

On TAP:

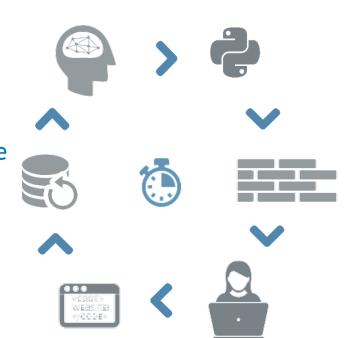
Module 1: Data Ingest

Kyle H. Ambert, PhD Intel Big Data Solutions, Datacenter Group



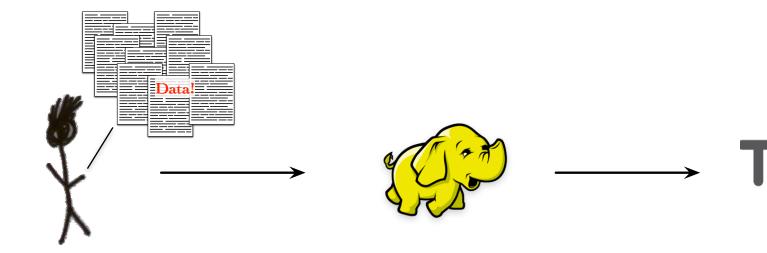
Ingest:

- Getting data into TAP
- Getting data into atk
- Moving from raw to the beginning of structure
- File types:
 - CSV
 - xml
 - json

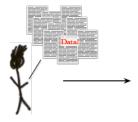




Data on TAP

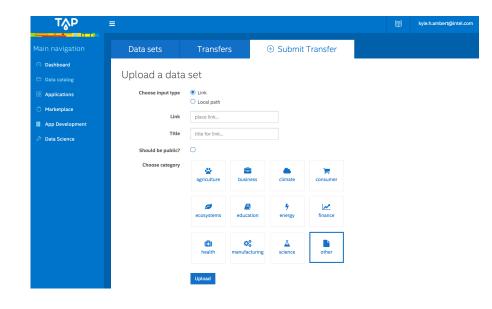




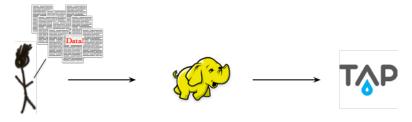




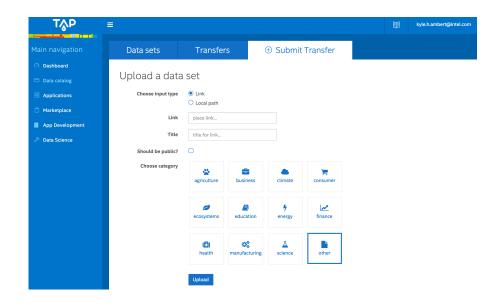




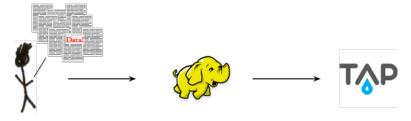




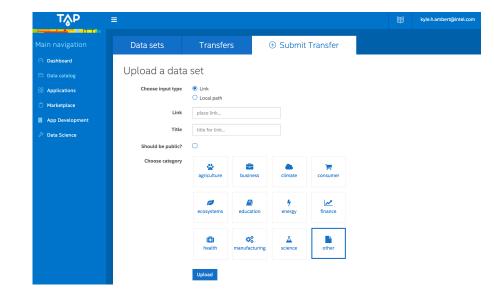
Getting data into the Toolkit



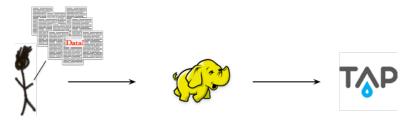




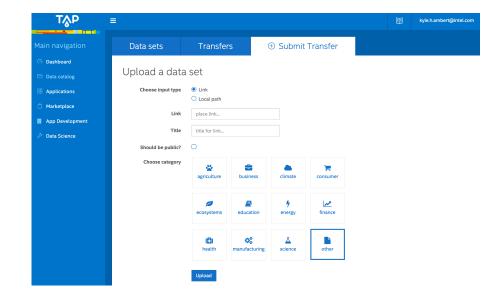
Use the console front-end



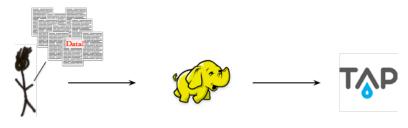




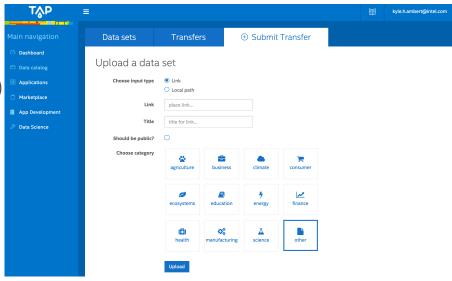
- Use the console front-end
 - Data Catalog page







- Use the console front-end
 - Data Catalog page
- Submit transfer from local path (or remote url)



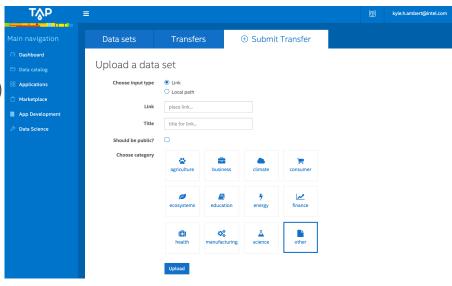








- Use the console front-end
 - Data Catalog page
- Submit transfer from local path (or remote url)
- Upload!



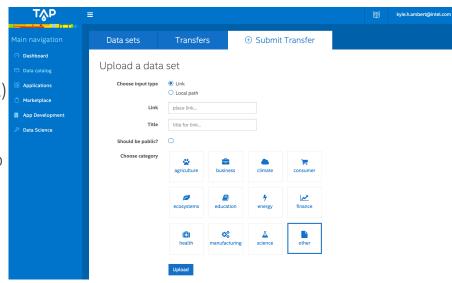








- Use the console front-end
 - Data Catalog page
- Submit transfer from local path (or remote url)
- Upload!
- Data can be previewed from the Data Sets tab





- [1] Define the **Data**
 - Iterative error-prone drudgery
 - One-off, ad hoc models in isolation
- [2] Define the Schema
 - A list of tuples: ('column_name', data_type)
- [3] Use the ia.UploadRows function



[1] Define the **Data**

- Iterative error-prone drudgery
- One-off, ad hoc models in isolation

[2] Define the Schema

- A list of tuples: ('column_name', data_type)
- [3] Use the ia.UploadRows function





[1] Define the **Data**

- Iterative error-prone drudgery
- One-off, ad hoc models in isolation

[2] Define the Schema

- A list of tuples: ('column_name', data_type)
- [3] Use the ia. UploadRows function

Example:

```
animals = ia.Frame(ia.UploadRows([['puppy', 'Harrison'], ['cat', 'wahlberg']], [('animal', str), ('name', str)]))
```



[1] Define the **Data**

- Iterative error-prone drudgery
- One-off, ad hoc models in isolation

[2] Define the Schema

- A list of tuples: ('column_name', data_type)
- [3] Use the ia. UploadRows function



Example:

```
animals = ia.Frame(ia.UploadRows([['puppy', 'Harrison'], ['cat', 'wahlberg']], [('animal', str), ('name', str)]))
```



All columns must be the same length! All data types must be **serializable**!



Character-separated Value Files (.csv)

[1] Define the schema

- Iterative error-prone drudgery
- One-off, ad hoc models in isolation

[2] ta.CsvFile

Single-threaded, single-node processing

```
__init__(self, file_name, schema, delimiter=',', skip_header_lines=0)
```

Example:

```
import trustedanalytics as ta
ta.CsvFile("raw_data.csv", schema=[("col1", ta.int32), ("col2", ta.float32)])
```

Character-separated Value Files (.csv)

[1] Define the schema

- Iterative error-prone drudgery
- One-off, ad hoc models in isolation

[2] ta.CsvFile

• Single-threaded, single-node processing

```
__init__(self, file_name, schema, delimiter=',', skip_header_lines=0)
```

Example:

```
import trustedanalytics as ta
ta.CsvFile("raw_data.csv", schema=[("col1", ta.int32), ("col2", ta.float32)])
```



Be sure you have the correct delimiter!

Parallel Ingest:

 Just point the toolkit to a directory containing multiple xml files to process them all at once.

Parallel Ingest:

 Just point the toolkit to a directory containing multiple xml files to process them all at once.

Identify the tag denoting separation of entities:

Parallel Ingest:

• Just point the toolkit to a directory containing multiple *xml* files to process them all at once.

Identify the tag denoting separation of entities:

```
<note>
    <to>Self</to>
    <from>Self</from>
    <heading>Eagles Dig</heading>
    <body>Chip Kelly was better at UO</body>
</note>
```

Parallel Ingest:

 Just point the toolkit to a directory containing multiple xml files to process them all at once.

Identify the tag denoting separation of entities:

Parallel Ingest:

 Just point the toolkit to a directory containing multiple xml files to process them all at once.

Identify the tag denoting separation of entities:



Schema-naïve processing is currently a WIP. It will make everyone's life easier!



```
Wrapping this sort of processing code in a
                                function will be useful when working with
                                complex xml files!
         Just point
         multiple xi
                                   def parse_xml_to_frame(path, tag, name):
                                      Helper function to convert an xml file on the hdfs into a data frame...
Identify the tag
                                      xml = ia.XmlFile(path, tag)
                                      # Check that the frame doesn't already exist. Drop it, if it does...
         <note>
                                      if name in ia.get_frame_names():
             <to>Self</
                                         sys.stderr.write("Dropping existing frame named {NAME}...\n".format(NAME=name))
                                          ia.drop frames(name)
             <from>Self
                                      frame = ia.Frame(xml. name=name)
                                      return frame
             <heading>E
             <body>Chir
         </note>
Example:
```

```
xml = ia.XmlFile(file_name='path', tag_name='tag')
frame = ia.Frame(source=xml, name='name')
```





Schema-naïve processing is currently a WIP. It will make everyone's life easier!



JSON Files

Parallel Ingest

- Just point the toolkit to a directory containing u'ADMITTING_DIAG_CODEs': [u'786.05'], multiple json files to process them all at once. u'ADMIT_REASONS': [u'COPD EXACERBATION', u'SICK098]
- Creating a frame from json will result in a frame with a single column: 'data_lines'

ia.JsonFile:

```
__init__(self, file_name)
```

Example:

```
json_file = ta.JsonFile("data/raw_data.json")
my_frame = ta.Frame(json_file)
```

```
u'ADMITTING_CCS_LVL_2_LABELs': [u'Other lower respiratory disease [133.]'],
u'ADMITTING_CCS_LVL_4_LABELs': [u' '],
u'ADMITTING_DIAG_CCS_LVL_4_LABELs': [u'nan', u' '],
u'ADMIT_REASONs': [u'COPD EXACERBATION', u'SICK0985'],
u'ADMIT_SOURCE_MEDVIEW': u'EMERGENCY OP UNIT',
u'ARITHMETIC_LOS': 4.2,
                      u'Diabetes mellitus without complication [49.]',
                      u'Other and unspecified lower respiratory disease'l
```



For large, nested json, data_lines can get quite large, causing some operations to take time



— code! —

Problems:

- [1] Create a 3x1 frame with the following contents:
 - (c1) "animals": five species of animal
- [2] Ingest a new character-separated (csv) file (offsides.tsv)
- [3] Create a 100 x 2 frame with the following contents:
 - (c1) "rn": 100 random integers between 0-1000
 - (c2) "evenodd": "EVEN" if the corresponding integer is even, otherwise "ODD"