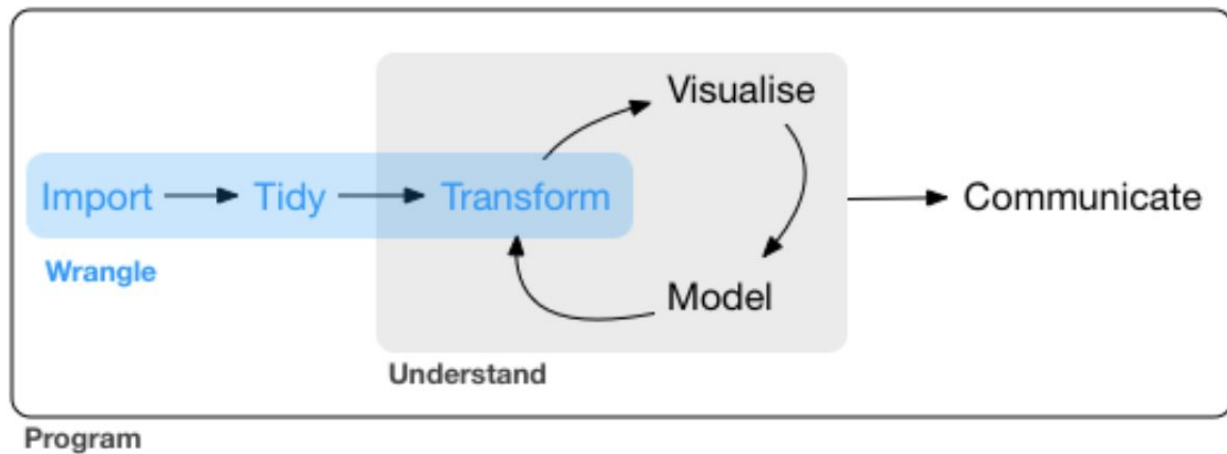# 1.2 Import

Applied Data Analysis (ADA)

Oxford DH Summer School - 2023
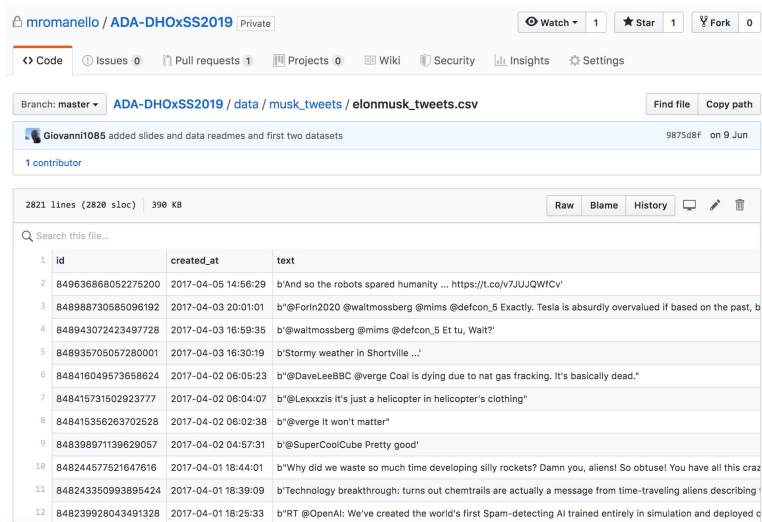
**Data formats**

In any data analysis project, a substantial chunk of time goes into preparing your data for analysis. This includes reading (legacy) data formats, storing data to intermediates files, saving data to a database system for longer-term storage.

We will focus on the following **key formats**:
- CSV
- XML
- JSON

**CSV**

- *CSV*: comma separated values

- *TSV*: tab separated values

- The *lingua franca* of tabular data

- First row: column headers

Some **drawbacks**:

- Characters in noisy text columns may interfere with separators (use TSV)

- Data type information not stored in the file (unlike e.g. Parquet)

**XML**

e**X**tensible **M**arkup **L**anguage

*Descriptive markup*: focus on **content** (data) rather than its **presentation**

Extensible: there is no pre-defined tagset

XML documents:
- must be well-formed (document has one root; all elements closing tag, etc.)
- must validate against a schema (structure first!)

# XML

## Elements

(1) `<l>The sun does arise, </l>`

(2) `<br />`

(3) `<lg>`

   `<l>The sun does arise, </l>`

   `</lg>`

## Attributes

- **Element**
- **Attribute**
- **attribute value**

`<lg type="stanza">`

```
<div2 type="poem" xml:id="d6">
    <head>The Ecchoing Green</head>
    <lg type="stanza">
        <l>The Sun does arise, </l>
        <l>And make happy the skies; </l>
        <l>The merry bells ring </l>
        <l>To welcome the Spring; </l>
        <l>The sky-lark and thrush, </l>
        <l>The birds of the bush, </l>
        <l>Sing louder around </l>
        <l>To the bells' chearful sound, </l>
        <l>While our sports shall be seen </l>
        <l>On the Ecchoing Green. </l>
    </lg>
    <lg type="stanza">
        <l>Old John, with white hair, </l>
        <l>Does laugh away care, </l>
```

Mixing vocabularies:

Declaration of namespaces in the root element

`<root xmlns:tei="http://www.tei-c.org/ns/1.0">`

`<tei:lg>`

**XML**

**Technologies:**

- *Presentation* → CSS

- *Transformation* → XSLT

- *Navigation* → XPath

- *Query* → XQuery

**XML in the wild:**

- TEI (Text Encoding Initiative) XML

- MARCXML (library catalogue data)

- RDF XML

- ALTO XML (OCR data)

- GraphML

- ...

**JSON**

**J**ava**S**cript **O**bject **N**otation

```
{
  '_id' : 1,
  'name' : { 'first' : 'John', 'last' : 'Backus' },
  'contribs' : [ 'Fortran', 'ALGOL', 'Backus-Naur Form', 'FP' ],
  'awards' : [
    {
      'award' : 'W.W. McDowell Award',
      'year' : 1967,
      'by' : 'IEEE Computer Society'
    }, {
      'award' : 'Draper Prize',
      'year' : 1993,
      'by' : 'National Academy of Engineering'
    }
  ]
}
```

https://www.mongodb.com/json-and-bson

**Data types**:

- Number

- String: any sequence of zero or more Unicode characters

- Boolean: true | false

- Array (list in Python)

- Object: unordered collection of name-value pairs {"title": "ADA DHOxSS"}

**JSON**

- Need for structure: JSON schema

- Databases that *speak* JSON: ElasticSearch, MongoDB, etc.

- JSON-LD: Web-friendly format to encoded RDF data

**Working with data formats**

(In ADA most of the times we don't get to choose the format of data we work with.)

But if you are the one to choose, consider:

- Target use (analysis, online presentation, etc.)
- Target community (wider public, scholars, data scientists)
- Multiple formats for multiple usage scenarios (internal usage, data publication, web application, etc.)