

# Data: Types & Roles

INFX 551

Foundation of Data Curation

# Agenda

- Data - Type vs Role distinctions; Structured vs Unstructured Data
- Curation - The premise of data curation; Workflows for Curation

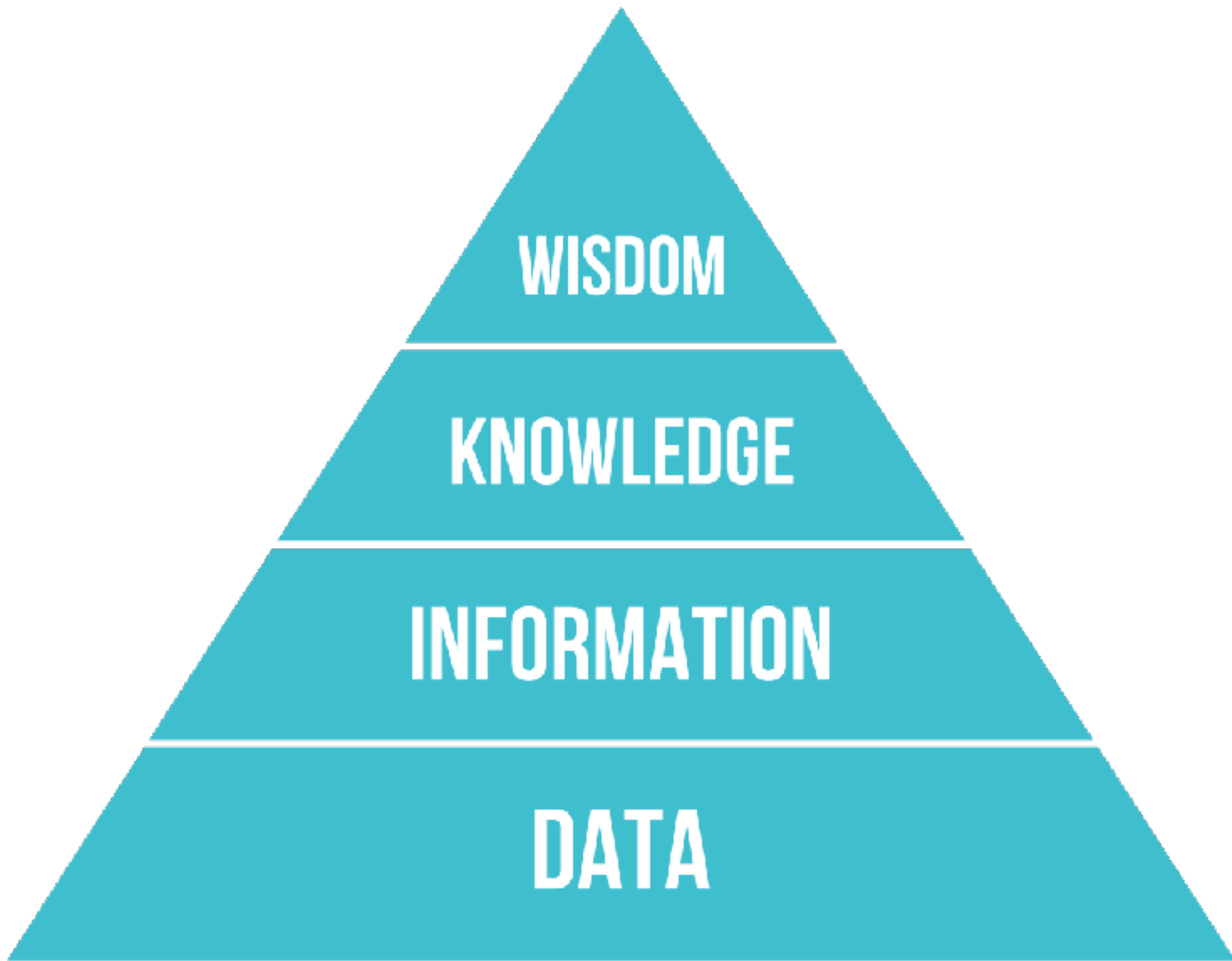
## Course Outline

This course consists of three modules.  
Each module consists of three topics.

Data	Data Systems	Policy, Privacy, & Ethics
<u>Types &amp; Roles</u>	Repositories	Policy
<u>Documentation</u>	Preservation	Privacy
<u>Standardization</u>	Cost Models	Ethics

For this course, let's assume that....

Data Curation is the active and ongoing management of **data** throughout a lifecycle of use, including reuse unanticipated contexts.



## **Research Data**

“The data, records, files or other evidence, irrespective of their content or form (e.g. in print, digital, physical or other forms), that comprise research observations, findings or outcomes, including primary materials and analysed data.”



# Open Data

“Open data is data that can be freely used, shared and built-on by anyone, anywhere, for any purpose.”



OPEN KNOWLEDGE

For this course, let's assume that....

Data are various **types** of information objects playing the **role** of evidence.

Evidence of...

Infection (Public Health)

Patterns of Consumer Behavior (Business)

Weather (Atmosphere)



Types vs Roles

Structured vs Unstructured

# Type vs Role distinctions

Type:  
Donald Trump is a person.



Role:  
Donald Trump is POTUS



# Types of Data

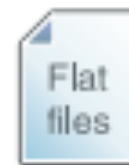
(by file format)



**XML**



**Databases**



**Flat Files**



**EDI**



**Excel**



**XBRL**



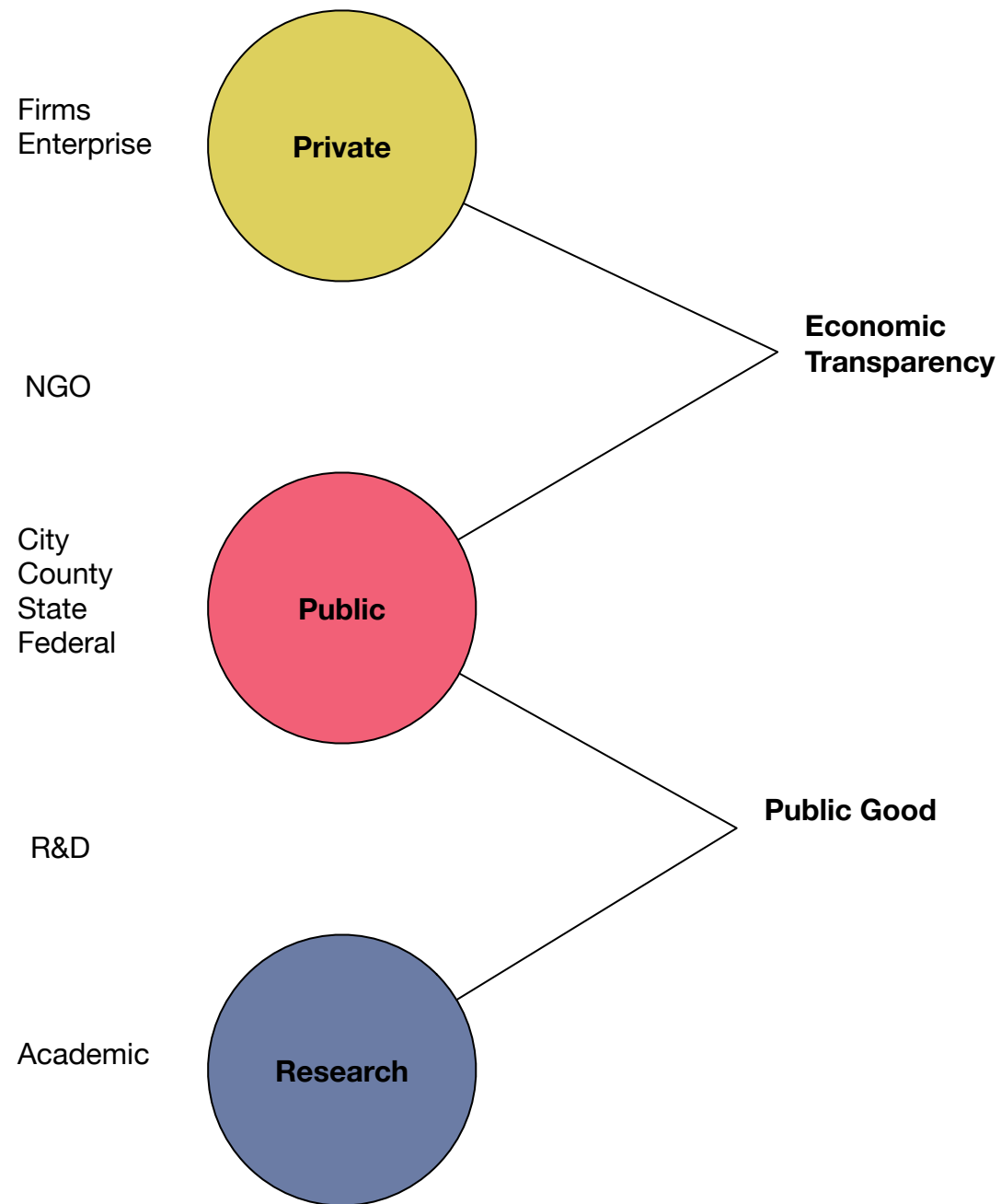
**JSON**



**Web Services**

# Types of Data

(by sector)



# Roles



# GBIF

## Global Biodiversity Information Facility

# Roles

```
<text xmlns="http://www.tei-c.org/ns/1.0" xml:id="d1">
  <body xml:id="d2">
    <div1 type="book" xml:id="d3">
      <head>Songs of Innocence</head>
      <pb n="4"/>
      <div2 type="poem" xml:id="d4">
        <head>Introduction</head>
        <lg type="stanza">
          <l>Piping down the valleys wild, </l>
          <l>Piping songs of pleasant glee, </l>
          <l>On a cloud I saw a child, </l>
          <l>And he laughing said to me: </l>
        </lg>
        <lg type="stanza">
          <l>"Pipe a song about a Lamb!" </l>
          <l>So I piped with merry chear. </l>
          <l>"Piper, pipe that song again;" </l>
          <l>So I piped, he wept to hear. </l>
        </lg>
        <lg type="stanza">
          <l>"Drop thy pipe, thy happy pipe; </l>
          <l>Sing thy songs of happy chear;" </l>
          <l>So I sung the same again, </l>
          <l>While he wept with joy to hear. </l>
        </lg>
        <lg type="stanza">
          <l>"Piper, sit thee down and write </l>
          <l>In a book, that all may read." </l>
          <l>So he vanis'd from my sight, </l>
          <l>And I pluck'd a hollow reed, </l>
        </lg>
        <lg type="stanza">
          <l>And I made a rural pen, </l>
          <l>And I stain'd the water clear, </l>
          <l>And I wrote my happy songs </l>
          <l>Every child may joy to hear. </l>
        </lg>
      </div2>
    </div1>
  </body>
</text>
```

# Roles



For this course, let's assume that....

Data are various **types** of information objects playing the **role** of evidence.

Evidence of...

Infection (Public Health)

Patterns of Consumer Behavior (Business)

Weather (Atmosphere)



**Unstructured**

vs

**Structured**

# Unstructured

Early Lead

## Boston Marathon 2017: Geoffrey Kirui tops American Galen Rupp and Edna Kiplagat cruises

By Cindy Beren April 17, 2017



Geoffrey Kirui, of Kenya, crosses the finish line. (Elise Amendola/Associated Press)

Geoffrey Kirui pulled away from Galen Rupp over the last few miles, winning the men's division of the 121st Boston Marathon on Monday and ending Rupp's bid to become the first American winner since 2014. In the women's division, Edna Kiplagat made the most of her first appearance in the Boston Marathon, separating from the pack at the 18-mile mark and cruising to a victory.

Kirui, a 24-year-old runner from Kenya, won in 2:09:37 on a warm day in Boston. Rupp, the Olympic bronze medalist in the event, finished second in his first big-city American marathon in 2:09:58 and Suguru Osako of Japan was third in 2:10:28. In his three marathons, Rupp has finished first (in the Olympic trials), third and second. "It lived up to and exceeded all my expectations," Rupp told NBC. "[The crowd] really lifted me those last three miles."

# Structured

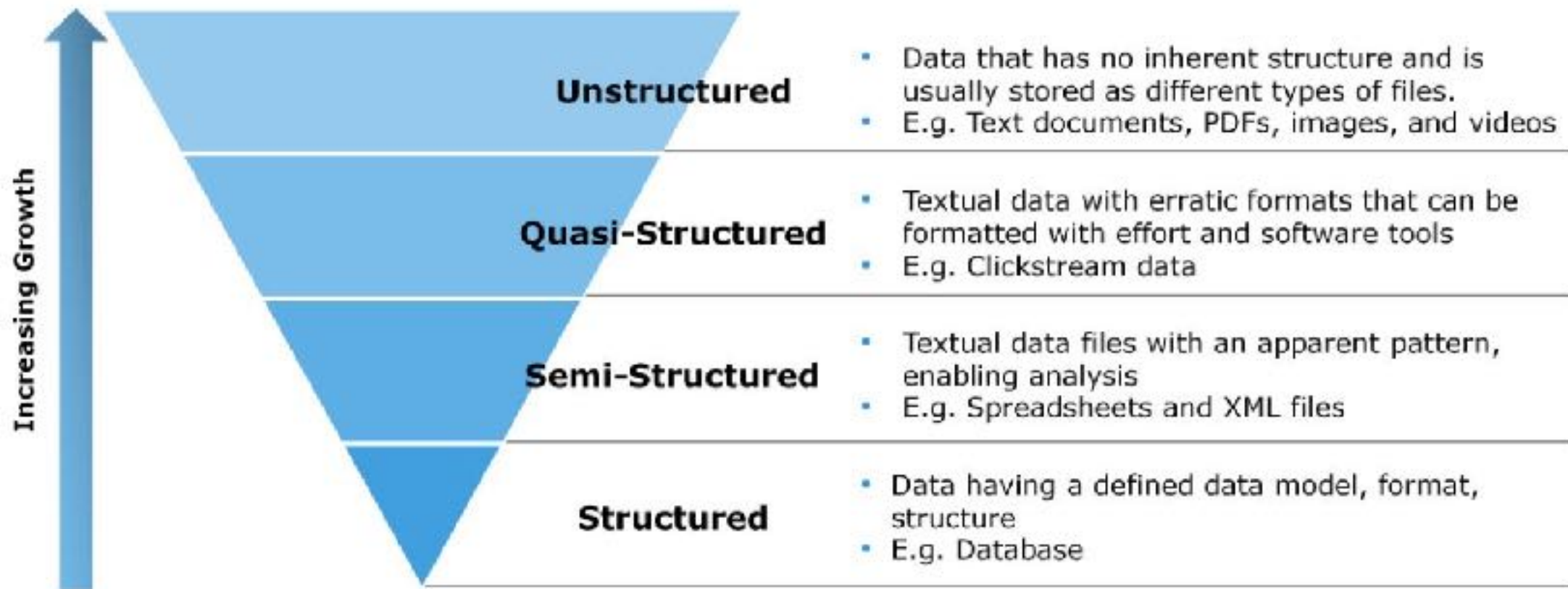
Name	Time	Nationality	Place
Geoffrey Kirui	2:09:37	Kenya	1
Galen Rupp	2:09:58	USA	2
Suguru Osako	2:10:28	Japan	3

# Unstructured

# Structured

?

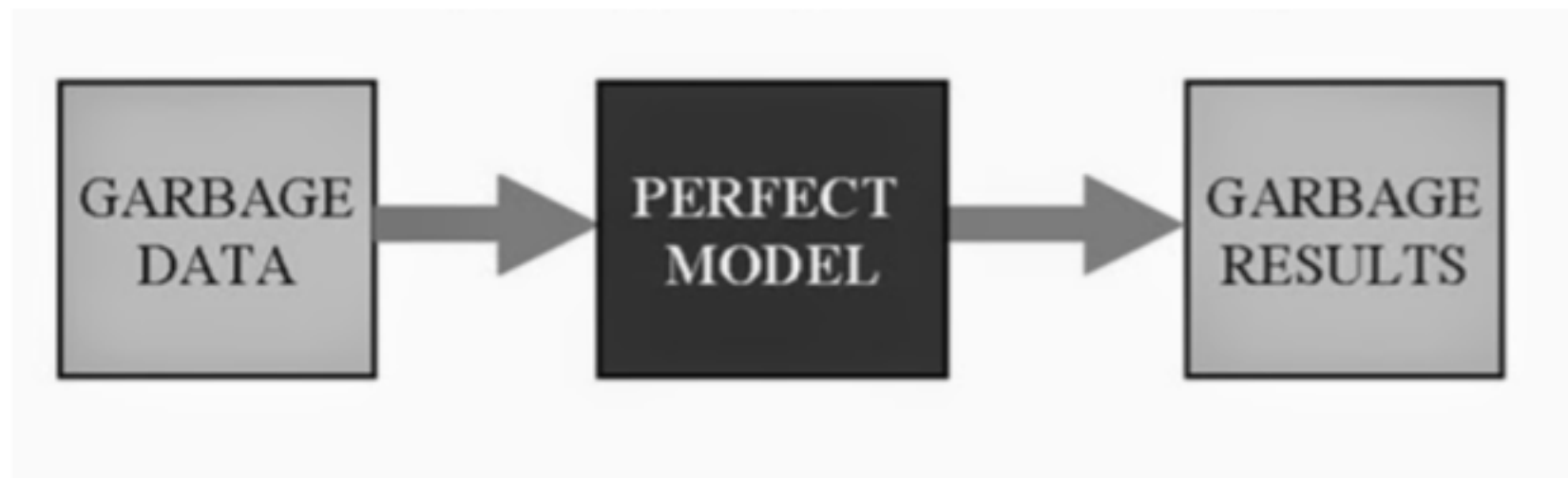
```
<text xmlns="http://www.tei-c.org/ns/1.0" xml:id="d1">
  <body xml:id="d2">
    <div1 type="book" xml:id="d3">
      <head>Songs of Innocence</head>
      <pb n="4"/>
      <div2 type="poem" xml:id="d4">
        <head>Introduction</head>
        <lg type="stanza">
          <l>Piping down the valleys wild, </l>
          <l>Piping songs of pleasant glee, </l>
          <l>On a cloud I saw a child, </l>
          <l>And he laughing said to me: </l>
        </lg>
        <lg type="stanza">
          <l>"Pipe a song about a Lamb!" </l>
          <l>So I piped with merry chear. </l>
          <l>"Piper, pipe that song again;" </l>
          <l>So I piped, he wept to hear. </l>
        </lg>
        <lg type="stanza">
          <l>"Drop thy pipe, thy happy pipe; </l>
          <l>Sing thy songs of happy chear;" </l>
          <l>So I sung the same again, </l>
          <l>While he wept with joy to hear. </l>
        </lg>
        <lg type="stanza">
          <l>"Piper, sit thee down and write </l>
          <l>In a book, that all may read." </l>
          <l>So he vanis'd from my sight, </l>
          <l>And I pluck'd a hollow reed, </l>
        </lg>
        <lg type="stanza">
          <l>And I made a rural pen, </l>
          <l>And I stain'd the water clear, </l>
          <l>And I wrote my happy songs </l>
          <l>Every child may joy to hear. </l>
        </lg>
      </div2>
    </div1>
  </body>
</text>
```



# Curation

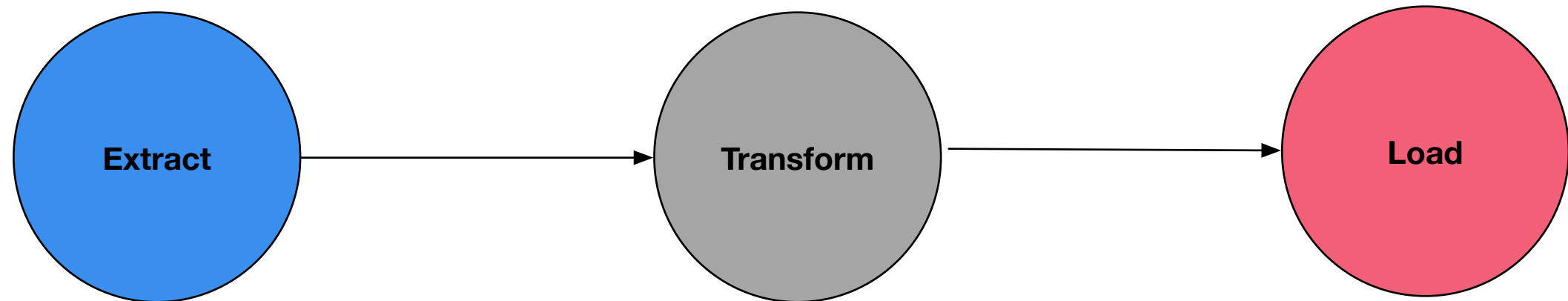
# Data Curation

Old computer science saying "Garbage in = Garbage out"



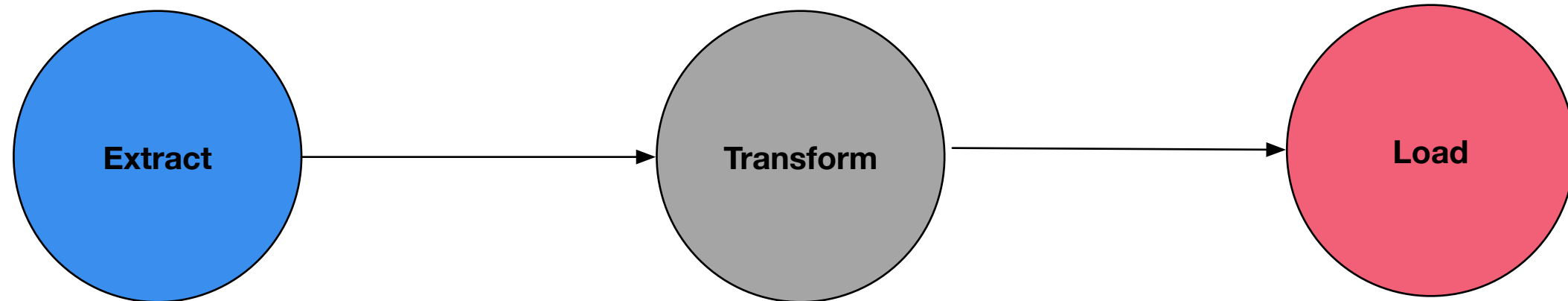
Data curation says "Quality in = Quality out"

# ETL: Simplest form of Curation Workflow





# ETL: Simplest form of Curation Workflow



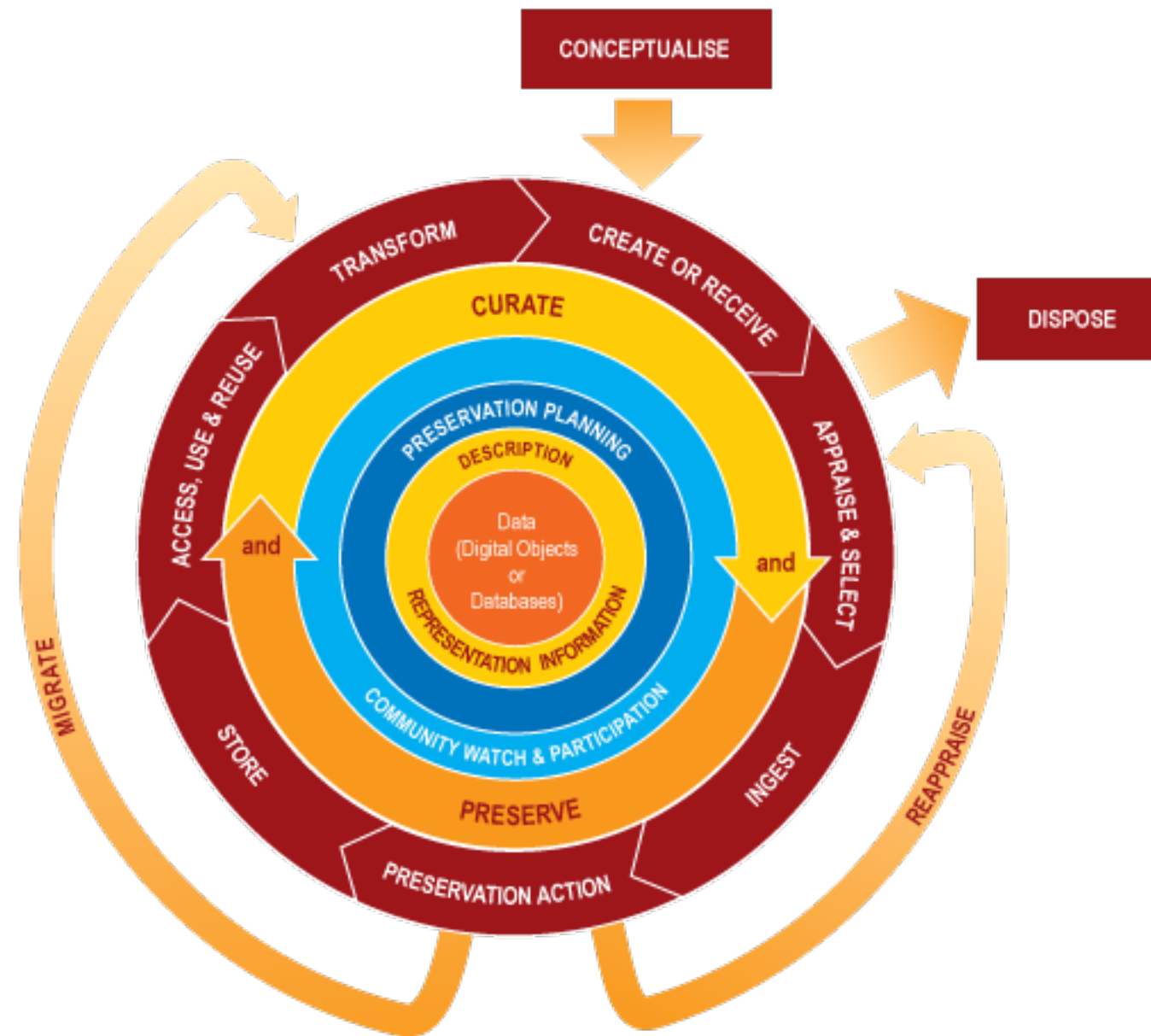
```
<text xmlns="http://www.tei-c.org/ns/1.0" xml:id="d1">
  <body xml:id="d2">
    <div1 type="book" xml:id="d3">
      <head>Songs of Innocence</head>
      <pb n="4"/>
      <div2 type="poem" xml:id="d4">
        <head>Introduction</head>
        <lg type="stanza">
          <l>Piping down the valleys wild, </l>
          <l>Piping songs of pleasant glee, </l>
          <l>On a cloud I saw a child, </l>
          <l>And he laughing said to me: </l>
        </lg>
        <lg type="stanza">
          <l>"Pipe a song about a Lamb!" </l>
          <l>So I piped with merry cheer. </l>
          <l>"Piper, pipe that song again;" </l>
          <l>So I piped, he wept to hear. </l>
        </lg>
        <lg type="stanza">
          <l>"Drop thy pipe, thy happy pipe;" </l>
          <l>Sing thy songs of happy cheer;" </l>
          <l>So I sung the same again, </l>
          <l>While he wept with joy to hear. </l>
        </lg>
        <lg type="stanza">
          <l>"Piper, sit thee down and write </l>
          <l>In a book, that all may read." </l>
          <l>So he vanish'd from my sight, </l>
          <l>And I pluck'd a hollow reed, </l>
        </lg>
        <lg type="stanza">
          <l>And I made a rural pen, </l>
          <l>And I stain'd the water clear, </l>
          <l>And I wrote my happy songs </l>
          <l>Every child may joy to hear. </l>
        </lg>
      </div2>
    </div1>
  </body>
</text>
```





# Data Curation Lifecycle

(more complex workflow)



# Data Curation Lifecycle

(more complex workflow)

Upstream

Downstream

