# Serveis de dades i computació per al suport del cicle de recerca a una universitat

**Cicle de Conferències sobre la gestió de dades de recerca**
**Consorci de Serveis Universitaris de Catalunya**
**Novembre 18, 2020**

**Mercè Crosas, Ph.D., Harvard University**
**scholar.harvard.edu/mercecrosas  @mercecrosas**

# This Talk

**Research Data and Computing Services:**

- A recent US Report

- Our efforts at Harvard

- Examples

- What's next

# A Recent US Report

**RESEARCH REPORT**

November 18, 2020

# Research Data Services in US Higher Education

Jane Radecki, Rebecca Springer

DOI: https://doi.org/10.18665/sr.314397

**Topics:** Digital scholarship and data management, Libraries, Research practices

**Tags:** Data services

Jane Radecki

Rebecca Springer

# Study Design and Definitions

- Review 120 US universities
- Three groups:
  - **R1:** doctoral universities: very high research activity
  - **R2:** doctoral universities: high research activity
  - **SLACs:** Baccalaureate colleges

- Consider research data and computing services from:
  - **Libraries**
  - **IT/Research Computing**
  - **Research center and facilities**
  - **Professional Schools** (e.g., Medical School, Business School)

"we defined research data services as any concrete, programmatic offering intended to support researchers in working with data."

# Key Findings

- **Libraries** are important providers of research data services

- **IT/research computing** provide fewer research data services than libraries, but are an important provider

- A wide variety of services are provided by **academic departments, research centers and facilities, and professional schools**

- **High performance computing** offered: 100% R1, 60% R2, 24% SLACs
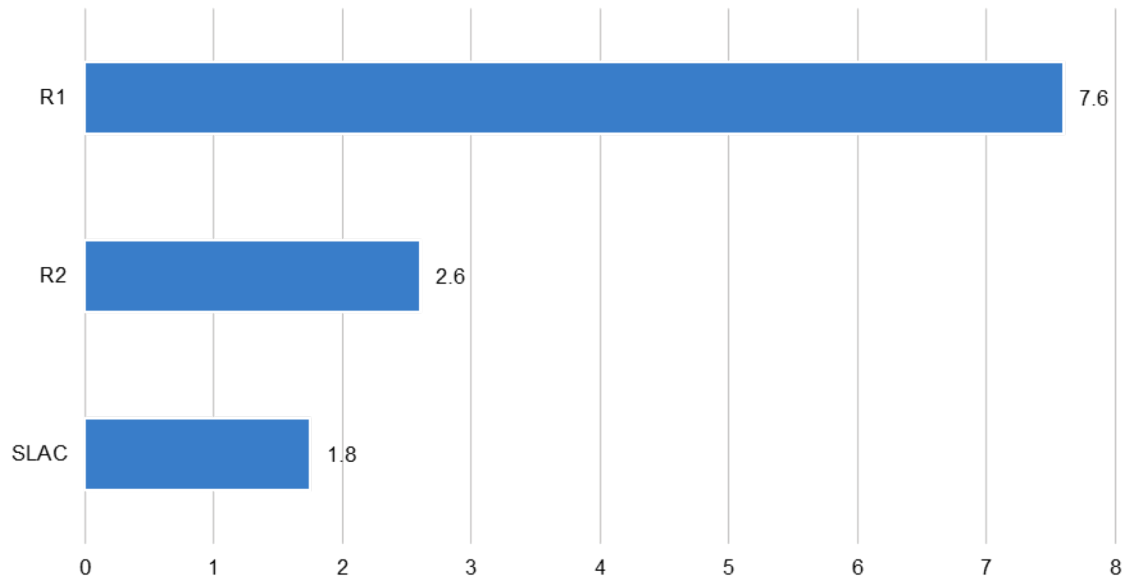
# Types of Research Data Services

## Within Libraries and IT

- Consulting
- Training events
- Backend work (data architecture, metadata design)
- Front end work (web development, data visualizations)
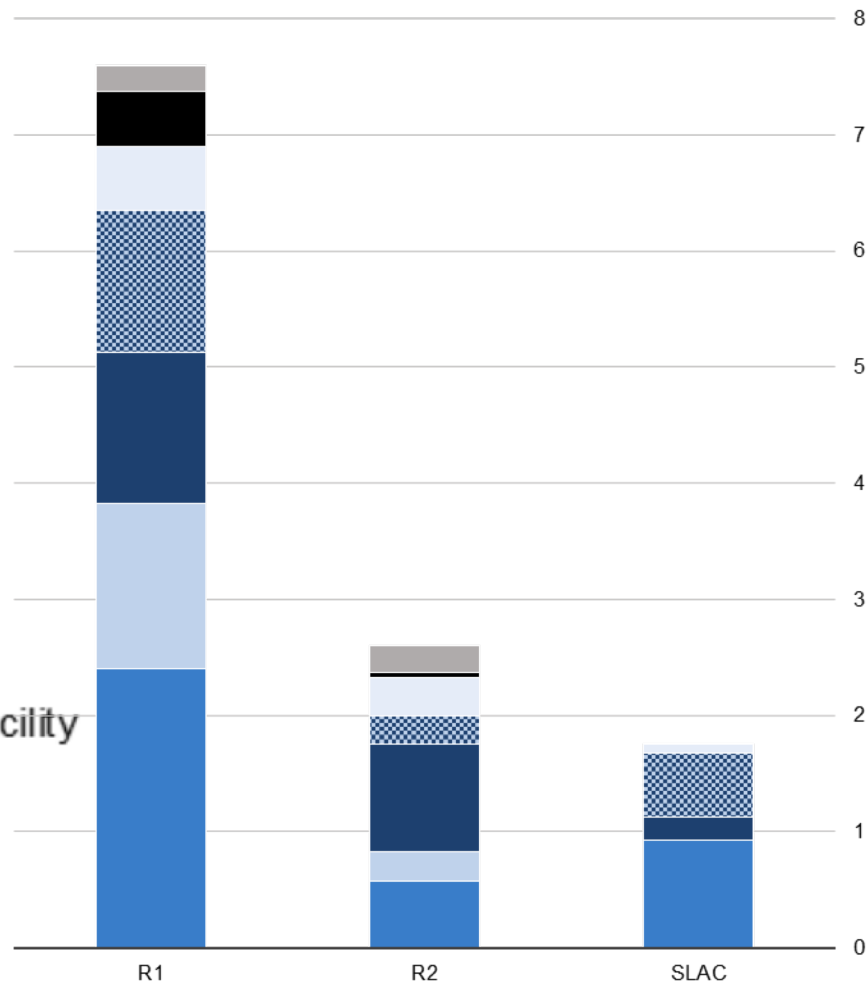
## Outside Libraries and IT

- Statistics
- Bioinformatics
- Geospatial
- Clinical data
- Business
- Social Science
- Visualizations

# R1s provide 2.5 times more services than R2s

# Libraries are the largest contributors to Research Data Services:
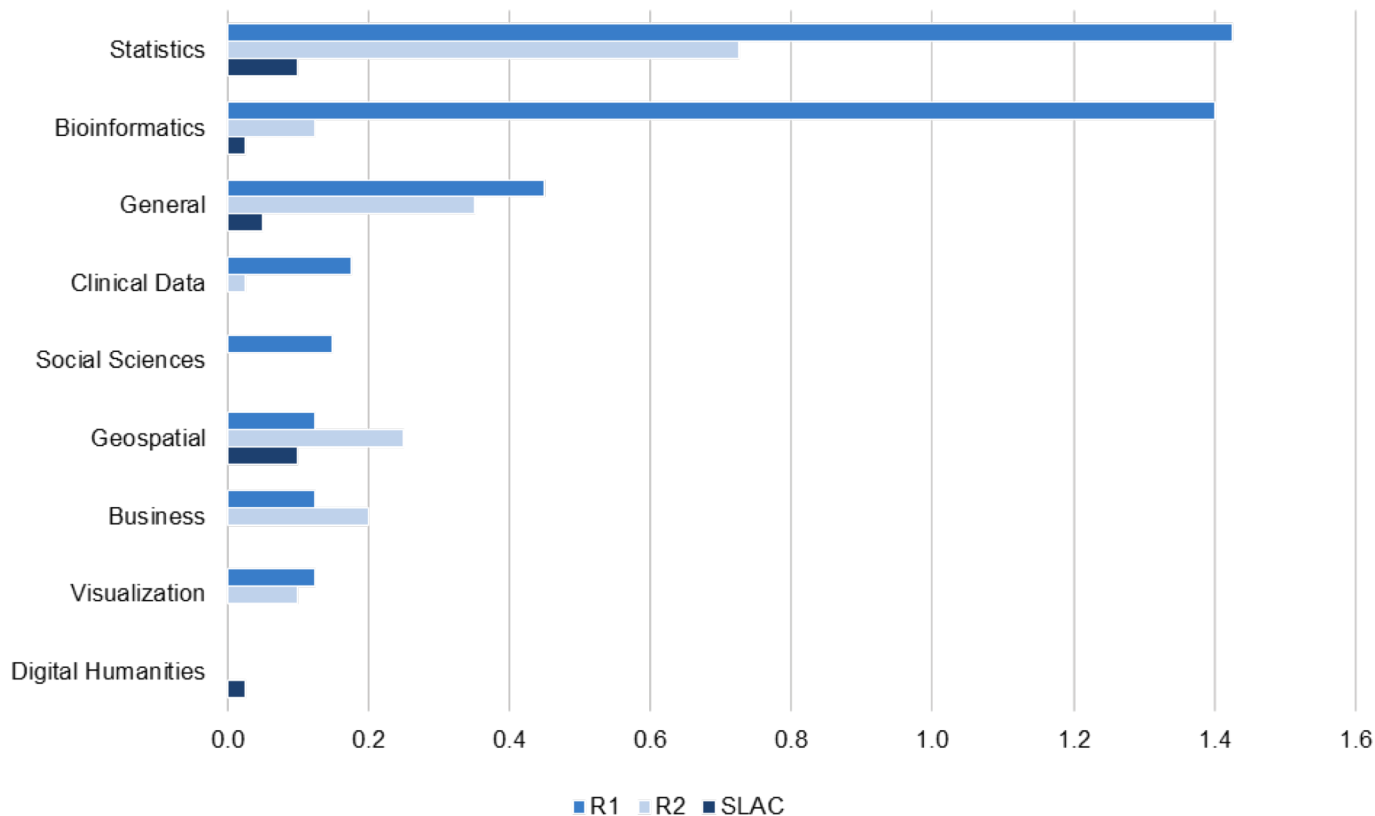
- 32% in R1s
- 53% in SLACs
- Only 22% in R2s



Legend:
- Library
- Medical School
- Independent Research Center or Facility
- IT Department
- Academic Department
- Other Professional School
- Business School

# Profile of types of Library Data Services

**Generalist consultation** is the most common service offered by the libraries

| | Consulting | Training Events | Front End Work | Back End Work | Total |
|---|---|---|---|---|---|
| **General** | 35.9% | 16.0% | 3.2% | 2.6% | **57.7%** |
| **Geospatial** | 16.7% | 9.0% | 0.0% | 0.0% | **25.6%** |
| **Statistics** | 7.1% | 1.3% | 0.0% | 0.0% | **8.3%** |
| **Digital Humanities** | 2.6% | 1.3% | 0.0% | 0.0% | **3.8%** |
| **Social Sciences** | 0.6% | 0.6% | 0.0% | 0.0% | **1.3%** |
| **Health Sciences** | 0.6% | 0.0% | 0.0% | 0.0% | **0.6%** |
| **Other** | 1.3% | 1.3% | 0.0% | 0.0% | **2.6%** |
| **Total** | **64.7%** | **29.5%** | **3.2%** | **2.6%** | **100%** |

# Average number of research data services per institution offered by centers and facilities, departments, and schools
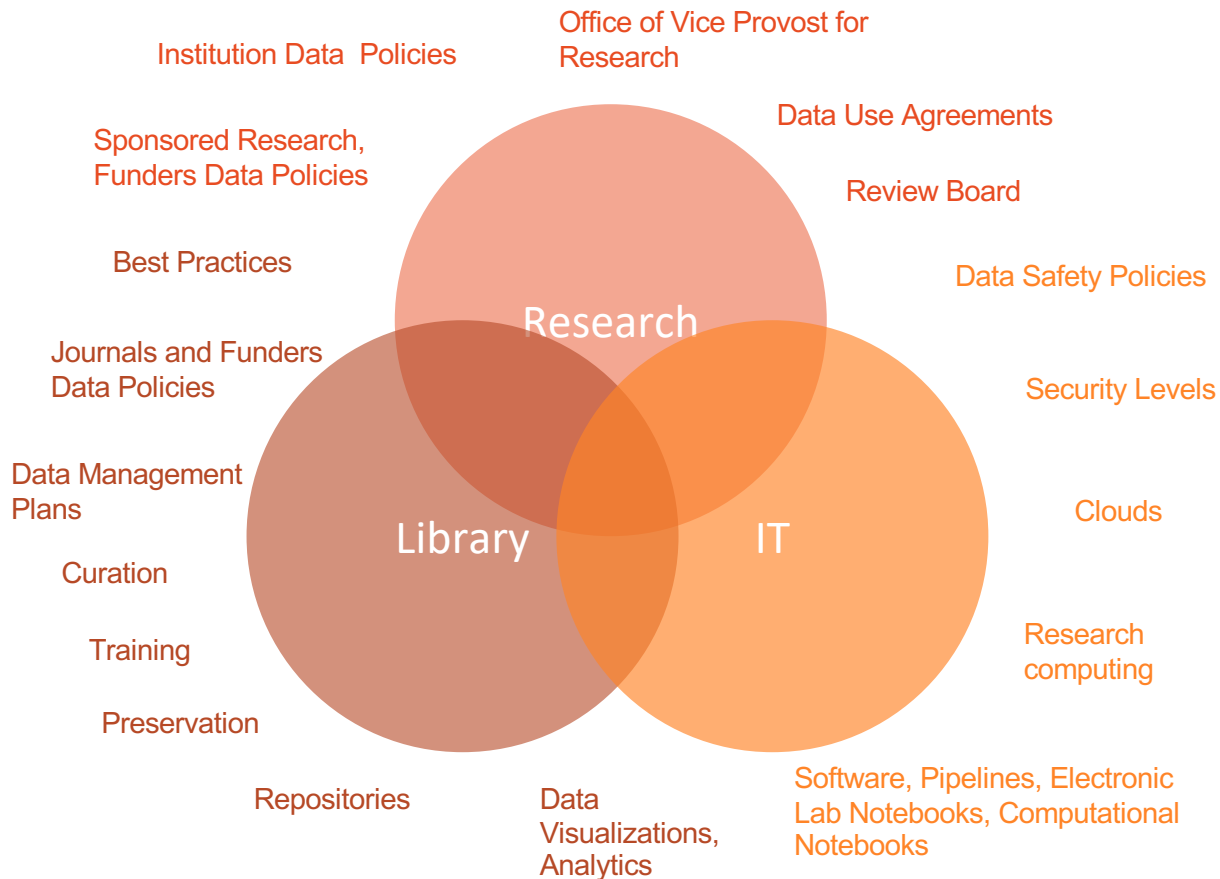
# Report Conclusions

- Research Data Services in US universities are **decentralized** and **uneven**

- **Data-driven research** is increasing, but services not funded at same pace

- Should **duplicated** services merge for efficiency?

- Should **collaborative models** across universities be used to share expertise?

**"As centralized points of contact on campus, libraries, IT departments, and research offices may be particularly well positioned to act as dispatchers, connecting scholars to the services that best meet their needs"**

# Our efforts at Harvard

A **wide variety** of research services across Harvard.

**Collaboration** between Research offices, the Library, and IT is key.

Institution Data Policies

Office of Vice Provost for Research

Data Use Agreements

Sponsored Research, Funders Data Policies

Review Board

Best Practices

Data Safety Policies

Research

Security Levels

Journals and Funders Data Policies

Data Management Plans

Library

IT

Clouds

Curation

Training

Research computing

Preservation

Repositories

Data Visualizations, Analytics

Software, Pipelines, Electronic Lab Notebooks, Computational Notebooks

# A common project

# A single resource to find all research services

- Collaborate in a project with a **common vision**

- Build a **research support website**:

    - *"To help faculty, researchers, and those who work with them to advance their research by easily finding and browsing the University's breadth of resources and services"*

- Sponsored by Harvard Library, IT, and Office of vice provost for research

- Initial launch planned for: Early 2021

# The method

# Build inventories of existing services

- **Collect data:** Phase I, 2018-2020

  - Harvard Library (11 libraries/units)

  - Research Computing (4 schools/units)

  - Research Administration and Compliance (2 units)

  - Ethics Board (IRB), Data Use Agreements, Data Safety (3 units)

- **Create a catalog** of a total of 36 service offerings

- Focus on the **service function**, not on the service provider

# Standardize services information across units

- Classify and describe the services offerings in a **unified and uniform** way

- Three main **Services**:

  - **Research Data and Scholarship Services:** 25 service offerings

  - **Research Computing:** 6 service offerings

  - **Research Administration and Compliance:** 5 service offerings

- Three phases of the **Research Lifecycle**:

  - **Planning:** 11 service offerings

  - **Active Research:** 19 service offerings

  - **Dissemination and Preservation:** 6 service offerings

# Coordinate through working groups

- In decentralized universities, **working groups** can help establish and **achieve unified goals** and **communicate** across schools and units

In the last couple of years, we created the following groups relevant to this project:

- Working group to help **coordinate research data management** efforts

- Working group for **building the research support website**

# The result (so far)

To be launched in January 2021

SUGGEST A RESOURCE

# Research Support at Harvard

Services ▾    Lifecycle ▾    Research Remotely    About    Get Help

# Explore the breadth of the University's resources and services

# Browse service offerings by three main services



**Services**

## Welcome

Harvard researchers have access to a wide range of service offerings across the University. From planning a project or study, to sharing and archiving methods or findings, our services span the entire research lifecycle. This website brings together Harvard's offerings across central units and schools, including support for research administration and compliance, data management and scholarship, and research computing.

## Browse by Services

### Research Administration and Compliance →

Harvard offers a full spectrum of resources to support and facilitate research and researcher compliance with internal and external regulations and policies. Services available

### Research Computing →

Research Computing at Harvard facilitates the advancement of research by providing leading-edge computing services including cluster computing, storage, software licenses, virtual instances, and

### Research Data and Scholarship →

Researchers at Harvard generate data and scholarship that changes the world. Services across the University are available to support data creation, curation, and transformation, as well as research publishing and

# Uniformity across research support offerings
## Same fields for each service offerings: Audience, Provider, Fee, Website, Contact

# Browse services offerings by Research Lifecycle phases



**Planning:**
Access & Reuse
Plan & Design

**Active Research:**
Collect & Create
Analyze & Collaborate

**Dissemination & Preservation:**
Evaluate & Archive
Share & Disseminate

# Planning: Access & Reuse Plan & Design

**11 service offerings:**

- Buying and Licensing Data
- Data Retrieval, Finding Data
- Data Safety and Regulated Data
- Data Use Agreement Processing
- Human Subjects & Animal Research Resources
- Pre- & Post-Award Resources
- Research Data Management Lifecycle
- Research Design
- Training, Workshop, Capacity Building
- Project Health Informationist

# Active Research



The active research phase of a project may include collecting or acquiring data, information, or sources, conducting quantitative or qualitative analysis, and/or using computation resources, data storage, quantitative or qualitative tools, visualizations, or information exploration.

**RESEARCH LIFECYCLE**

Planning

Active Research

Dissemination & Preservation

**Cluster Computing →**

Doing computations at scale allows a researcher to test many different variables at once, thereby shorter time to outcomes, and also provides the ability to ask larger, more complex problems (i.e., larger data sets, longer simulation time, more degrees of freedom, etc.). Researchers can take advantage of the scale of the cluster by setting up workflows to split many different tasks into large batches, which are scheduled across the cluster at the same time....

**Data Cleaning →**

Data Cleaning services and consultation support for cleaning, reformatting, merging, and scraping data for analyzing, visualization and reporting...

**Data Curation →**

Specialists throughout Harvard Library are available to consult about data curation, organization, and integration....

**Data Handling →**

Consultation, instruction, and support for practices and procedures involving data (e.g. reformatting)...

**Data Science and Research Software Engineering Collaboration →**

Data Science and Software Engineering play an important role in research by creating new capabilities to process and analyze data, helping ensure reproducibility, and aiding researchers in extracting knowledge and insight for the data. The

**Data Security Support →**

Consultations and/or instruction on ensuring data security during the research lifecycle, including compliance with University policies...

---

## Active Research:
## Collect & Create
## Analyze & Collaborate

**19 service offerings**

- Cluster Computing, Virtual Instances
- Research Data Storage, Database, Data Security
- Software and Platforms
- Research Computing Consulting & Facilitation
- Data Science and Research Software Engineering, Statistical Analysis, Text Analysis
- Dataset Creation, Data Cleaning, Data Curation, Data Handling, Metadata creation
- Data Visualization
- Geospatial data
- Qualitative Data Support
- Microbiological Safety

# Dissemination & Preservation:
## Evaluate & Archive
## Share & Disseminate

**6 service offerings:**

- Copyright and Intellectual Property

- Archiving data

- Data Deposit

- Data Sharing and Publishing

- Harvard Dataverse Curation

- Harvard Dataverse Repository

| Services | Planning | Active Research | Dissemination & Preservation |
|---|---|---|---|
| **Research Administration & Compliance** | • Data Safety and Regulated Data<br>• Data Use Agreement Processing<br>• Human Subjects & Animal Research Resources<br>• Pre- & Post-Award Resources | • Microbiological Safety | |
| **Research Computing** | | • Cluster Computing, Virtual Instances<br>• Research Data Storage, Database, Data Security<br>• Research Computing Consulting & Facilitation<br>• Data Science and Research Software Engineering | |
| **Research Data & Scholarship** | • Buying and Licensing Data<br>• Data Retrieval, Finding Data<br>• Research Data Management Lifecycle<br>• Research Design<br>• Training, Workshop, Capacity Building<br>• Project Health Informationist | • Statistical Analysis, Text Analysis<br>• Dataset Creation, Data Cleaning, Data Curation, Data Handling, Metadata creation<br>• Data Visualization<br>• Geospatial data<br>• Qualitative Data Support<br>• Software and Platforms | • Copyright and Intellectual Property<br>• Archiving data<br>• Data Deposit<br>• Data Sharing and Publishing<br>• Harvard Dataverse Curation<br>• Harvard Dataverse Repository |

# Examples

Data Use Agreement Processing

Data Science and Research Software Engineering

Harvard Dataverse Curation

**SERVICES**

▾ Research Administration & Compliance

Data Safety & Regulated Data

Data Use Agreement Processing

eSupport - Committee on Microbiological Safety (eCOMS)

Human Subjects and Animal Research Resources

Pre- & Post-Award Resources

▸ Research Computing

▸ Research Data and Scholarship

# Data Use Agreement Processing

The transfer of data between organizations is common in the research community. When the data is confidential, proprietary, or otherwise considered sensitive or protected, the organization providing the data, whether that is Harvard or a third party, may require a Data Use Agreement (DUA) to govern the exchange of data. This section includes guidance on the DUA-Agreements Application which supports the review, approval and management process for DUAs, and other related resources.

Data Use Agreement review and compliance application for researchers at Harvard interested in requesting data from a third party, or providing data to a third party. This application supports the review, approval and management process for Data Use Agreements.

A DUA is a binding contract governing access to and treatment of nonpublic data provided by one party (a "Provider") to another party (a "Recipient"). DUAs are often required by external parties, and may also be necessary for Harvard data to be disclosed to another organization. DUA terms and conditions vary depending on the laws and regulations governing the specific type of data to be shared, as well as the policies and/or requirements of the Provider and Recipient. If you are unsure whether a DUA is necessary, feel free to reach out to your sponsored research office.

## Details by Provider

⊕ HUIT Administrative Technology Services, Research Administration and Compliance

⊕ Harvard Medical and Dental Schools

⊕ Harvard T. H. Chan School of Public Health

⊕ Harvard University Area

**A Research Administration and Compliance service offering in the Planning phase:**

**Data Use Agreement Processing**

A Research Computing service offering in the Active Research phase:

Data Science and Research Software Engineering

A Research Data service offering in the Dissemination and Preservation phase:

Harvard Dataverse Curation

Harvard researcher requests a Dataset from a third party

↔ 1. **Review** by Harvard DUA team (iterate with researcher and third party)

→ 2. **Approve** DUA. Dataset can be used for research

→ 3. **Manage** DUA for compliance

Does it needs a Data Use Agreement (DUA)?

Yes. Enter info into System →

Harvard DUA System

Researcher obtains dataset

DUA

# Metrics on Data Use Agreement Processing

- DUA system launched in 2018

- 879 active DUAs

- Average time to review a DUA:

  - 85 days in 2019

  - 65 days in 2020

- Top 5 departments with DUAs:
  Health Care, Education,
  Bioinformatics, Epidemiology,
  Economics

| Top 20 Departments with Active Agreements | # of DUAs |
|---|---|
| Health Care Policy | 214 |
| Center for Education Policy Research [CEPR] | 91 |
| Biomedical Informatics | 63 |
| Epidemiology | 57 |
| Economics | 34 |
| Environmental Health | 26 |
| Nutrition | 24 |
| Global Health and Population | 23 |
| Education Policy Research, Center for | 22 |
| Psychology | 21 |
| Genetics | 16 |
| Social and Behavioral Sciences | 13 |
| Accounting and Management | 12 |
| Biostatistics | 12 |
| Health Policy and Management | 10 |
| Government | 9 |
| Joint Center for Housing Studies [GSD] | 9 |
| Other[HLS] | 9 |
| Finance | 9 |
| Other[GSE] | 9 |

# Data Science and Research Software Engineering Collaboration

Data Science and Software Engineering play an important role in research by creating new capabilities to process and analyze data, helping ensure reproducibility, and aiding researchers in extracting knowledge and insight for the data.  The term software here is used broadly to include all the ways in which one creates and analyses data. Researchers utilize software in their research by using scripts, tools, open-source software, and licensed software.  Data science also covers a wide range of skills and techniques applied to cleaning (aka wrangling), processing, and statistics that are typically beyond what a researcher from a specific domain might have. Due to the rapidly evolving nature of research, there are not always codes for all functions needed, nor are their clean data sources; therefore, the software or data pipelines are developed specifically for a given project. Traditionally, this development was done with researchers (graduate students and postdocs) or independent contractors. This approach poses several issues in terms of maintenance, optimization, reproducibility, and cost.  RSE or Data Scientist team can work closely with other Research Computing Systems teams to design, develop, deploy, optimize, and maintain software packages/tools and data pipelines that are paired with specific hardware architectures to accelerate cutting-edge research at Harvard University.

## Details by Provider

⊕ Faculty of Arts and Sciences, Research Computing

⊕ Institute for Quantitative Social Sciences

⊕ Harvard Business School

A Research Administration and Compliance service offering in the Planning phase:

Data Use Agreement Processing

**A Research Computing service offering in the Active Research phase:**

**Data Science and Research Software Engineering**

A Research Data service offering din the Dissemination and Preservation phase:

Harvard Dataverse Curation

# Data Science Services

- Offered by the Institute for Quantitative Social Science (IQSS)

- Focuses on **social science** support, but includes other scientific domains

- **Consulting:** short term

- **Collaboration:** longer project (fee)

- **Workshops** (in collaboration with Harvard Business School):
  - **Python:** Introduction, web scraping
  - **R:** Introduction, regressions models, graphics, data wrangling
  - **Stata:** Introduction, data management, regression models, graphics
  - **Other:** Introduction to programming, SAS introduction, data science tools

# More than 20 workshops per year attended by 700 scholars

# Consultations and requests

- Majority of consultations on statistics support

- More common languages are R and Stata

- Questions about web technology, text scraping, and text mining are fairly frequent



Number of Help Requests by Type (2012 - 2020)

# Harvard Dataverse Curation

The Harvard Dataverse data curation team, staffed by member of IQSS and the Harvard Library (and separately, the Harvard Kennedy School Library), provides fee-based curation services and free consultations to researchers around the world who are depositing data into the Harvard Dataverse.

Research data replication datasets, data for related publications, and all file types and domains are welcomed in the Harvard Dataverse. Through this engagement, the curation services team will ensure that deposited datasets are discoverable, accessible, interoperable, and reusable (FAIR). (IQSS)

## Details by Provider

**⊖ Institute for Quantitative Social Sciences**

**Audience**
- All Affiliates
- All Faculty
- All Graduate Students
- All Undergraduate Students
- Public

**Service Provider**
Institute for Quantitative Social Sciences (IQSS)

**Service Fee**
Yes

**Service Website**
https://support.dataverse.harvard.edu/curation-services

**Contact Information**
support@dataverse.harvard.edu

---

**SERVICES**

- ▸ Research Administration & Compliance
- ▸ Research Computing
- ▾ Research Data and Scholarship
  - Archiving Faculty Research Data and Archiving Data
  - Buying and Licensing Data
  - Copyright and Intellectual Property
  - Data Cleaning
  - Data Curation
  - Data Deposit
  - Data Handling
  - Data Retrieval
  - Data Security Support
  - Data Sharing and Publishing
  - Data Visualization
  - Dataset Creation
  - Finding Data
  - Geospatial Library, Data Analysis, Creation, Visualization

---

A Research Administration and Compliance service offering in the Planning phase:

Data Use Agreement Processing

A Research Computing service offering in the Active Research phase:

Data Science and Research Software Engineering

**A Research Data service offering in the Dissemination and Preservation phase:**
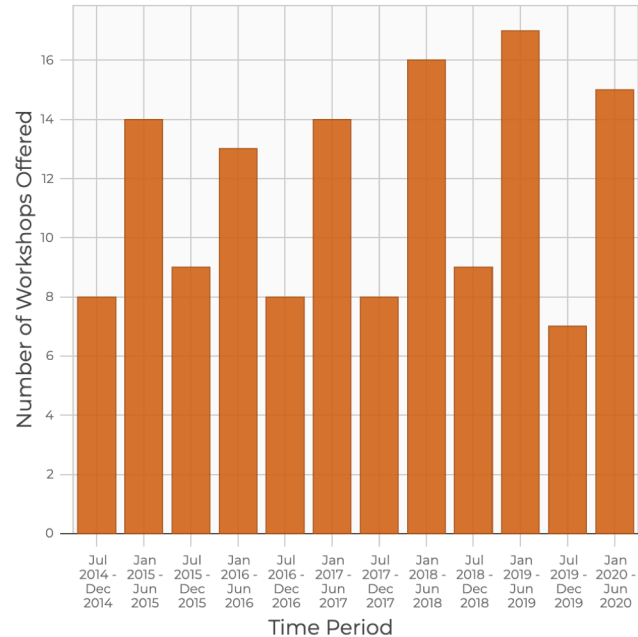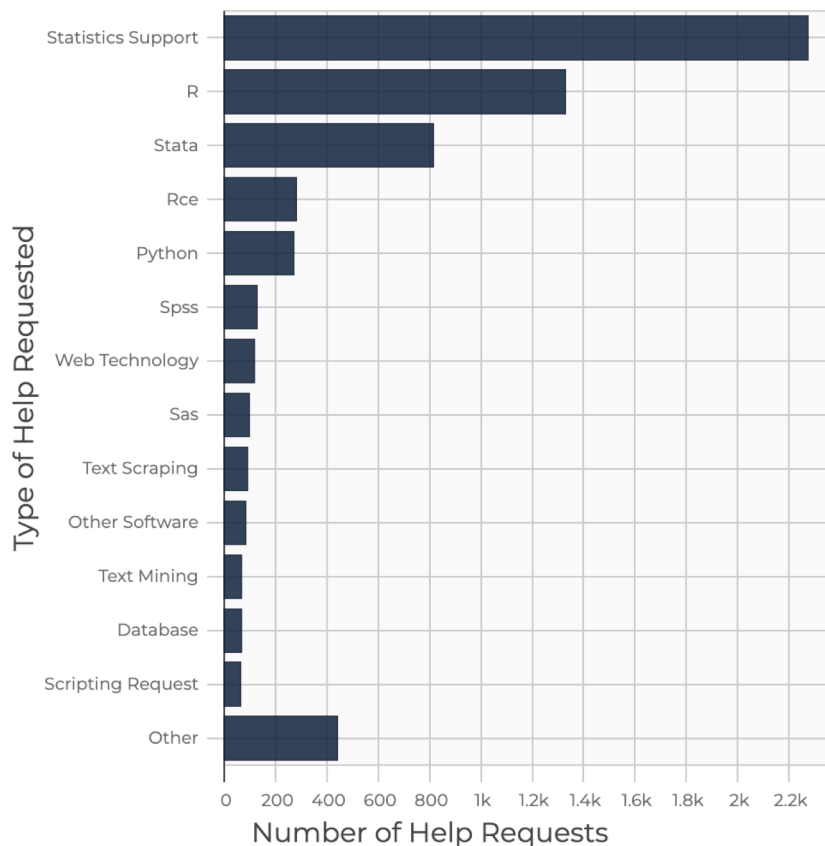
**Harvard Dataverse Curation**

# Harvard Dataverse Curation Services

- Launched in 2020
- A **collaboration** between IQSS and the Harvard Library
- **Surveyed** Harvard affiliates and Journals using Harvard Dataverse repository
- **Tiered service offerings:**
  - Free consultation (limited to 3 hours)
  - Extended consultation services
  - Dataverse collection set-up services and dataset and data file ingest
  - Ongoing Dataverse collection administration and curation services
  - Custom services for existing Dataverse collections
- In 2021, new service for supporting "managed collections" interested in receiving **Core Trust Seal** certification.

# What's Next

# What we are learning

- Increase in **data science** and **data-centric research** is transforming the way we need to provide research services to our universities

- Data science and data handling are becoming an integrated part of the **education**

- Research support services are often **distributed across schools, centers and facilities**, but they might be duplicative or depend on each other

- Some services benefit from being **centralized**, but others work better close to subject expertise

- **Collaboration and communication** are key and must be constant

# Towards an integrated solution

- Research data, computing, and compliance services should be **more integrated** to each other and to the research work

- We need **integrated technology and research tools** to support the services

- Whenever possible, we should **automate and streamline** the steps. *For example: machine-actionable Data Management Plans and Data Use Agreements; Electronic Lab Notebooks, Computational Notebooks, and Workflows integrated with repositories*

**A Data Commons can be part of the solution by providing the interoperability and tooling needed and connecting the services with the technology.**
*(see https://scholar.harvard.edu/mercecrosas/presentations/el-data-commons)*

# Thank you

**Mercè Crosas, Ph.D., Harvard University**

**scholar.harvard.edu/mercecrosas  @mercecrosas**