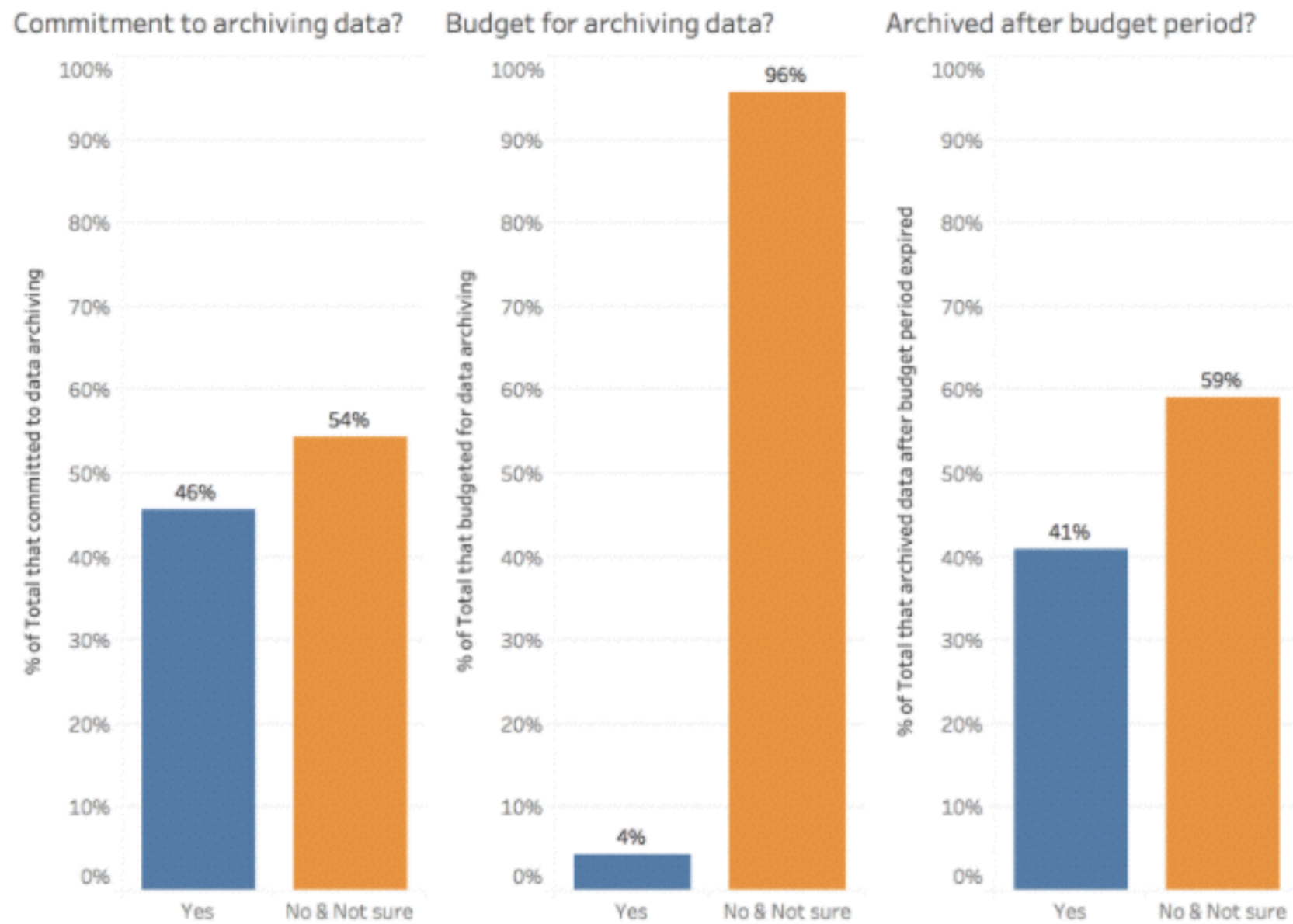
### 6.3 Additional charges

In all instances (other than where a Waiver is granted), Submitters will be charged excess data fees for Data Packages greater than 20 Gigabytes. Submitters also are responsible for paying any third-party costs associated with Submission of the Content (e.g. fee to use a large file transfer service external to Dryad).

### 6.4 Refunds

Data Publishing Charges and excess storage fees are not refundable. Please see Purchase Agreements for information regarding refunds for Payment Plans.

(back to top)

# $120 per archived dataset

- Back of the envelope math:

  - 2016 budget = $350,000

  - Cost for data deposit $120

  - Total data deposited 4258

    - Minus waiver rate of 1/3 of all submissions

  - Total 2016 revenue = $357,720. ($7000 margin!)

| | Self-Service $0 | Guided Service $3,000 + per dataset fee* | Lifecycle Service $5,000 + per dataset fee* |
|---|---|---|---|
| UNC Dataverse tool access | X | X | X |
| Data citation generation | X | X | X |
| Persistent identification (DOI) | X | X | X |
| Basic utilization reporting | X | X | X |
| Long-term preservation | X | X | X |
| Standardized metadata | X | X | X |
| User support | limited | standard | dedicated |
| Introductory Dataverse software training | | X | X |
| Dataset collection arrangement | | X | X |
| Metadata template development | | X | X |
| Data Management Plan implementation | | | X |
| File format normalization | | | X |
| Data file review | | | X |
| Access policy enforcement | | | X |
| Education and training program development | | | X |