# DEFINITIONS

I535: MANAGEMENT, ACCESS, AND USE OF BIG AND COMPLEX DATA

INDIANA UNIVERSITY BLOOMINGTON

1

---

2

## A BIT OF HISTORY…

▶ 1995
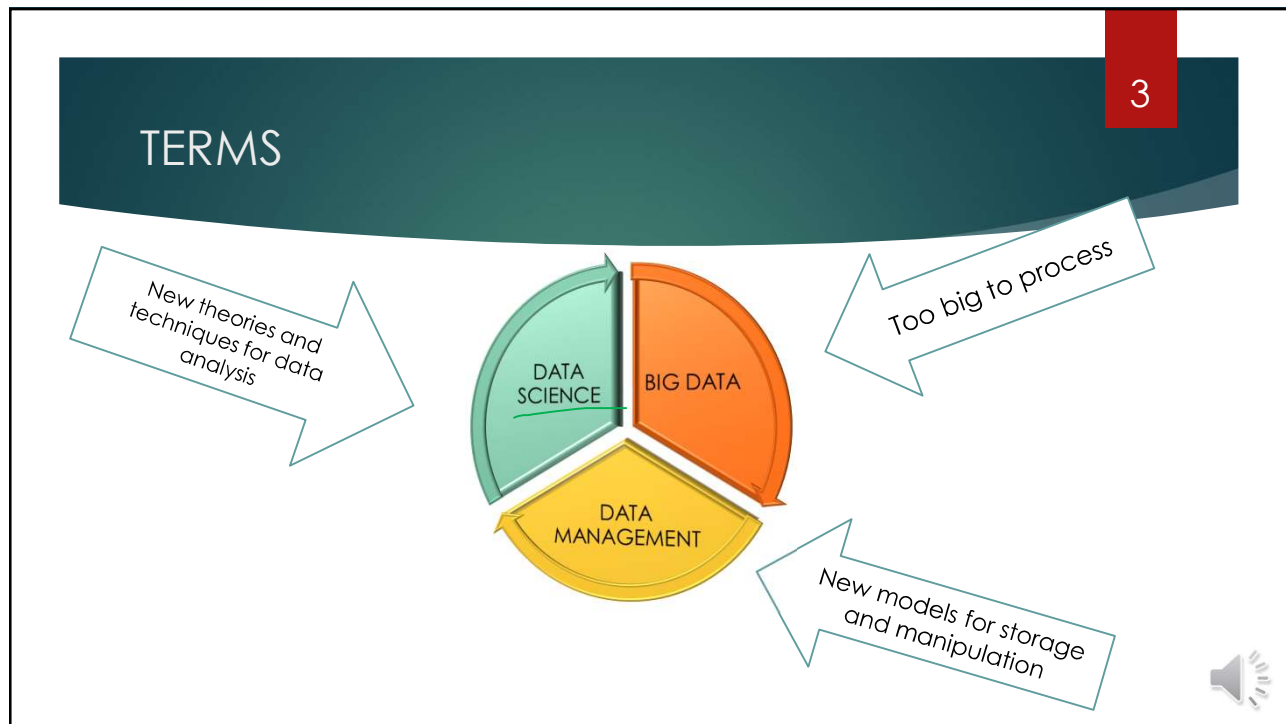
"... massive data sets pose a great challenge to scientific research … Today's data sets … have now outstripped the capability of previously developed data measurement, data analysis, and data visualization tools."
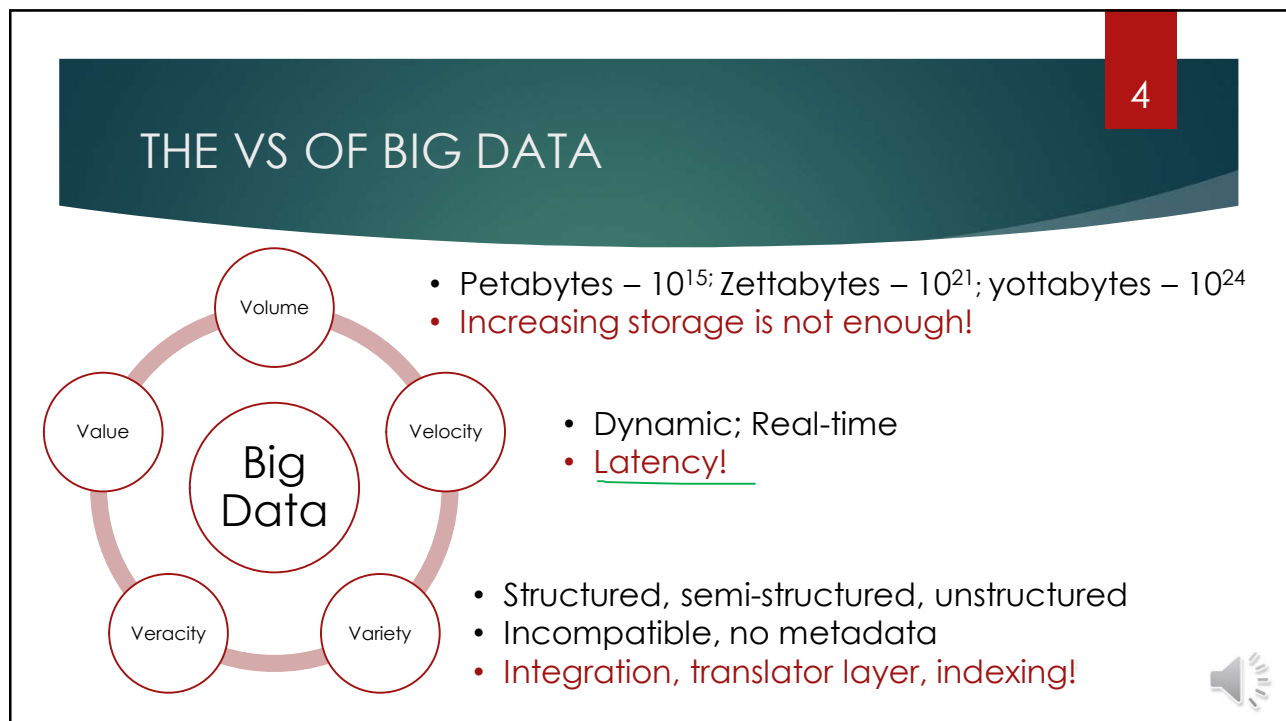
▶ 1997

"Big data objects are just that -- single data objects (or sets) that are too large to be processed by standard algorithms and software on the hardware one has available."
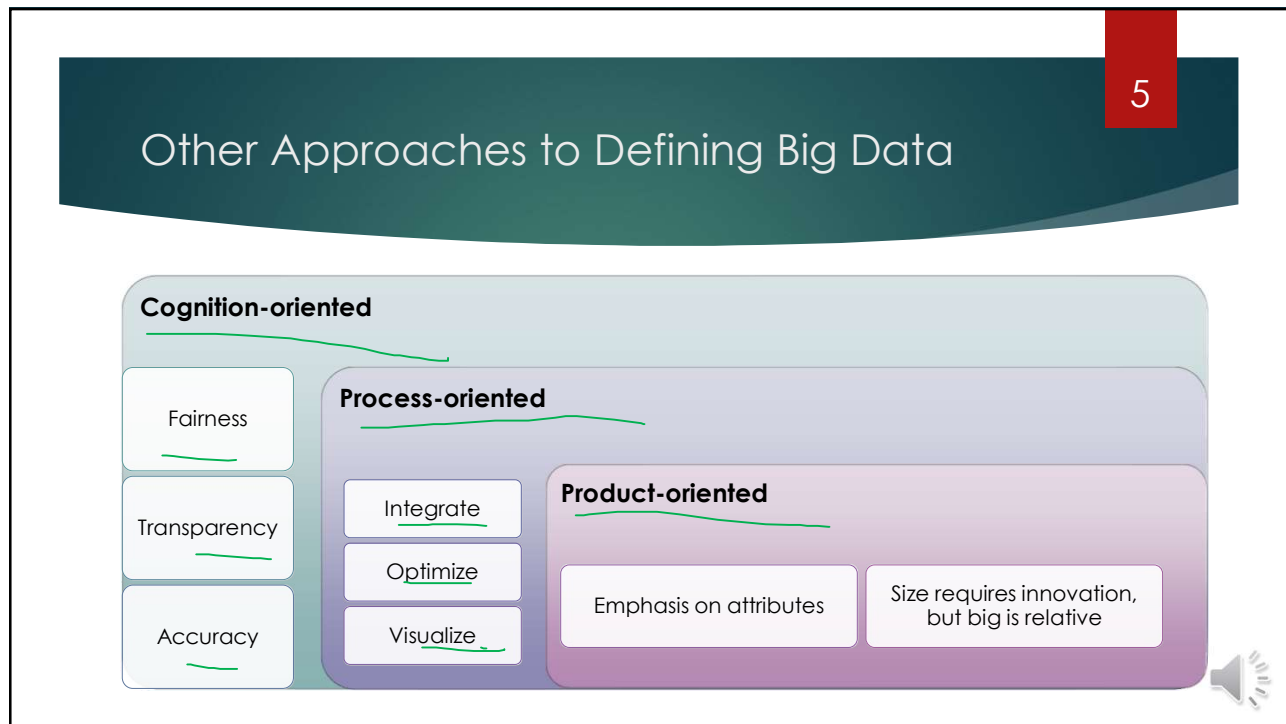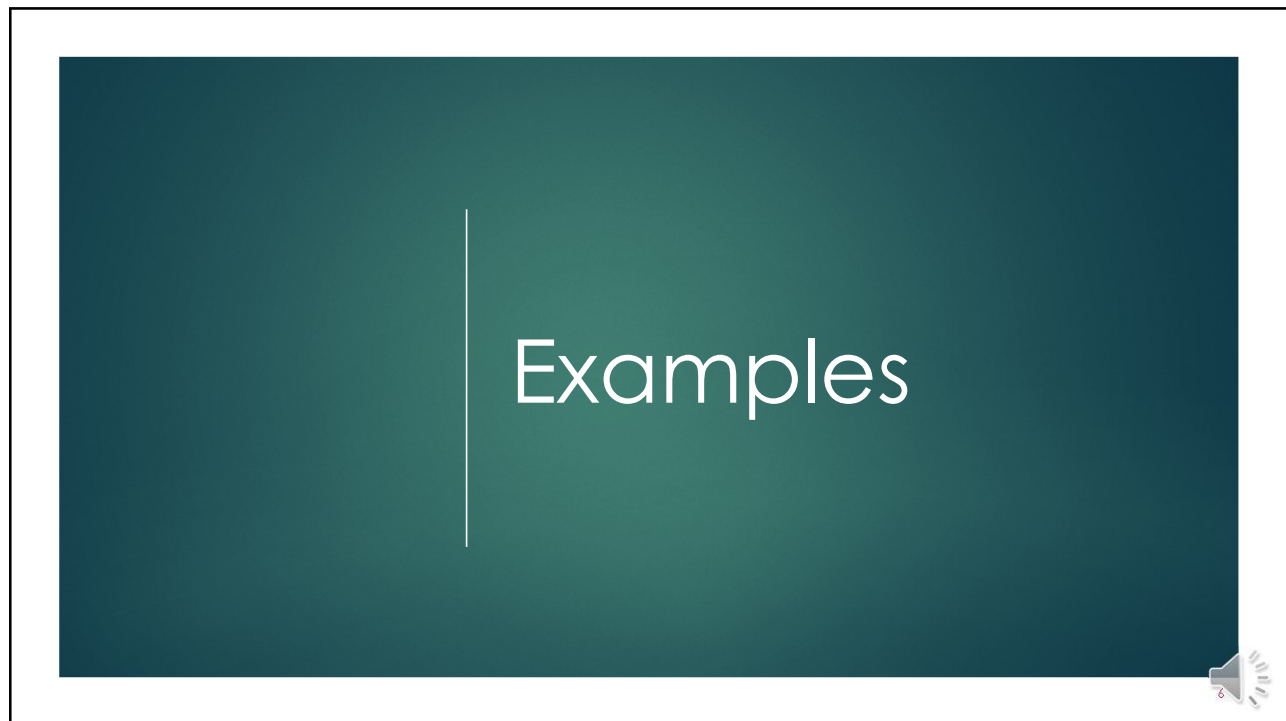
2

## TERMS

3

New theories and techniques for data analysis

DATA SCIENCE

BIG DATA

Too big to process

DATA MANAGEMENT

New models for storage and manipulation

3

## THE VS OF BIG DATA

4

Volume

Value

Velocity

Big Data

Veracity

Variety

- Petabytes – $10^{15}$; Zettabytes – $10^{21}$; yottabytes – $10^{24}$
- Increasing storage is not enough!

- Dynamic; Real-time
- Latency!

- Structured, semi-structured, unstructured
- Incompatible, no metadata
- Integration, translator layer, indexing!

4

I535: Management, Access, and Use of Big
and Complex Data

Other Approaches to Defining Big Data

5

**Cognition-oriented**

Fairness

Transparency

Accuracy

**Process-oriented**

Integrate

Optimize

Visualize

**Product-oriented**

Emphasis on attributes

Size requires innovation, but big is relative

5



Examples

6

I535: Management, Access, and Use of Big and Complex Data

## Google Maps

7

- 170 billion images collected in 87 countries
- Quantitative information like road length, terrain, and distances
- Information coming from individuals and local governments
- Images to form a layered and detailed map
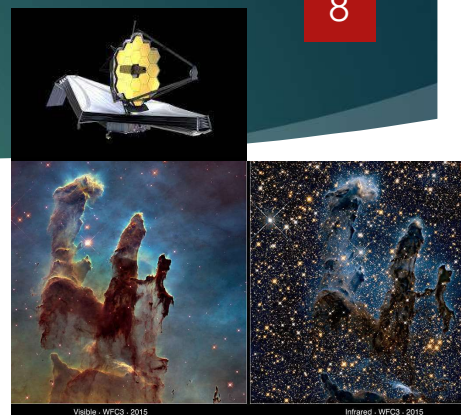- Machine learning algorithms to automate mapping, routing, and other tasks

7

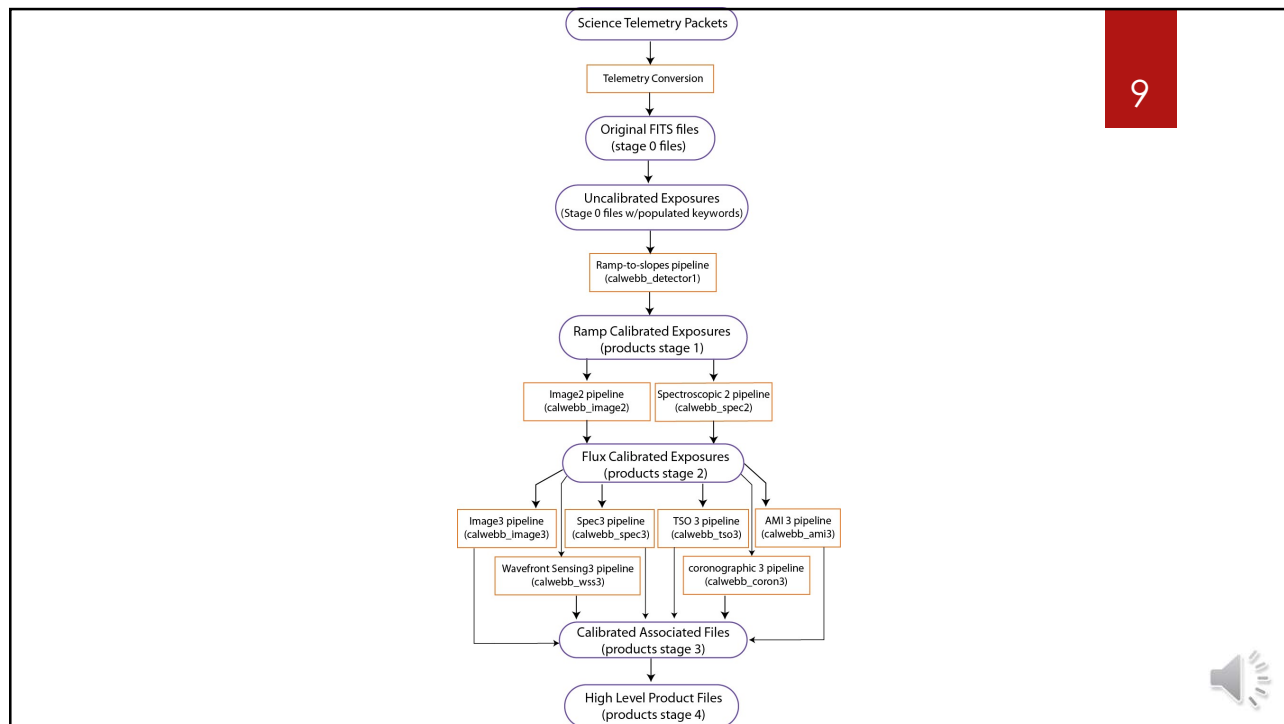## The James Webb Space Telescope (JWST)

8

https://webb.nasa.gov

➢ Volume: 57GB per day

➢ Variety: multiple instruments and measurements

➢ Velocity: constant stream, with intermediate local storage

8

I535: Management, Access, and Use of Big and Complex Data

9



## More examples on Canvas

▶ CERN
▶ Experian
▶ Facebook
▶ Netflix
▶ Olympic cycling team
▶ SKA telescope

10

## Goals of Data Management

11

Planning, implementation, and oversight
of acquiring, delivering, and enhancing the value of data

Understand **the needs** of the organization and its stakeholders
Store and **protect** the data assets

Improve **the quality** of data and information

Ensure **privacy** and confidentiality of the data

Maximize the effective **use and value** of the data

DAMA: Data Management Association https://dama.org/

11

## Tools and Skills of Big Data
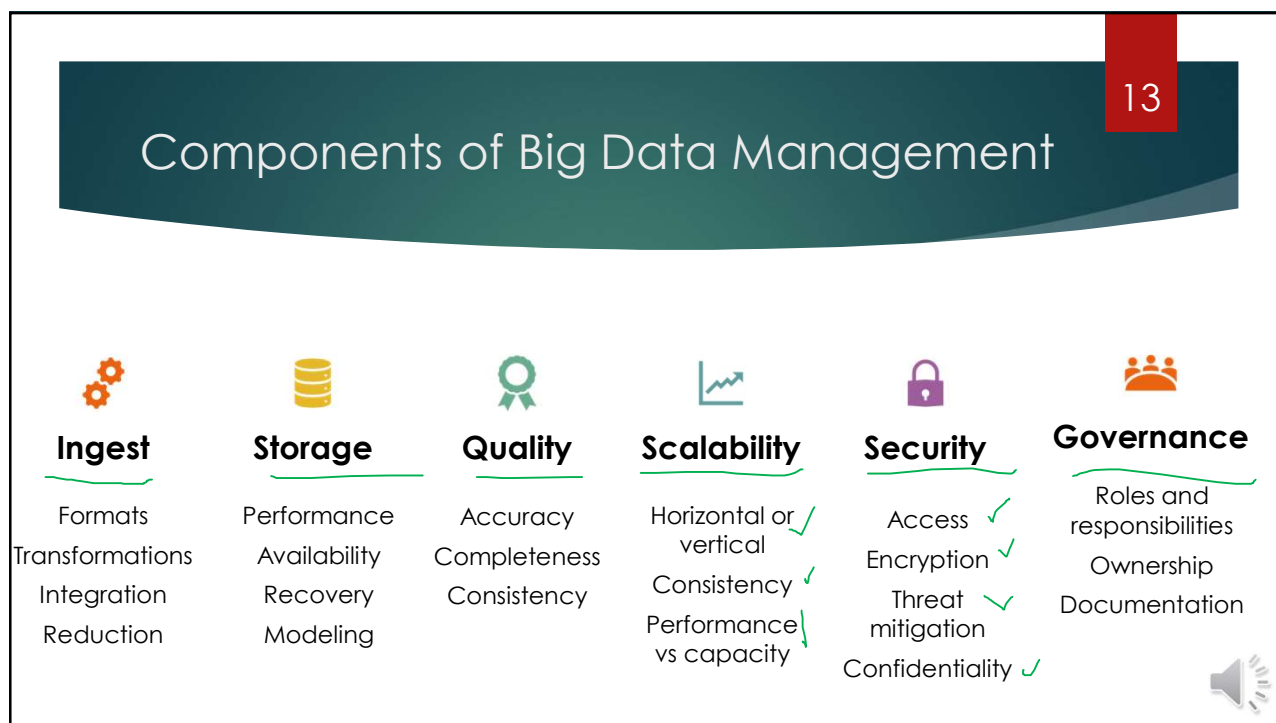
12

**Data Analyst** ✓
- Strong understanding of RDBMS tools
- Programming languages: Python, R, Java, C++, Matlab
- VBA and SQL skills
- Tensorflow and/or Keros
- Experience with Salesforce CRM ✓
- Experience with data visualization software (Tableau or Spotfire)
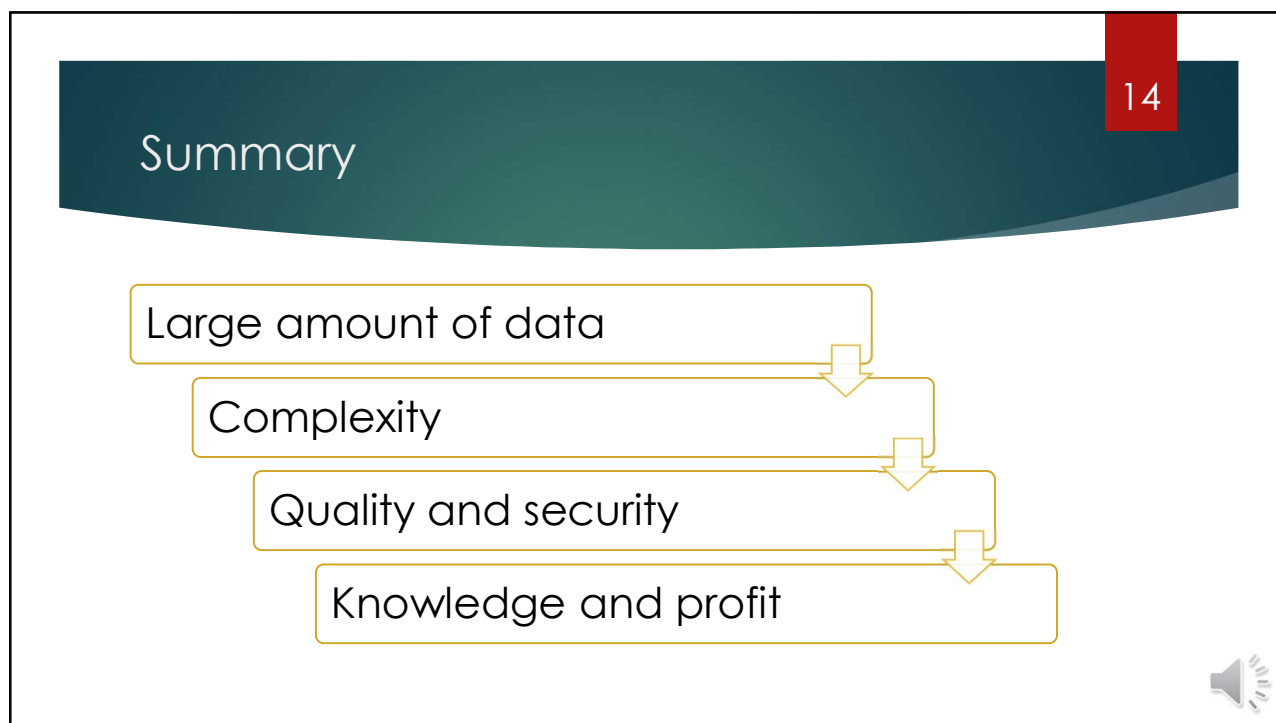
**Data Engineer** ✓
- Compute cluster/high-performance computing environment
- Linux/Unix/MacOS as software development platform
- *Architecting distributed systems, creating reliable pipelines, combining data sources*
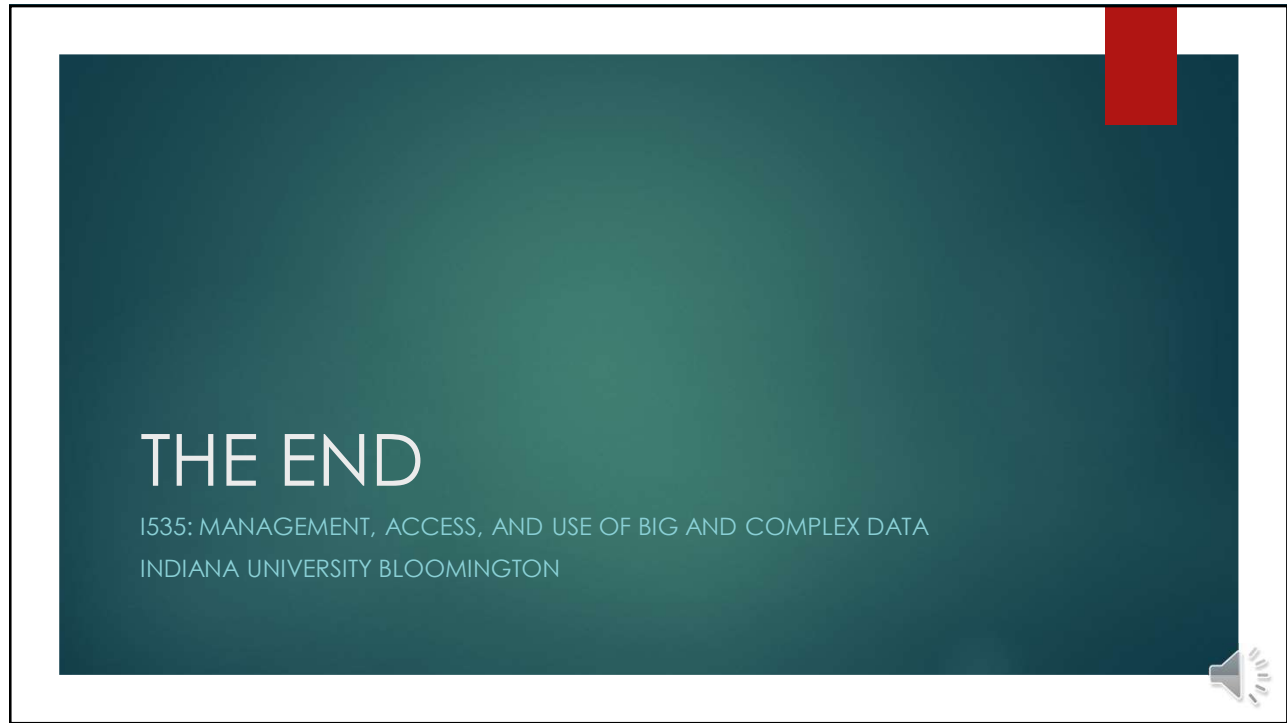- Microsoft Azure
- Data Bricks experience

12

I535: Management, Access, and Use of Big
and Complex Data

# Components of Big Data Management

**Ingest**
Formats
Transformations
Integration
Reduction

**Storage**
Performance
Availability
Recovery
Modeling

**Quality**
Accuracy
Completeness
Consistency

**Scalability**
Horizontal or vertical ✓
Consistency ✓
Performance vs capacity

**Security**
Access ✓
Encryption ✓
Threat mitigation
Confidentiality ✓

**Governance**
Roles and responsibilities
Ownership
Documentation

13

# Summary

Large amount of data

Complexity

Quality and security

Knowledge and profit

14

I535: Management, Access, and Use of Big and Complex Data

THE END

I535: MANAGEMENT, ACCESS, AND USE OF BIG AND COMPLEX DATA

INDIANA UNIVERSITY BLOOMINGTON

15