# Visual Data Science

**Andrea Julca**
College of Information Studies
University of Maryland, College Park

# MOTIVATING QUESTIONS:
How do I go from nothing to my hypothesis?
Designing interactive visualizations for exploratory analysis

# OUTLINE

- Defining data science

  - Extract, transform, load (ETL)

- Exploratory analysis and modeling
  - NLP – Natural Language Processing

- Streaming visualization

# What is a "data scientist?"

"Data Scientist (n.): Person who is better at statistics than any software engineer and better at software engineering than any statistician." - Josh Wills

- Something of a marketing term, but careers and formal data science programs have sprung up around the concept

# Data science competencies

Anderson et al. (2014):
- **Information retrieval**
- **Large or streaming data sets**
- **Databases**
- AI and statistical techniques
- Software development and algorithms
- Mathematics
- Communication
- Social, ethical, and legal awareness

# Data science competencies

Anderson et al. (2014):

- Information retrieval
- Large or streaming data sets
- Databases
- **AI and statistical techniques**
- **Software development and algorithms**
- **Mathematics**
- Communication
- Social, ethical, and legal awareness

# Data science competencies

Anderson et al. (2014):
- Information retrieval
- Large or streaming data sets
- Databases
- AI and statistical techniques
- Software development and algorithms
- Mathematics
- **Communication**
- **Social, ethical, and legal awareness**

# Data science workflow

**1. Scope out the problem or question**
**2. Knowledge search: Research and sensemaking**
3. Data retrieval; extract, transform, load (ETL)
4. Exploratory analysis
5. Modeling
- System-building [sometimes]
- Versioning/archival
- Communication

# Data science workflow

1. Scope out the problem or question
2. Knowledge search: Research and sensemaking
3. **Data retrieval; extract, transform, load (ETL)**
4. Exploratory analysis
5. Modeling
- System-building [sometimes]
- Versioning/archival
- Communication

# Extract: Information retrieval

- Information systems: Get data from a database

- Information studies: "Everything is data"
  - Tables
  - Text
  - Images
  - Media files (video, audio)
  - Interviews?
  - Artifacts??
  - Other examples?

# Extract: "Webscraping"

# Transform [Info Systems]

Reshaping and restructuring data for the target database

- Clean
- Filter
- Apply models
- Business rules
- Aggregate
- Et cetera

# Transform [Mathematics]

- Geometry:
  - Reflect
  - Rotate
  - Scale (resize)
  - Translate (shift position)

- Generally:
  - An invertible function mapping one domain to another

# Why not both?

# Data science workflow

1. Scope out the problem or question
2. Knowledge search: Research and sensemaking
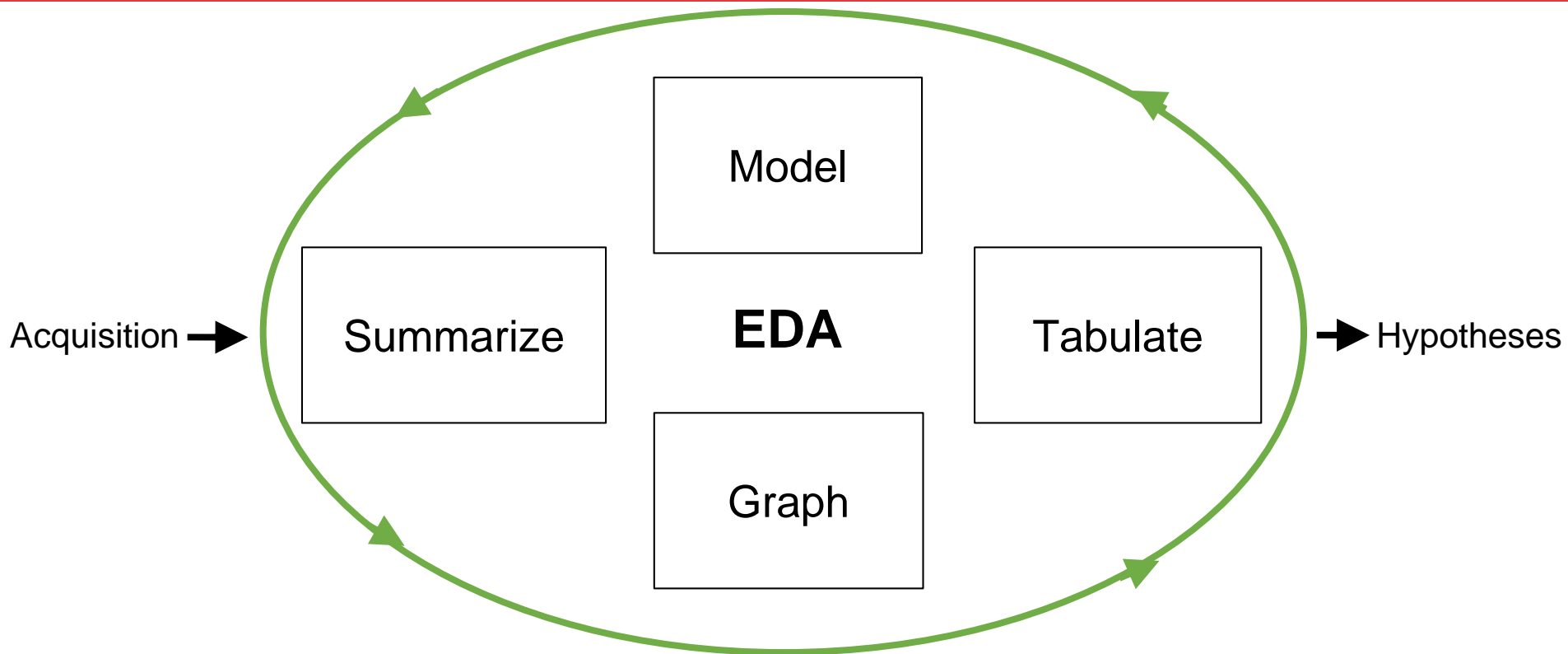3. Data retrieval; extract, transform, load (ETL)
4. **Exploratory analysis**
5. **Modeling**
   - System-building [sometimes]
   - Versioning/archival
   - Communication

# Exploratory Data Analysis (EDA)



Acquisition → 

**EDA**

Model

Summarize

Tabulate → Hypotheses

Graph

Based on Tukey, J.W. *The Future of Data Analysis*. 1962.
*and* Exploratory data analysis. 1977.

[demonstration]

# In-Class Exercise 1

1. **git clone/DL github repo: …/DreaJulca/streamvis-lecture**
2. **Run Bind.R**
3. **(On paper) design and sketch a better (but also transitioning) visualization**

# Load (Stage/Publish) & Archive

- We've "loaded" the data from our chosen website into the *R* environment
  - Not a reliable way to warehouse. Why? *Low permanence*
  - Also not a great publication / communication platform
- In a more complete information or business system, we would:
  - Perform further transformations
  - Load into database with well-defined schema (*higher permanence*)
  - We're skipping that today

# Data science competencies

Anderson et al. (2014):

- Information retrieval
- **Large or streaming data sets**
- Databases
- AI and statistical techniques
- Software development and algorithms
- Mathematics
- Communication
- Social, ethical, and legal awareness
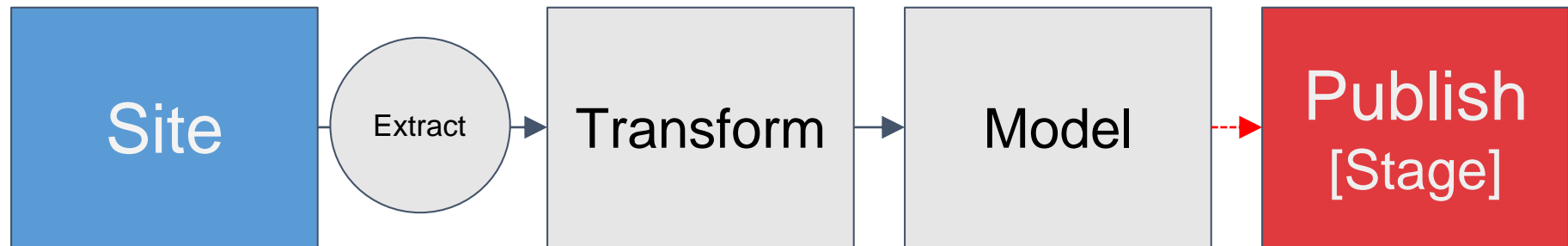
# Streaming Visualization

What is "streaming?"
- Transfer of continuously-generated data in real time
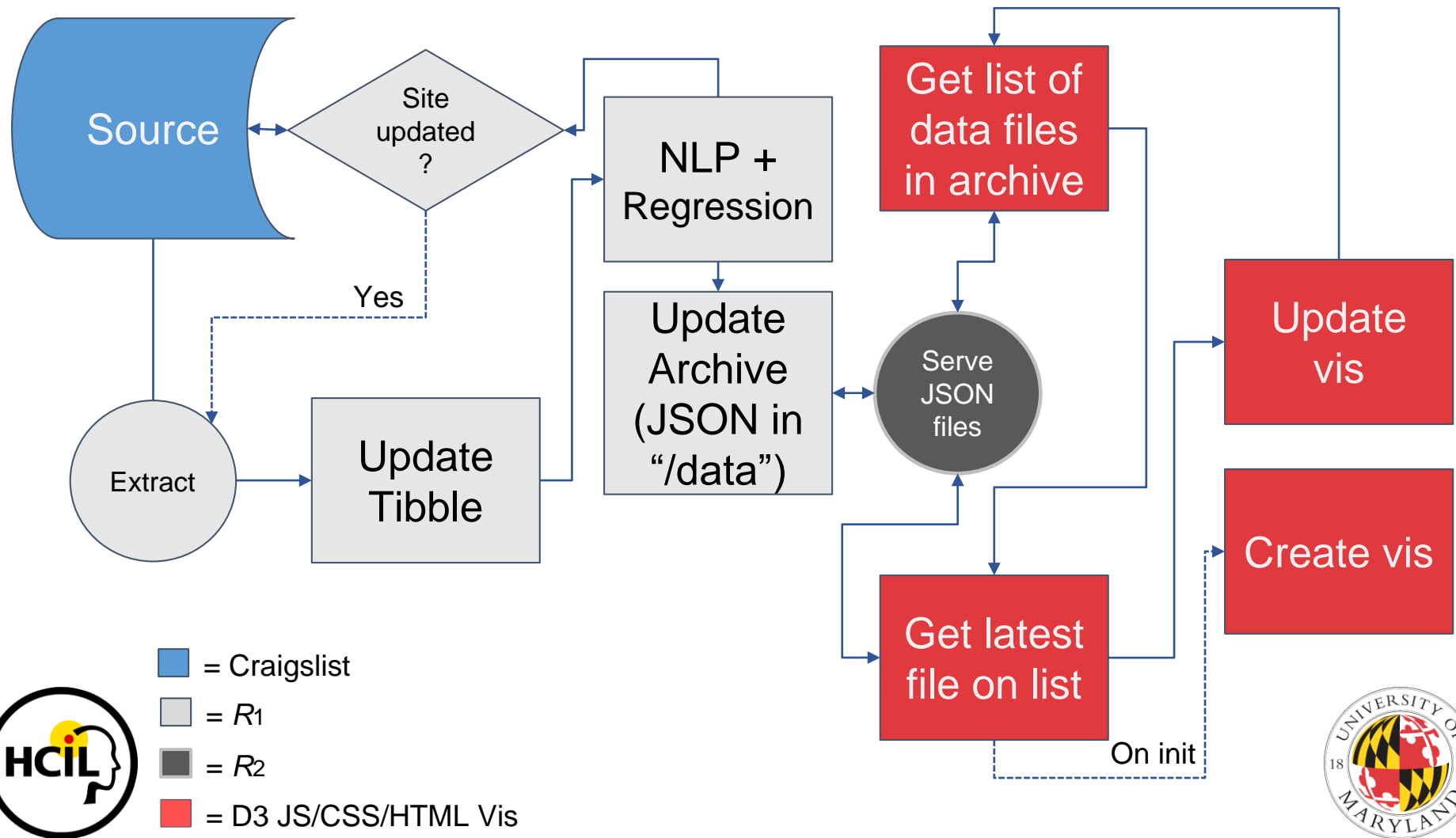- "Real-time" somewhat subjective, contextual

"Streaming visualization," then, is any vis that is continuously updated based on newly-generated, high frequency data

# Our process so far & next step

Site → Extract → Transform → Model → Publish [Stage]

# Our finished network of continuous processes

# In-Class Exercise 2

**Modify index.html to reflect your changes (Note: I don't necessarily expect you to finish this now)**

# Data science workflow

1. Scope out the problem or question
2. Knowledge search: Research and sensemaking
3. Data retrieval; extract, transform, load (ETL)
4. Exploratory analysis
5. Modeling

- **System-building [sometimes]**
- **Versioning/archival**
- **Communication**

CLOSING REMARKS

# QUEST IONS: