# Next Generation Machine Learning: The Evolution of the Library as Research Partner, Project Catalyst and Digital Integrator

Xuemao Wang
Vice Provost for Digital Scholarship
Dean and University Librarian

James Lee
Associate Vice Provost for Digital Scholarship
Associate Dean of Libraries
Director, Digital Scholarship Center

Kristen Burgess
Operational Manager
Digital Scholarship Center

# University of Cincinnati

**Founded in**
# 1819

**Enrollment**
## 46,798

**Degree Programs**
## 414

**Student:Teacher Ratio**
## 16:1

## University of Cincinnati Quick Facts

**Location:** Cincinnati, Ohio
**Number of Buildings:** 118 facilities on 476 acres
**Majors & Programs:** 414 degree programs, 262 minors and certificates
**Athletics:** NCAA Division I; American Athletic Conference
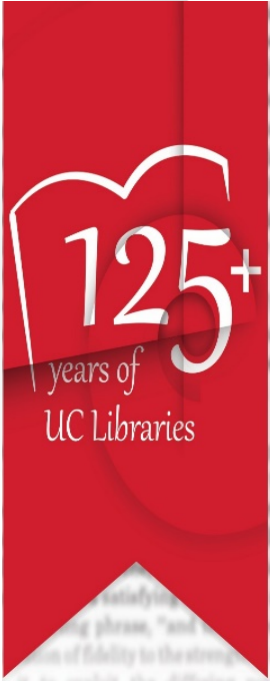**Mascot:** Bearcat
**Colors:** Red and Black
**Famous Alumni & Faculty:** Astronaut Neil Armstrong; President and later Chief Justice William Howard Taft; Eula Bingham, environmental scientist and one-time head of OSHA; Albert Sabin, developer of the oral polio vaccine; and prima ballerina Suzanne Farrell
**Students from:** 50 states and 114 countries
**Living Alumni:** over 300,000 with approximately half (more than 140,000 residing in the greater Cincinnati region).

Additional information available on the UC Fact Sheet.

# University of Cincinnati Libraries

**125+ years of UC Libraries**

*Archives & Rare Books Library*

*Chem-Bio Library*

*Classics Library*

*CCM Library*

*DAAP Library*

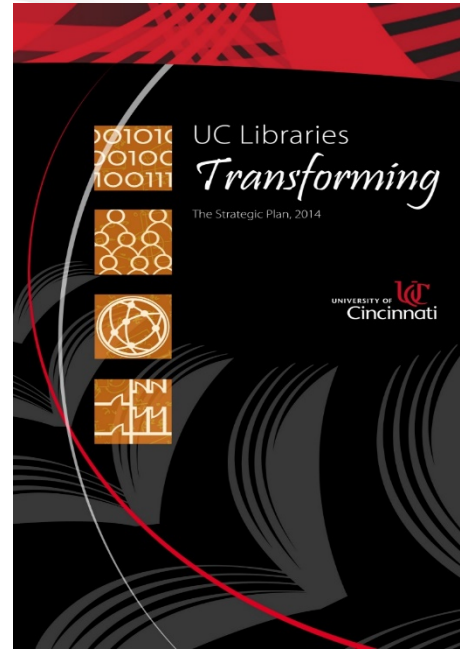*CECH Library*

*CEAS Library*

*Geo-Math-Phys Library*

*Health Sciences Library*

*Regional Libraries: Blue Ash, Clermont & Law*
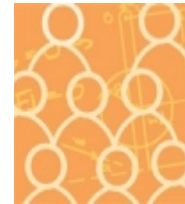
www.libraries.uc.edu

# UC Strategic Sizing Plan

University of CINCINNATI | LIBRARIES

## ENROLLMENT GROWTH 2018-2028

**Strategic Sizing**
Headcount Enrollment Projections Fall 2018 - Fall 2028

College Projection
IR Projection

45,949 | 46,388 | 46,798 | 50,602 | 52,367 | 53,953 | 55,602 | 57,316 | 59,101 | 60,957 | 62,890
48,000 | 50,686 | 52,122 | 53,486 | 54,740 | 55,918 | 56,916 | 58,001

2018 | 2019 | 2020 | 2021 | 2022 | 2023 | 2024 | 2025 | 2026 | 2027 | 2028

## GROWTH IN ACCESS & QUALITY
### TRENDS IN COLLEGE-LEVEL PROGRAM EFFORTS

COMPUTING & BIG DATA

HEALTH & WELLNESS

DIGITAL

INTERDISCIPLINARY

RETENTION

www.libraries.uc.edu

## GROWTH IN EDUCATION
### INSTITUTION-WIDE EFFORTS TO INCREASE ACCESS & QUALITY

BEARCAT PROMISE

CO-OP 2.0

CPS STRONG

LAUNCH UC

DIRECT ADMIT

NEXT/NOW SCHOLARSHIPS

VETERANS

ADULT LEARNING INSTITUTE

UC ONLINE

PROFESSIONAL MASTERS

PHD

URM PROGRAMS

RECRUITMENT, RETENTION, & COMPLETION

www.libraries.uc.edu

## GROWTH IN URBAN/COMMUNITY IMPACT
### TRENDS IN COLLEGE-LEVEL EFFORTS

CINCINNATI INNOVATION DISTRICT

CPS READINESS

SKILLING & CREDENTIALING

CLINICAL PARTNERSHIPS

# UC Research 2030 Plan

RESEARCH2030

**UC'S 10-YEAR STRATEGIC PLAN FOR RESEARCH**

The University of Cincinnati is the leading R1 urban university in our region with unrivaled talent solving problems that matter. We are rigorously pursuing diversity, equity and inclusion in research and actively transforming society through the creation of game-changing new knowledge and application of disruptive discoveries.

SEE THE PLAN

OBJECTIVE
## NATIONAL PROMINENCE

GOAL
**TOP 25 PUBLIC RESEARCH UNIVERSITY**

GUIDING PRINCIPLE
**Galvanizing our mission to serve the public good**

## INVEST IN SUCCESS

- **RECRUIT AND RETAIN TOP TALENT**
- **INNOVATE THE RESEARCH INFRASTRUCTURE**
- **CULTIVATE & GROW PROGRAMS OF EXCELLENCE**

## INVEST TO ADVANCE

- **URBAN FUTURES PATHWAY**
- **RETHINKING THE WHERE**
- **COALITION FOR CHANGE**

OBJECTIVE
## IMPACTFUL RESEARCH

GOAL
**IMPROVING PEOPLE'S LIVES**

GUIDING PRINCIPLE
**Foundational Partnerships to Solve Real World Problems**

www.libraries.uc.edu

# DSC: Mission

- Core mission - To break silos and cross wires across the university. We work at the intersection of data science, the arts and humanities, and the libraries.

- Academic Center (Libraries + Arts & Sciences)

- Mellon Digital Integrator (One of six, projects with eight colleges).

- UC Digital Futures Anchor Team

www.libraries.uc.edu

# DSC: What We Do

- We are a technical catalyst: technology to activate new research

- Machine Learning and Human-Interpretable Data Visualization on Large Unstructured Datasets (Text, Image, Sound, Video)

- We translate between disciplines that rarely interact in order to connect content experts with technical experts.

- We provide resources and infrastructure to nurture research questions and collaborations that slip between the cracks of colleges and funding agencies.

# UC's Digital Catalyst



**Digital Scholarship Center**

**Mission:** To create new areas of research and discovery across the university by applying AI + data visualization to human problems.

**Vision:** The DSC will be a distinctive global leader in transdisciplinary digital research.

**TEACHING**
- Intro coursework & workshops
- Student research in faculty team labs
- Job placement & digital training

**CORE SERVICES**
- Digital tools development & training
- Faculty development
- Digital skills for analog expertise & archives

**RESEARCH**
Networked Research: Leverage different funding models and different research outputs for common team goals.

CEAS
Allied Health Sciences
Social Science
Planning
Languages
DAAP
Architecture
CECH
Business
Digital Futures
CCM
CCHMC
Politics
Natural Science
Social Justice
Design
A&S
Law
Pharmacy
Office of Research
Education
1819 Innovation Hub
Libraries
Nursing
IT@UC
Technology
Engineering
Medicine
Computers

# Digital Integration

UC Libraries provides access to a wide range of Research Data and GIS services and resources for the campus community. Informationists and librarians are available to assist researchers in managing and preserving research data, finding and acquiring external data, and in utilizing GIS techniques and software. The library also provides a variety of computing and collaboration spaces to support researchers.

# DSC: Who

- **DSC**


James Lee, Assc Vice Provost for Digital Scholarship and Assc Dean of Libraries, Director of DSC


Kristen Burgess, Operational Manager


Lindsay Nickels, Program Coordinator


Ezra Edgerton, Data Visualization Developer


Erin McCabe, Digital Scholarship Library Fellow

- **RDS**


Amy Koshoffer, Asst Director of Research & Data Services


Tiffany Grant, Asst Director of Research & Data Services


Rebecca Olson, Business and Social Science Informationist


Don Jason, Health Informationist


Ted Baldwin, Director, Science & Engineering Libraries


Dorcas Washington, Data Analyst Specialist

- **Graduate Students (English, Computer Science, Business Analytics)**

www.libraries.uc.edu

# Digital Scholarship Center, Phase 1: From Initiative to Center

The DSC has assembled research groups that genuinely span multiple disciplines, with people trained to think very differently about every step in the research process.

Teams are composed of true partners across entire research lifecycle:

- o  Formulation of research questions
- o  Pitching grant proposals
- o  Dataset cleanup and manipulation
- o  Data analysis and visualization
- o  Argument formation
- o  Publication of findings

In 2017, the DSC received $900,000 from the Andrew W. Mellon Foundation to expand this mission.

www.dsc.uc.edu

## UC awarded a $700K grant from The Andrew W. Mellon Foundation

The renewal grant will advance and expand the Digital Scholarship Center's "catalyst" model

# How?

1. <u>Technical</u>:  Our platform adapts machine learning approaches to any text and image dataset for research projects. We apply these methods in a discipline-specific way.

2. <u>Human</u>:  We assemble teams to nurture these unconventional transdisciplinary research questions and partnerships – every collaborator has different goals and culture.

# Digital Integration: Technical Mission

University of CINCINNATI | LIBRARIES

## Data Management

- Open Science Framework
- Open Datasets
- Machine Actionable Data Management
- "Gentle Introduction to Data": Human Welcome

## Data Analysis: Digital Humanities to Bioinformatics

- Multimodal Deep Learning / Machine Learning
- Data Visualization Interfaces
- Mixed Methods: EDA + CDA
- Quantitative and Qualitative Data

## Digital Outcomes and Products

- Publications
- Conference Presentations
- Reference Datasets
- Project Websites / Apps
- Grants
- Machine Actionable Data Management

- High Performance Computing
- Storage
- Data Structures:
    - Format,
    - Fields,
    - Metadata,
    - Machine Readable

- Team Science Culture
- Writing + Environmental Scans
- Method Translation
- Student Training and NextGen MA / PhD

www.dsc.uc.edu

# Human Centered ML Needs Data Visualization

Machine Learning Platform: Model of Models (MoM)

# Machine Learning Platform: Model of Models (MoM)

- Two machine learning strategies used to observe the latent patterns in large corpora.
    - Topic modeling (Latent Dirichlet Allocation – LDA)
    - Word embeddings (word2vec, BERT)
- Aggregates multiple models in parallel to compare word usage across the parallel models.
- Clusters integrate topics from the 6 models into an aggregated "model of models":
    - Confirm consistent topics across all models
    - Reveal underrepresented topics that may not have appeared in a single model representation.
- Distributed parallel approach increases user confidence and interpretability of our models by bringing the most stable topics to the top tier of the model results



Anesthestic Action and Biochemistry

- Internal validation: Topic Coherence

- External validation:
  - Subject matter expert tagging of randomized 20% (N>1000) corpus
  - Blind human coder panel – Percentage agreement

- Replicability:
  - Pattern recognition capabilities of NLP methods as an information retrieval – and not a black-box classification – approach.
  - Provide models capable of evaluation by our panel of multiple independent coders.
  - Parallel replicates in each model (6-20 runs)

- Hybrid ML approach:
  - Human judgment of subject matter experts to verify and tag the model result – outperforms a purely machine-based analysis.
  - Semi-supervised learning for classification.

# Interoperability with Datasets, Medium and Large

- HTRC Extracted Features, JSTOR Data for Research, Chronicling America, Text Creation Partnership, Harvard Case Law.
- PubMed + PubMed Central, US Patent Claims, EPIC EHR FHIR Data Structures (IRB Approved).
- Social Media (Twitter, Instagram, Reddit).
- Small corpora: we'll help you read them.

model | of | models    | Get Started |   Learn    Analyze                              Account ⌄

## Select Database

Search Text

Filter Docs

Explore Docs

Select Vis

Set and Run Model

### Select your Database ⓘ

| COVID-19 Articles: 13.2k docs | Pubmed Abstract: 29.4M docs | Pubmed Central: 2.15M docs |
|---|---|---|
| 1.27 s/100 docs | 1.27 s/100 docs | 2.36 s/100 docs |

| Jstor Life Science: 825k docs | CaseLaw: 325k docs | Archaeology: 2.39k docs |
|---|---|---|
| 28.6 s/100 docs | 18.4 s/100 docs | 41.7 s/100 docs |

| Iowa Latin Canon: 2.98k docs | Ehealth Alzheimer: 129 docs | Text Creation: 69.9k docs |
|---|---|---|
| No Time Data | 1.30 s/100 docs | 72.8 s/100 docs |

| AC Justice: 297 docs | Anesthesiology: 28.0k docs | Ted Talks: 992 docs |
|---|---|---|
| No Time Data   3.60 s/100 docs | | |

## Documents ⓘ

| | |
|---|---|
| The Impact of Respiratory Viral Infection on Wheezing Illnesses and Asthma Exacerbations | 5.803449511528015 |
| Epidemiology of viral respiratory tract infections in a prospective cohort of infants and toddlers attending daycare | 5.163539171218872 |
| Population-based hospitalization incidence of respiratory viruses in community-acquired pneumonia in children younger than 5 years of age | 5.11656254529953 |
| Respiratory Syncytial Virus Coinfections With Rhinovirus and Human Bocavirus in Hospitalized Children | 4.919344365596771 |
| Virus Etiology of Airway Illness in Elderly Adults | 4.881579369306564 |
| Respiratory viral infections in a cohort of children during the first year of life and their role in the development of wheezing☆ | 4.8710708022117615 |

### The Impact of Respiratory Viral Infection on Wheezing Illnesses and Asthma Exacerbations

Overview ::: Respiratory viral-induced wheezing illnesses in young children Viral bronchiolitis is a LRTI typically associated with cough, tachypnea, retractions, and diffuse wheezing and rales [8], [9]. Bronchiolitis is a leading cause of hospitalizations in the first year of life, accounting for an estimated 120,000 infant hospitalizations annually [10]. In infants, the etiologic agents of bronchiolitis and other viral respiratory infections associated with wheezing include respiratory syncytial virus (RSV), rhinovirus, influenza, parainfluenza (PIV), adenovirus, and more recently identified viruses, such as human metapneumovirus (hMPV) and human boca virus (hBoV) [11], [12], [13], [14]. RSV causes epidemics of bronchiolitis and typically circulates in temperate climates during November to April with peaks in the winter months [11], [15], [16]. In tropical climates, peaks are related to temperature and level of rainfall [17]. RSV infects the majority of children during their first year of life and essentially all children show evidence of RSV infection by age 3 years [18]. The initial RSV infection is typically the most severe, causing lower respiratory tract disease, such as bronchiolitis, in 20% to 30% of infants [11], [18], [19]. Other viruses such as rhinovirus, PIV, and adenovirus circulate nearly year round with seasonal peaks of illness [10], [11], [19]

Tree ⓘ    Circle ⓘ    Network ⓘ

Doc-link
limit:1818

Circle – Paragraph Level/Square – Article Level

Public Health Neighborhood

Children/Asthma/Treatment Cluster

Treatment Neighborhood

Bench Science Neighborhood

cluster 17
increased expression levels patients production

## Clusters ⓘ

| cluster | #para, #docs | # topics | terms |
|---|---|---|---|
| 12 | 0,7872 | 10 | children patients asthma rsv age |
| 2 | 0,7190 | 19 | patients patient cases days years |
| 3 | 0,7100 | 8 | information countries people health research |
| 20 | 0,4774 | 13 | binding proteins protein peptides viruses |
| 17 | 0,4586 | 10 | expression production cells activation increased |
| 10 | 0,3609 | 7 | detection sensitivity samples detect positive |
| 7 | 0,3257 | 8 | infected model individuals transmission parameters |
| 29 | 0,3176 | 7 | cells expression infected incubated infection |
| 23 | 0,2567 | 8 | participants studies respondents information people |
| 6 | 0,2566 | 8 | samples positive detected tested patients |
| 11 | 0,2453 | 7 | human humans species viruses virus |
| 16 | 0,2350 | 5 | covid-19 patients |

# Model of Models: User Interface

Adaptable visualization outputs based on a single underlying model.

# Model of Models: User Interface
Adaptable visualization outputs based on a single underlying model.

# MoM: Digital Humanities

MoM: Social Sciences

Android Android Android Android Android Android Android Android Android Android Android Frame

FOR JONATHAN FOR JONATHAN FOR JONATHAN

Filter Select Review Curate Preview

# Digital Scholarship and Academic Health

- Partnership with College of Medicine, CCHMC Biomedical Informatics.
- Text Mining Electronic Health Records, Scientific Literature, Grant Databases, Social Media, Imaging.
- A Two-Way Street: Teaching STEM about qualitative data.

There are 69 labels, including -1.
Selected label: 2.0; # Edges: 43
Show By Cc / By Node.

| CC | Model/Topic | Positive Terms |
|---|---|---|
| 1.0 | 6/42 | **expression 0.07**, **gene 0.04**, stress 0.03, tissue 0.02, response 0.02, promoter 0.0 regulate 0.02 |
| | | evolution 0.03, specie 0.02 |
| | | dopsis 0.02, expression 0.02 |
| | | mutant 0.02, **function 0.02** |
| | | nt 0.06, control 0.03, **arabidopsis 0.03**, yeast 0.02, signal 0.0 |
| | | t 0.05, **function 0.03**, expression 0.02 |
| | | 4, family 0.02, **function 0.02** |
| | | clade 0.04, evolutionary 0.04, origin 0.03, loss 0.03, |
| | | 3, phylogeny 0.03 ... ncient 0.02, su |
| | | ence 0.03 ..., motif 0.0 |
| | | tation ... 0.03 |
| | | me ... 0.03, stress 0.0 |
| | | fy 0.02 |
| | | ion 0.02, study 0.02, **arabidopsis** |
| | | al 0.02 |
| | | plex 0.02, **show 0.02**, sequence 0.0 |
| | | 03, cell 0.03, **expression 0.02** |
| | | ucture 0.03 |
| | | 03, **function 0.03**, **show 0.02**, cell 0.02 |
| | | tion 0.06, **arabidopsis 0.04**, **gene 0.04**, conserve 0.03, divergence 0.02, spe |
| | | function 0.03 |
| | | ge ... **arabidopsis 0.04**, function 0.03 |
| | | evolve 0.03, evolutionary 0.02, angiosperm 0.02, dive |
| | | **protein 0.03**, unit 0.02, mammalian 0.02 |
| | | 0.06, mutant 0.04, nction 0.03, development 0.03, regu |
| | | ress 0.02 |
| | | ain 0.02 |
| | | **gene 0.04**, angiosperm 0.02 mily 0.02 |
| | | xpression 0.06, **show 0.03** ..., nction 0.03, express 0.03, **arabidops** |
| 1.0 | 5/86 | **function 0.06**, **arabidopsis 0.04**, pathw 0.03, development 0.03, regulate 0.03, 0.03, control 0.02, role 0.02 |
| 1.0 | 5/100 | ene 0.09 |
| 1.0 | 15/159 | **protein 0.26**, **domain 0.08** |
| 1.0 | 19/113 | pr ... yeast 0.02 |
| 1.0 | 10/27 | **protein 0.04**, **domain 0.03**, family 0.02, sequence 0.02 |
| 1.0 | 10/26 | **gene 0.06**, expression 0.03, **arabidopsis 0.02** |
| 1.0 | 16/98 | **protein 0.13** |

University of
CINCINNATI | LIBRARIES

# Uncertain Diagnosis Project
CCHMC Hospital Medicine Division

MetaMap

classify based on...

symptoms
disorders/diagnoses
tests/labs
specialties
treatments

| Tier 1 | Tier 2 |
|---|---|
| differential | if |
| etiology | likely |
| uncertain | consult* |
| unclear | could |
| possib* | may |
| consider | unknown |
| vs | although |
| abdominal/abdominal pain | suggest* |
| | however |
| | elevated |
| | broad |
| | further workup/work up |

Uncertain Case

Unclear diagnosis at this time, but differential would include post-viral gastroparesis/ileus, although severe intermittent abdominal pain would not be consistent with that diagnosis. Could have intermittent intussusception or volvulus, with a lead point of an enlarged lymph node in the setting of recent viral gastroenteritis. Renal colic is a possibility with the description of "writhing" in pain, but pain is not localized to the back or flanks and no blood of other abnormality seen on UA. Biliary colic could be considered, although would be unusual in her age range and without associated with food. Appendicitis remains on the differential, as was not visualized on ultrasound, but exam findings not consistent with the diagnosis.

Vent Notes Prototype

| Corpus : Covid | Term : severe_acute_respiratory_syndrome_coronavirus_2-OR-covid-19-OR-sars_cov_2 | # Docs : 31,818 | # Topics : 40 | Stopwords : | Years : - |

## Documents ⓘ

| Title | Score |
|---|---|
| Response to | 10.3162498474 |
| Combating Devastating COVID -19 by Drug Repurposing | 4.48861750960 |
| Favipiravir versus Arbidol for COVID-19: A Randomized Clinical Trial | 3.87015905976 |
| Title: What do we know about remdesivir drug interactions? | 3.79288643598 |
| | 3.75096887350 |
| Interferon beta-1a for COVID-19: critical importance of the administration route | 3.61853513121 |
| Chloroquine and hydroxychloroquine for COVID -19: A word of caution | 3.58210587501 |

## Doc View



Tree ⓘ  Circle ⓘ  Network ⓘ

Doc-link limit:4545

Circle – Paragraph Level/Square – Article Level

Public Health 'Neighborhood'

Treatment 'Neighborhood'

Hydroxychloroquine topics

Bench Science 'Neighborhood'

| cluster | #para, #docs | # topics | terms |
|---|---|---|---|
| 16 | 0,23418 | 15 | patients patient pandemic resources severe |
| 1 | 0,18905 | 17 | patients severe study treatment cases |
| 7 | 0,16931 | 13 | pandemic people crisis impact students |
| 8 | 0,16355 | 10 | cases data infected model population |
| 0 | 0,8368 | 6 | cases china outbreak transmission infected |
| 19 | 0,6891 | 7 | cases patients patient diagnosis images |
| 10 | 0,4979 | 9 | patients increased il-6 inflammation cytokine |
| 37 | 0,4670 | 5 | ppe patients risk patient masks |
| 3 | 0,4287 | 7 | positive patients negative samples days |
| 20 | 0,4140 | 7 | treatment drugs patients efficacy hydroxychloroquine |

# Second Phase: Multimodal ML



anishkapoor

Removing Hashtags with fewer than 10 occurrences

# Second Phase: Jupyter Notebooks Pipeline

## Organization of Data

What do you notice about the format of the data above?

Each sentence is already *tokenized* split into a series of word and punctuation strings, with whitespace removed. This saves a lot of time having to do this work ourselves, manually.

To start to organize our data, let's put these sentences into a pandas DataFrame, an object which has a format very similar to an Excel spreadsheet. We will first make two spreadsheets (one for news, and one for romance), and then combine them into one. We will also add the category each sentences came from, which will be our labels for each sentence and its associated feature representation (which we will build ourselves).

```
In [ ]: ndf = pd.DataFrame({'sentence': news_sent,
                           'label':'news'})
        rdf = pd.DataFrame({'sentence':romance_sent,
                           'label':'romance'})
```

```
In [ ]: # combining two spreadsheets into 1
        df = pd.concat([ndf, rdf])
```

Let's see what this DataFrame looks like

```
In [ ]: df
```

```
In [ ]: df.head()
```

### So how many texts are there of each type?

```
In [ ]: df['label'].value_counts()
```

### What if we want to visualize that information?

We first create a figure and axes on which to draw our charts using plt.subplots(). Each chart is one axes, and a figure can contain multiple charts. Our data is encapsulated in df['label'].value_counts(), which is itself a dataframe. We then tell the Pandas to visualize the dataframe as a bar chart using .plot.bar(ax=ax, rot=0). The ax keyword tells Pandas which chart in the figure to plot, and the rot keyword controls the rotation of the x axis labels.

```
In [ ]: fig, ax = plt.subplots()
        _ = df['label'].value_counts().plot.bar(ax=ax, rot=0)
        fig.savefig("categories_counts.png", bbox_inches = 'tight', pad_inches = 0)
```

We have slightly more news texts that romance texts, which we should keep in mind as we go ahead with classification.

## Extracting Features

### Defining Features

What should we use as features for the datset? What did we use for the fruit example before?

| Object | Height | Width | Color | Mass | Round? |
|--------|--------|-------|-------|------|--------|
| Apple | 6cm | 7cm | Red | 330g | TRUE |

# Let's Work Together

- Model of models platform: https://modelofmodels.io
- DSC website: http://dsc.uc.edu

Second Grant Objectives:

- 15 Subgrants through Mellon Foundation Award using our technologies.

- Expand use cases and projects deploying MoM and Jupyter library.

- Experimental use of MoM for data services.

- Partnerships with external collaborators engaged in digital scholarship.

Acknowledgements:
We are grateful to the Andrew W. Mellon Foundation, Public Knowledge Program for their support.