Leveraging HDF5 infrastructure by ADF to build an interoperable package & contextualized data in Pharma using semantic technology

Amnon Ptashek Allotrope Foundation





- The *RDF* Resource Description Framework is a framework for expressing knowledge about resources.
- Resources can be anything: devices, people, physical objects, and abstract concepts
- We can visualize knowledge as a connected graph, consisting of nodes and arcs



- *RDF* statement expresses a relationship between two resources
- A statement always has the following structure called triples:
 <subject> <predicate> <object>
- The *subject* and the *object* represent the two resources being related; the *predicate* represents the nature of their relationship.
- The relationship is phrased in a directional way from *subject* to *object* and is called in *RDF* a *property*

* Copyright © 2015 W3C® (MIT, ERCIM, Keio, Beihang). This software or document includes material copied from or derived from https://www.w3.org/TR/rdf11-prime





 This ability to have the same resource be in the *object* position of one triple and the *subject* position of another triple makes it possible to find connections between triples, which makes *RDF* very powerful

<Bob> <is a> <person>.
<Bob> <is a friend of> <Alice>.
Person
4
8ob> <is born on> <14th of July 1990>.
<Bob> <is interested in> <the Mona Lisa>.
<the Mona Lisa> <was created by> <Leonardo da Vinci>.
<the video 'La Joconde à Washington'> <is about> <the Mona Lisa>



- IRI International Resource Identifier are used to identify resources uniquely
- **Objects** can be **literal**s that are basic values and not **IRI**s
- IRIs are typically used in combination with RDF vocabularies that provide semantic context about these resources
- There are universal *RDF vocabularies* such as ("Friend of a Friend" (FOAF), Person vocabulary for describing social networks
 There are universal *RDF vocabularies* "1990-07-0"





ardo Da Vi

- Use SPARQL (Protocol and RDF Query Language) to query RDF graph e.g. "people interested in the Mona Lisa"
- Query for multiple triple patterns

PREFIX rdf: <<u>http://www.w3.org/1999/02/22-rdf-syntax-ns#</u>> PREFIX foaf: <<u>http://xmlns.com/foaf/0.1/</u>> PREFIX wd: <<u>http://www.wikidata.org/entity</u>>

SELECT ? name WHERE { ?name rdf: type foaf: Person. ?name foaf: topic_interest wd:Q12418.







©2020 Allotrope Foundation

- ADF (Allotrope Data Format) File contains several components:
 - Data Description
 - Data Cubes
 - Data Package
 - Graph Store
 - Audit Trail (Audit/compliance features)
 - Checksums (Audit/compliance features)
- ADF File uses HDF5 as the underlying file format to store scientific data in a performant, consistent, and <u>self-describing way</u>



- Graph Store
 - Stores user defined graphs
 - Stores it in a form of *RDF* graph
 - Based on Apache Jena
- **RDF** Graph



ADF file		
	Semantic web Framework Graph Store	
HJF		

* Copyright © 2015 W3C® (MIT, ERCIM, Keio, Beihang). This software or document includes material copied from or derived from https://www.w3.org/TR/rdf11-primer.

- Graph Store
 - Stores user defined graphs
 - Stores it in a form of *RDF* graph
 - Based on Apache Jena
- Data Description may contain
 - O Contextual metadata about the Data Cube content
 - O Contextual metadata about the Data Package content
 - Specifically for a Lab environment:
 - Results description
 - Sample description
 - Process description
 - Method description
 - Instrument description
 - Etc.
- The Data Description utilizes semantic web technology for a rich metadata and contextual information vs. HDF5 attributes



ADF file semantic we Data Description ena Graph Store Data Cubes Data Package

- Data Cubes
 - Contains multi-dimensional arrays of data
 - Specifically for a Lab environment it may contain measurements, results of an experiment or process data
 - Data Cube context (metadata) is described in the Data Description





HDF Users Group Meeting - October 15, 2020

- Data Package
 - General purpose virtual file storage system
 - Can (optional) archive the same data in the Data Cubes but in its original native format and/or any other data
 - Data Package context (metadata) is described in the Data Description
 - Ensure data consistency and integrity of files during storage and transfer.





ADF API Architecture

Data Packago API	Data Cuba ADI	Data Description API (Apache Jena)	SS
	Data Cube API	Quad Store API	ntologie
Platform independent file format (HDF 5)			





ADF internal structure utilizes HDF5 structure:

- Each of the three ADF layers is a toplevel group in the HDF5 file
 - Data Description
 - Data Cubes
 - Data Description
- Auxiliary features are also stored as a top-level group in the HDF5 file
 - Named Graphs
 - Audit Trail





©2020 Allotrope Foundation

HDF Users Group Meeting - October 15, 2020

ADF structure – Data Package:

- ADF uses a top-level data package group to hold the file system layer
- Underneath, folders are represented as HDF Groups and files as HDF datasets
- Key metadata about individual files and folders are stored in the Data Description automatically by the ADF API



ADF structure – Data Package:



ADF file



ADF structure – Data Cubes:

- RDF Data Cubes consist of <u>measures</u>, <u>dimensions</u>, and <u>scales</u>
 - <u>Measures</u> are the observed values for a given dataset, e.g. intensity or m/z
 - <u>Dimensions</u> are the controlled variables e.g. time or wavelength
 - <u>Scales</u> are the values of the dimensions, which may be linear, nonlinear, or entirely arbitrary
- All three are referenced in the Data Description such that metadata can be linked to them





©2020 Allotrope Foundation



©2020 Allotrope Foundation

HDF Users Group Meeting - October 15, 2020

ADF structure – Data Description:

- HDF5 does not yet have native graph storage
- ADF emulates this through a combination of:
 - A dictionary
 - A list of quads (actually quints, including a deletion flag)
 - Multiple index listings for fast search











AFO: Allotrope Taxonomies and Ontologies

- To address the laboratory analytical processes:
 - The *AFT* Allotrope Foundation Taxonomy formalizes the hierarchical relationships of terms
 - The AFO Allotrope Foundation Ontology is an ontology suite that provides a standard vocabulary and semantic model representation



AFO: Terminology

• Entities within the vocabulary are uniquely identified with standard *IRI*s

device identifier

(http://purl.allotrope.org/ontologies/result#AFR_0002018) synonyms: instrument id, machine id, device id, equipment id, equipment identifier, instrument identifier, machine identifier

A device identifier is an identifier that identifies some device. [Allotrope]

equipment serial number

(http://purl.allotrope.org/ontologies/result#AFR_0001119) synonyms: instrument serial number, device serial number, instrument id, machine id, machine serial number

Equipment serial number is measurement metadata that identifies an equipment used in the measuring by its serial number. [Allotrope]

balance identifier

(http://purl.allotrope.org/ontologies/result#AFR_0001986) synonyms: balance id

A balance identifier is a device identifier that identifies some balance. [Allotrope]



©2020 Allotrope Foundation

AFO: Terminology, Taxonomy



AFO: Terminology, Taxonomy and Ontology

• An ontology is a logic model that captures a domain in a machine understandable way including entities, relations and logic that controls relations.

located in — RO:0001025 — http://purl.obolibrary.org/obo/RO_0001025

has part — BFO:0000051 — http://purl.obolibrary.org/obo/BFO_0000051

 Ontologies provide an unconstrained vocabulary we can use to describe SubClass Of things (instances) in our open world and give them a meaning (= what it is)



AFO: Terminology, Taxonomy and Ontology



©2020 Allotrope Foundation

ADM: Interoperable ADF.File

- A "Model" represents one or more "Use Case(s)"
- The "Model" uses RDF graph to describe the "Use Case(s)"







MODEL	
	Γ
	RDF



ADM: Interoperable ADF.File

 An Allotrope standardized "Model" for the pharmaceutical laboratory analytical processes is called ADM – Allotrope Data Model





	ADM	
		R D F
-		





©2020 Allotrope Foundation

ADF file

ADM: Interoperable ADF File

 Heterogenous vendor systems can seamlessly exchange *ADF* files (read and write) and process its *Data* that adheres ^{va} to an associated standardized *ADM*





ADM: Shape File

- Contains a set of conditions
- Uses SHACL syntax:
 Shapes Constraint Language
 - Example of <u>a single condition</u> (single shape) on the Equipment serial number:

S purl.allotrope.org/ontologies/pro X

has value

synonyms: value

+

(http://purl.allotrope.org/ontologies/property#AFX_0000690)

C O Not secure purl.allotrope.org/ontologies/property#AFX_0000690

The literal piece of information as part of some information entity object. [Allotrope]

- af-s:AFS_0000611
 - rdf:type sh:PropertyShape ;
 - sh:path <http://purl.allotrope.org/ontologies/property#AFX_0000690>;

sh:message "Equipment serial number MUST have exactly one value of type xsd:string."

- sh:datatype xsd:string ;
- sh:maxCount 1;
- sh:minCount 1;
- Evaluate conformance of the ADM Data instance with the set of conditions in associated Shapes (SHACL)

Serial Number

Interoperable ADF File

Conductivity Measurement ADF (.adf)



©2020 Allotrope Foundation

HDF Users Group Meeting - October 15, 2020

ADF file

Thanks for your attention!

Allotrope Foundation Product Team

Amnon Ptashek <u>amnon.ptashek@allotrope.org</u> Allotrope Foundation <u>www.allotrope.org</u> email: <u>more.info@allotrope.org</u>