# POSTER 23

**CMAC** FUTURE MANUFACTURING RESEARCH HUB

## Making Pharmaceutical Manufacturing Data Ready for AI

Tabbasum Naz[1]*, Blair Johnston[1], Murray Robertson[1], Antony Vassileiou[1], Sophie Bailes[2], Neil Dawson[3], Simone Zomer[4], Tiffany Lai[3], Rachel Findlay[5], Gavin Reynolds[2], Amy Robertson[2]

**tabbasum.naz@strath.ac.uk**

1 CMAC Future Manufacturing Hub, University of Strathclyde, Glasgow, UK, 2 AstraZeneca, Macclesfield, UK, 3 Pfizer, Sandwich, UK, 4 GlaxoSmithKline, Ware, UK, 5 Centre of Process Innovation, Sedgefield, UK

## Introduction

Large volumes of pharmaceutical manufacturing data have been generated in recent years. A lot of time and effort has been spent producing data but they are, for the most part, scattered, unstructured, not machine readable and in heterogeneous formats. The work presented here provides integrated management and access to these valuable datasets. The Digital Design Accelerator Platform (DDAP) Extract-Transform-Load (ETL) tool has been developed to derive maximum value from the data acquisition effort to date and to allow future data to be integrated easily. DDAP ETL with multiple components can be used for automatic extraction, transformation and loading of heterogeneous pharmaceutical manufacturing data from multiple instruments. It is a collaborative effort to digitalise and make data Findable, Accessible, Interoperable and Reusable (FAIR). It also provides an opportunity to explore semantic heterogeneity across partners for standardisation efforts and ontology development in the medicine manufacturing domain. DDAP ETL can help domain experts to reap the benefits of the digital age and extract more value from organised data. It provides a foundation for future analytics and data-driven projects across the sector. In future AI, predictive analysis, statistical analysis, data visualization, data mining and machine learning techniques can be applied on the extracted data.

## Methodology

### ETL Extractor

DDAP Extractor is responsible for extracting schema/data from different raw data sources i.e., Brookfield Powder Flow Tester (PFT), Freeman Technology FT4 Powder Rheometer, Dr Dietmar Schulze Ring Shear Tester (RST-XS) and Micromeritics GeoPyc. There is heterogeneity in data at format, schema and data level and data is not machine readable. Data is available in multiple formats i.e., MS Excel, text files etc. Format heterogeneity makes it difficult to integrate so our DDAP-Extractor component resolves the problem of format heterogeneity by accessing multiple data formats and convert the data to XML format.
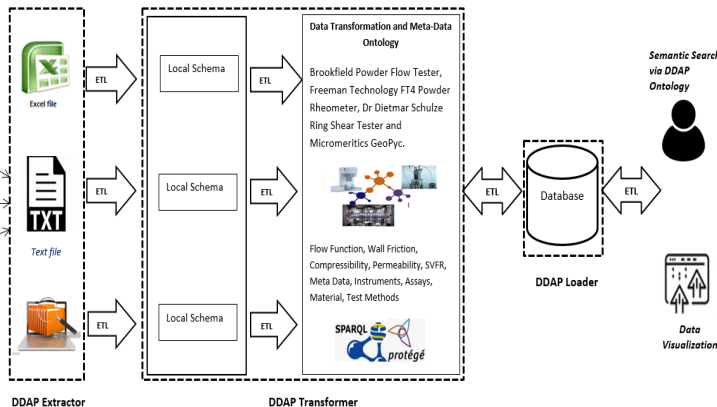
### ETL Transformer

With the help of DDAP Transformer, our key requirement is to provide automatic techniques for schema/data transformation. DDAP Transformer is responsible for resolving syntactic and semantic heterogeneity.
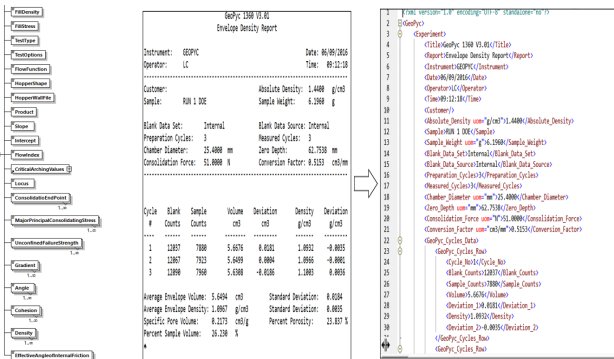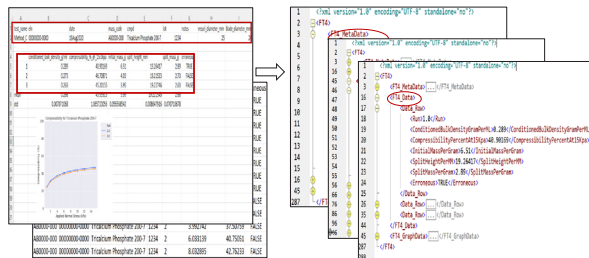
### ETL Loader

Once the extraction and transformation for schema and data is completed, the data is ready to send to the central repository. The DDAP-Loader component is responsible to load the data into DDAP repository.
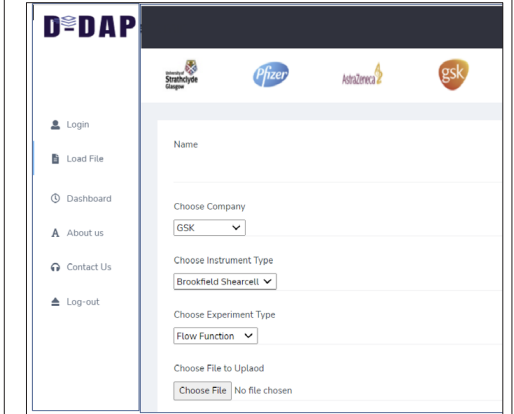
## Software Architecture
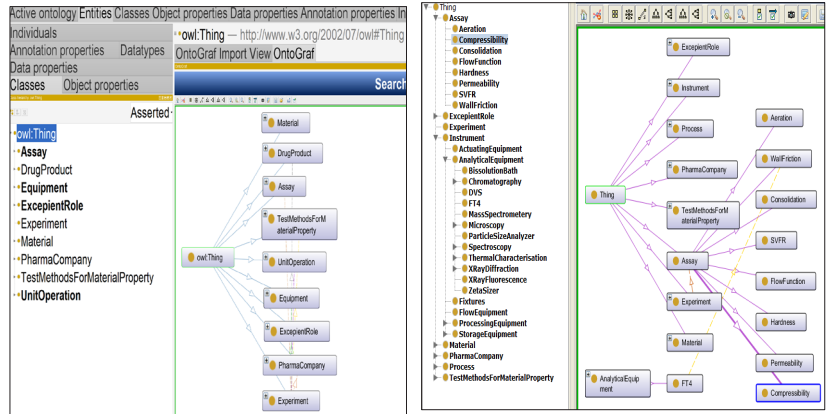


## Schema and Data Extraction Process

- Freeman Technology FT4 Powder Rheometer's Compressibility schema and data from MS Excel to xml by DDAP ETL.
- Schema for Flow Function from Brookfield Powder Flow Tester.
- Schema and data from Micromeritics GeoPyc text file to xml by DDAP ETL.



## DDAP ETL Interface



## DDAP Ontology for Pharmaceutical Manufacturing Domain



## Conclusion and Future Work

We have developed DDAP ETL (Extract, Transform and Load) – a tool that can extract pharmaceutical manufacturing data from different sources and develop a mechanism to translate between different concepts and data from multiple schemas. We have started to develop domain ontology in the pharmaceutical manufacturing domain as a way to represent meta-data and semantics. In future, data from DDAP ETL can be reused by experts in AI, predictive analysis, statistical analysis, data visualization, data mining and machine learning.