# Collaboration and Re-Use
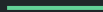
## Experiences with institutional data catalogs

Nicole Contaxis, MLIS
Lead, NYU Data Catalog

# Agenda

1. Background on the NYU Data Catalog

2. Background on the Data Discovery Collaboration Project

3. Stories of Data Re-Use

# Background on the NYU Data Catalog

# The NYU Data Catalog: An Overview

- Provides a standardized metadata schema to describe data



NYU HEALTH SCIENCES LIBRARY — NYU Data Catalog

## Neurological Emergencies Outcomes at NYU

NYU Dataset

Alternate Titles(s): NEON

UID: 10330

Author(s): Ariane Lewis*, Aaron Lord
* Corresponding Author

### Description

This dataset was collected as part of a combined retrospective and prospective cross-sectional study to establish risk factors for infection after intracerebral hemorrhage and subarachnoid hemorrhage and to determine the impact of those infections on long-term outcomes. Data was harvested from Tisch Hospital records from January 2013 to December 2014 retrospectively and from January 2015 to the present prospectively, and the study aims to recruit an additional 1,000 patients by 2027.

Patients are included in the study if they are over 18 years of age and have a new diagnosis of intracerebral hemorrhage or subarachnoid hemorrhage requiring admission to or consultation by acute neurology faculty members at NYU Langone Medical Center, and for prospective patients, if the patient or next of kin consent to participate in follow-up phone interviews at 3 months and 12 months.

Data that will be collected from both retrospectively and prospectively enrolled patients include:

- Admission data (hospital admission information, history of present illness)
- Admission vital signs (BMI, weight, height, temperature, heart rate, respiratory rate, blood pressure)
- Admission labs (serum chemistries, blood count, coagulation)
- Baseline data (demographics, medications, past medical history, social history, family history)
- Admission examination (Hunt/Hess grade, Glasgow Coma Scale (GCS), NIH Stroke Scale (NIHSS), premorbid Modified Rankin Scale (MRS)
- Admission CT scan and angiogram results
- Hospital procedures, surgical treatments, medical treatments

**Access via Data Request Form**
Form to request access

**Access Restrictions**
Application Required
Author approval required

**Access Instructions**
Please contact Dr. Ariane Lewis for information on how to apply for access to this dataset.

**Data Type**
Administrative
Clinical Measures
Imaging
Interviews

**Study Type**
Observational

**Dataset Format(s)**
SPSS, Stata, Microsoft Excel, CSV

**Data Collection Instruments**
Glasgow Outcome Scale
Modified Rankin Scale
Barthel Index
Neuro-QOL

# The NYU Data Catalog:
# An Overview

- Provides a standardized metadata schema to describe data

- Makes research data discoverable regardless of where it is stored



**NYU** HEALTH SCIENCES LIBRARY   **NYU Data Catalog**

## Neurological Emergencies Outcomes at NYU

**NYU Dataset**

Alternate Titles(s): NEON
UID: 10330
Author(s): Ariane Lewis*, Aaron Lord
* Corresponding Author

### Description

This dataset was collected as part of a combined retrospective and prospective cross-sectional study to establish risk factors for infection after intracerebral hemorrhage and subarachnoid hemorrhage and to determine the impact of those infections on long-term outcomes. Data was harvested from Tisch Hospital records from January 2013 to December 2014 retrospectively and from January 2015 to the present prospectively, and the study aims to recruit an additional 1,000 patients by 2027.

Patients are included in the study if they are over 18 years of age and have a new diagnosis of intracerebral hemorrhage or subarachnoid hemorrhage requiring admission to or consultation by acute neurology faculty members at NYU Langone Medical Center, and for prospective patients, if the patient or next of kin consent to participate in follow-up phone interviews at 3 months and 12 months.

Data that will be collected from both retrospectively and prospectively enrolled patients include:

- Admission data (hospital admission information, history of present illness)
- Admission vital signs (BMI, weight, height, temperature, heart
- Rankin Scale (MRS)
- Admission CT scan and angiogram results
- Hospital procedures, surgical treatments, medical treatments

**Access via Data Request Form**
Form to request access

**Access Restrictions**
Application Required
Author approval required

**Access Instructions**
Please contact Dr. Ariane Lewis for information on how to apply for access to this dataset.

**Data Type**
Administrative
Clinical Measures
Imaging
Interviews

**Study Type**
Observational

Barthel Index
Neuro-QOL

## Does Not Store Data

# The NYU Data Catalog: An Overview

- Provides a standardized metadata schema to describe data

- Makes research data discoverable regardless of where it is stored

- Open source
  - Code on GitHub
  - Documentation on OSF

# Background on the Data Discovery Collaboration Project

"To enhance discovery of data and other research products in order to maximize their value"

# Background: Example Data Catalog Record



**External Repository**

# Background: Example Data Catalog Record



**Restricted Access Data**

# Background: Example Data Catalog Record



**Datasets Via Author Only**

# Background: Example Data Catalog Record



**Electronic Health Record Data**

# Background: Example Data Catalog Record



**Datasets in Supplemental Files of Published Articles**

# What this means for re-use

- Examples of facilitating re-use with discovery metadata

- Librarians on the ground gaining experience with facilitating re-use

# Stories of Re-Use

# Local Experts & Re-Use

- Metadata element designed to facilitate collaboration and re-use

- At NYU, generally used on large, third party datasets

**NYU HEALTH SCIENCES LIBRARY** — **NYU Data Catalog**

## National Health and Nutritional Examination Survey

Alternate Titles(s): NHANES

UID: 10003

**Description**
The National Health and Nutrition Examination Survey (NHANES) is a program of studies designed to assess the health and nutritional status of adults and children in the United States. The survey is unique in that it combines interviews and physical examinations. The NHANES interview includes demographic, socioeconomic, dietary, and health-related questions. The examination component consists of medical, dental, and physiological measurements, as well as laboratory tests administered by highly trained medical personnel.

**Publisher**
United States - Centers for Disease Control and Prevention (CDC)

**Timeframe**
1957 - Present

**Geographic Coverage**
United States

**Local Expert for NYU**
Heather Gold
James Slover
Jiyoung Ahn
Judith Goldberg
Leo Trasande
Lorna Thorpe
Michael Weitzman
Niyati Parekh
Terry Gordon

**Subject Domain**

Access via NHANES

**Access Restrictions**
Free to All

**Access Instructions**
NHANES data is available on the website and is organized by year. Each year of NHANES data provides users with analytic guidelines, response rates, population totals and a web tutorial. Users can download demographics, examination, laboratory, questionnaire, and limited access data directly from the website. Selecting a dataset using the DOC file will take users to a description of that dataset including the data documentation, codebook, frequencies. Selecting the Data file will download the dataset in .XPT or RDC format.

**Data Type**
Surveys

**Dataset Format(s)**
SAS, PDF, SUDAAN

**PubMed Search**
View articles which use this dataset

**Related Datasets**
New York City Health and Nutrition Examination Study

# Local Experts & Re-Use

- Metadata element designed to facilitate collaboration and re-use

- At NYU, generally used on large, third party datasets



**Local Expert for NYU**
Heather Gold
James Slover
Jiyoung Ahn
Judith Goldberg
Leo Trasande
Lorna Thorpe
Michael Weitzman
Niyati Parekh
Terry Gordon

**NYU** HEALTH SCIENCES LIBRARY  **NYU Data Catalog**

## National Health and Nutritional Examination Survey

Alternate Titles(s): NHANES

UID: 10003

**Description**

The National Health and Nutrition Examination Survey (NHANES) is a program of studies designed to assess the health and nutritional status of adults and children in the United States. The survey is unique in that it combines interviews and physical examinations. The NHANES interview includes demographic, socioeconomic, dietary, and health-related questions. The examination component consists of medical, dental, and physiological measurements, as well as laboratory tests administered by highly trained medical personnel.

**Publisher**

United States - Centers for Disease Control and Prevention (CDC)

**Timeframe**

1957 - Present

**Geographic Coverage**

United States

**Local Expert for NYU**
Heather Gold
James Slover
Jiyoung Ahn
Judith Goldberg
Leo Trasande
Lorna Thorpe
Michael Weitzman
Niyati Parekh
Terry Gordon

**Subject Domain**

**Access via NHANES**

**Access Restrictions**
Free to All

**Access Instructions**
NHANES data is available on the website and is organized by year. Each year of NHANES data provides users with analytic guidelines, response rates, population totals and a web tutorial. Users can download demographics, examination, laboratory, questionnaire, and limited access data directly from the website. Selecting a dataset using the DOC file will take users to a description of that dataset including the data documentation, codebook, frequencies. Selecting the Data file will download the dataset in .XPT or RDC format.

**Data Type**
Surveys

**Dataset Format(s)**
SAS, PDF, SUDAAN

**PubMed Search**
View articles which use this dataset

**Related Datasets**
New York City Health and Nutrition Examination Study

# Successes

- Local experts report being contacted about those datasets

- Local experts report becoming co-authors on papers generated from that data

# Concerns

- Local experts complain that they are contacted by people outside of the institution that are not viable collaborators

- Local experts express concern about the amount of time their volunteer work takes

# Takeaways

- Researchers have questions, even on incredibly well-documented data

- Researchers have limited time so answering questions can be difficult

- Researchers enter into collaborations based on conversations on datasets and re-use

- Addressing this gap in responsibility (e.g., acting as a local expert) is a key part of ensuring future success of initiatives like the NYU Data Catalog

# Residents Research Practicum

- Third year residents worked with researchers and a librarian to develop original research through re-used datasets

- Residents worked as a group and located a dataset for re-use through through the NYU Data Catalog

Fred LaPolla, Research and Data Librarian, Lead of Data Education and Course Director

# Takeaways

- Residents and practicum leads needed to discuss the datasets with the original creators, even as the dataset was well-documented

- A poster was accepted at a conference based on the resident's work

- Residents were not able to publish on their research due to limitations from one of the funders from the original research

- We only know of this interaction due to librarian participation - further investigation into user tracking is necessary

# Future Efforts

- Further investigation into re-use use cases with the Data Discovery Collaboration Project

- Creation of infrastructure to help "Local Experts" manage their requests

# Thank you

Nicole Contaxis, MLIS
Lead, NYU Data Catalog