

# Foundations of Data Curation Review!!

# Data Curation

## What is it?

Data science = (Data Curation + Data Analytics)

Data science has two components:

<b>Data curation:</b>	Ensuring that data can be <i>efficiently</i> and <i>reliably</i> found and used
<b>Data analytics:</b>	Employing specific techniques to extract knowledge from data

**Data curation** is concerned primarily with the *management of data*

in order to better support *the analysis of data*

# Data curation is the larger part of data science

Not only is data curation essential for reliable efficient analysis,  
but most of the cost associated with using data is, by far, in curation, not analysis,  
and most of the workforce needs are, also by far, in curation, not analysis.

Ask any data manager in industry will tell you:

it is curatorial work where they make the largest investment,  
of money, staff, time, and effort.

# Areas of curatorial activities

<b>Collection:</b>	Support the collection and acquisition of data (and documentation of that, <i>throughout list</i> )
<b>Organization:</b>	Employ an appropriate data model and use appropriate standards
<b>Storage:</b>	Support reliable and effective storage
<b>Preservation:</b>	Ensure that data will be understandable and useable in the future
<b>Discoverability:</b>	Support the ability to search for and locate relevant data
<b>Access:</b>	Support the ability to retrieve and distribute data
<b>Workflow:</b>	Support the ability to systematize data workflows
<b>Identification:</b>	Support the ability to identify, authenticate, and validate data
<b>Integration:</b>	Support integration of data from different sources using different data models
<b>Reformatting:</b>	Support reformatting for use by different tools or to match new format standards
<b>Reproducibility:</b>	Support ability to reproduce results, ensuring scientific validity
<b>Sharing:</b>	Support sharing data between researchers, teams, and institutions.
<b>Communication:</b>	Support representation, publishing, and visualizations that provide insight
<b>Provenance:</b>	Support identifying what inputs and calculations are responsible for data values
<b>Modification:</b>	Support management of corrections and updates
<b>Compliance:</b>	Ensure compliance to legal, regulatory, and local policy requirements
<b>Security:</b>	Ensure that data is secure from tampering or inappropriate access and distribution

# Methods of curatorial action

Five categories stand out as particularly important:

## **Analysis**

To determine needs, develop relevant data models and *metadata*, and reformat, correct, or update data.

## **Documentation**

To record essential information (typically via *metadata*)

## **System design and implementation**

To support all data curatorial activities

To support the generation and use of data documentation and processing documentation

## **Policy**

To specify objectives, procedures, practices, and formats.

## **Process**

To ensure success and efficiency by managing the development of appropriate organizational units and roles, providing training, advocating for change, and managing curatorial activities.