
Building a cross-disciplinary platform for smart meter data

Darren Bell

Director of Technical Services
UK Data Service

CLOSER Webinar
24 Feb 2021

UK Data Service



UKDA “Traditional” data infrastructure

- Approx. 8000 survey datasets outputted in SPSS/STATA/TAB with PDF documentation
- Data held in files on Windows folders (with different levels of security) as binary files
- Metadata held in MS-SQL databases
- [Discovery](#) (search and download) website in C# asp.net, based on Umbraco CMS
- Other data exploration tools like [Nesstar](#) (based on MySQL, JSP)
- Very much a classical architecture disseminating “pre-designed” datasets



DDI Codebook metadata

- Flat and simple. Easy to understand. Limited in sophistication
The “Study” is the core unit of management
- Supports discovery operations (this is the primary use case) but less expressive around variable relationships or rights management.

```
<codeBook xsi:schemaLocation="http://www.ddialliance.org/Specification/DDI-Codebook/2.5/XMLSchema/codebook.xsd">
  <docDscr>
    <citation>
      <titlStmt>
        <titl xml:lang="en">DDI2.5 XML CODEBOOK RECORD FOR STUDY NUMBER 7481</titl>
      </titlStmt>
      <prodStmt>
        <prodDate>2014-04-04T14:42:58Z</prodDate>
      </prodStmt>
    </citation>
  </docDscr>
  <studyDscr>
    <citation>
      <titlStmt>
        <titl xml:lang="en">Integrated Census Microdata (I-CeM), 1851-1911</titl>
        <subTitl></subTitl>
        <altTitl>I-CeM</altTitl>
        <IDNo agency="UKDA" xml:lang="en">7481</IDNo>
        <IDNo agency="datacite" xml:lang="en">10.5255/UKDA-SN-7481-2</IDNo>
      </titlStmt>
      <rspStmt>
        <AuthEnty xml:lang="en">Schurer, K., University of Essex, Department of History</AuthEnty>
        <AuthEnty xml:lang="en">Higgs, E., University of Essex, Department of History</AuthEnty>
      </rspStmt>
      <prodStmt>
        <copyright>Copyright K. Schürer, E.J. Higgs and BrightSolid</copyright>
        <fundAg>Economic and Social Research Council</fundAg>
        <grantNo>RES-062-23-1629</grantNo>
      </prodStmt>
      <distStmt>
        <distrbtr xml:lang="en">UK Data Service</distrbtr>
        <depositr>Schurer, K., University of Essex, UK Data Archive</depositr>
        <depDate date="2018-09-06T16:25:04Z">06 September 2018</depDate>
        <distDate date="2014-04-04T14:42:58Z">04 April 2014</distDate>
      </distStmt>
    </citation>
  </studyDscr>
</codeBook>
```

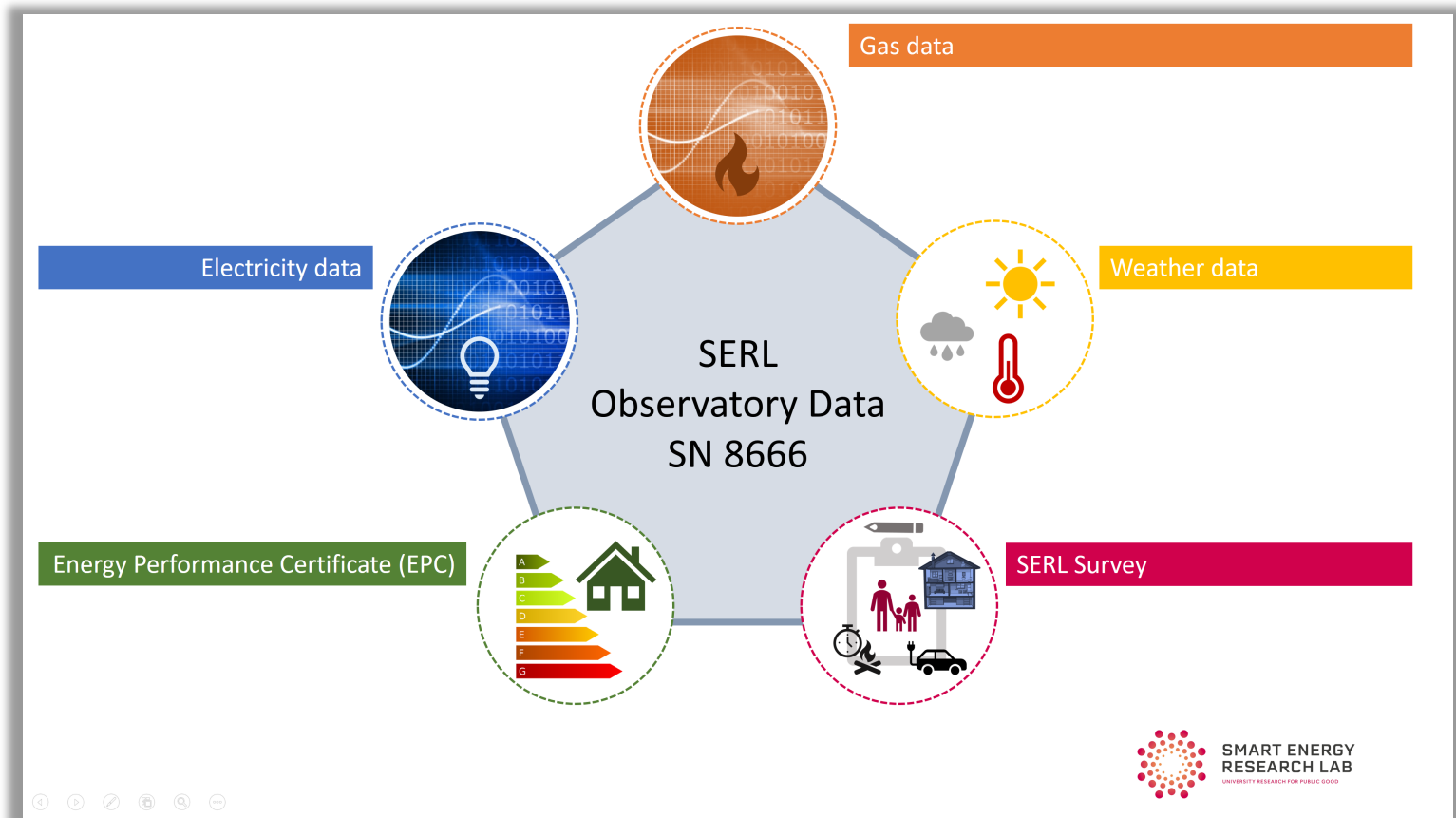


Smart Energy Research Lab (SERL) Project

- Smart Energy Research Lab (<https://www.serl.ac.uk>) is a data resource for UK research community
- High-quality smart meter and linked contextual data for innovative research
- Still very difficult for researchers to access high quality energy data in the UK
- However substantial barriers to accessing smart meter data
 - Technical
 - Legal
 - Financial



Cross-disciplinary by design

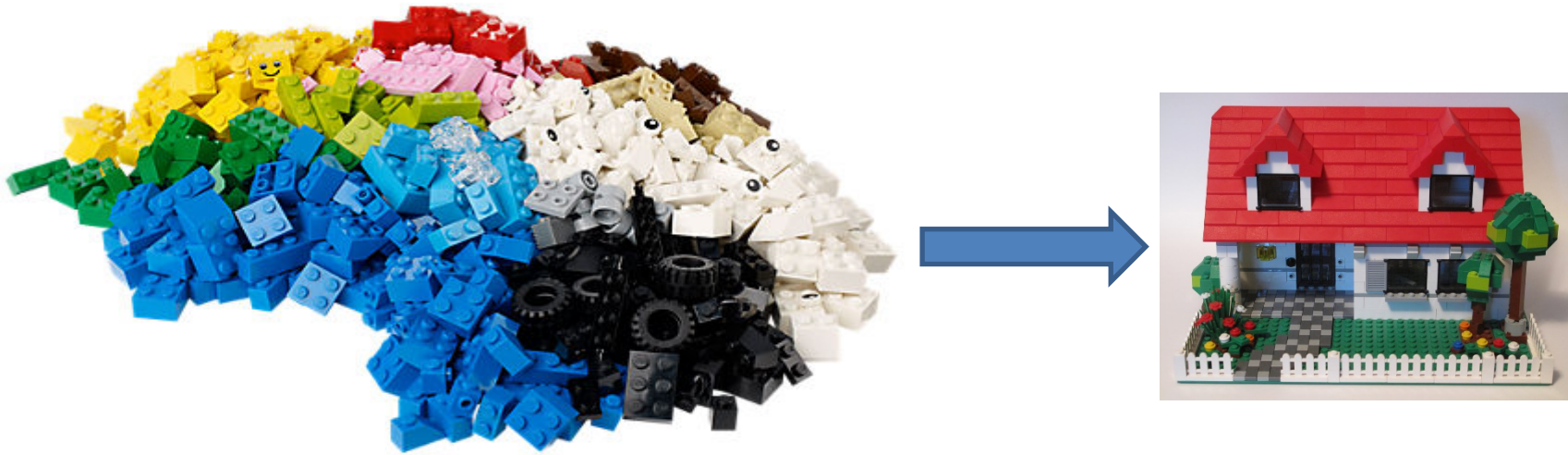


We need to move from this...



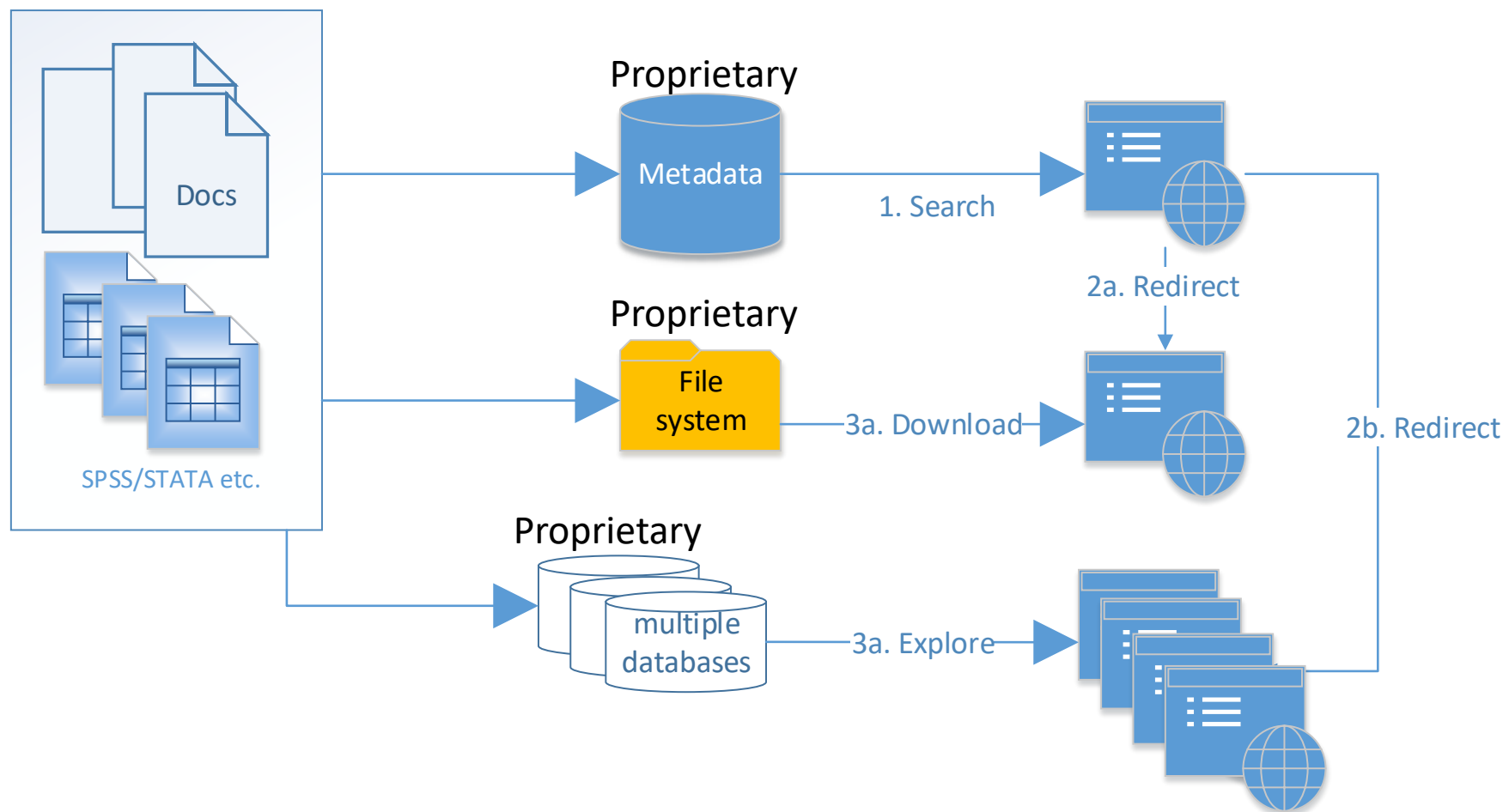
Pick pre-built datasets from the catalogue

To this...

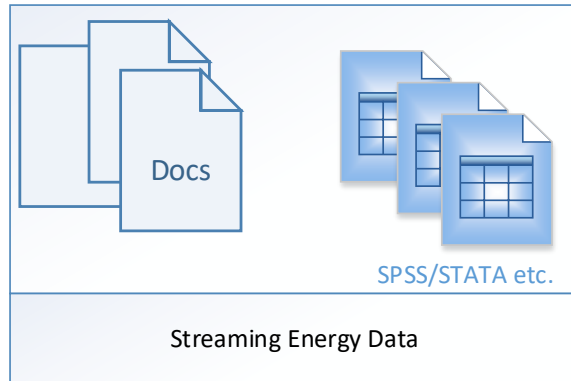


Build your own

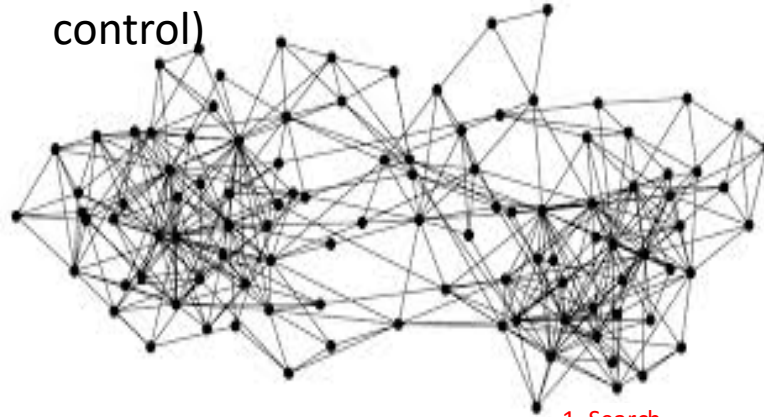
In more detail – from this...



In more detail – to this...



DDI-3/4 (data/metadata/access-control)



1. Search
2. Select (and link/aggregate)
- 3a. Simple viz
- 3b. Generate bespoke data product



Drill down to virtual lab
Researchers tools of choice

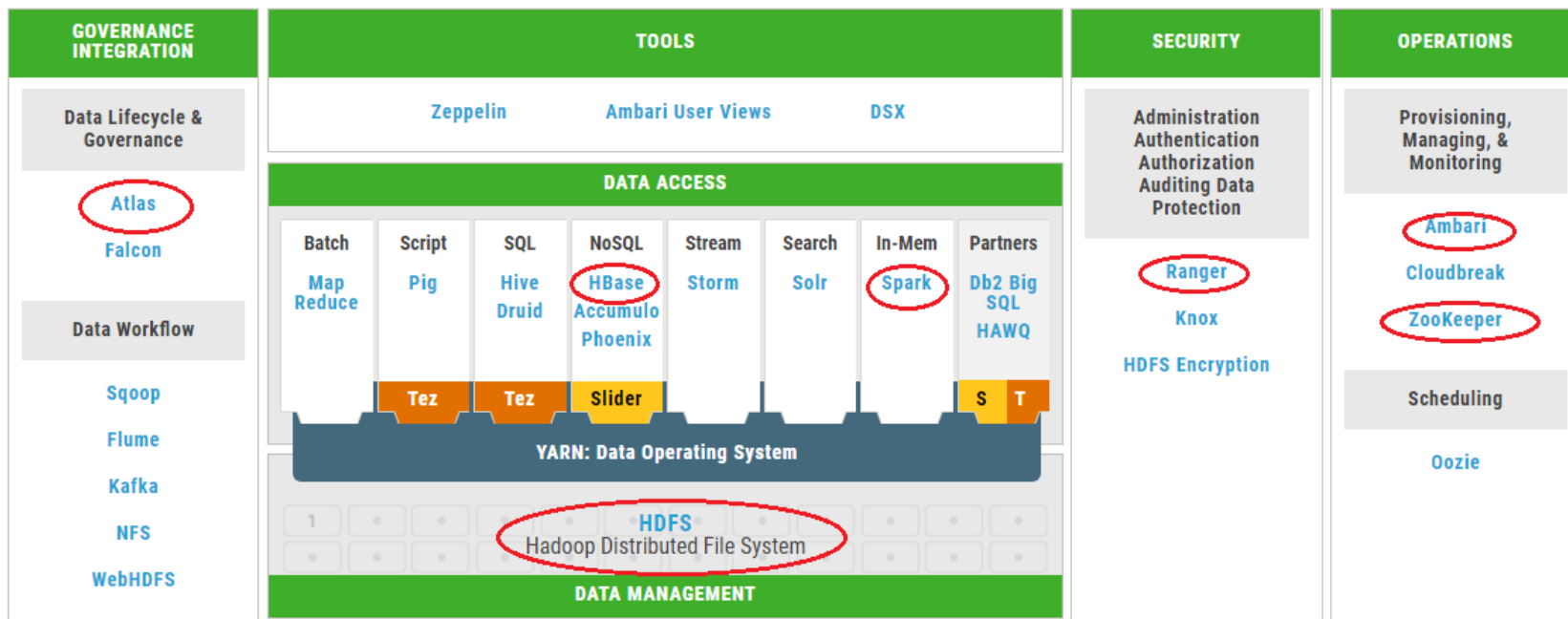


UK Data Service



The infrastructure bit - Big data and Hadoop

- Hadoop is a *suite* of different products (like Office is a suite of Excel, Access, Word, Powerpoint, Publisher etc.)





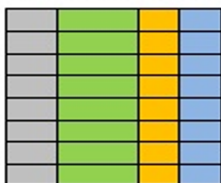
- 



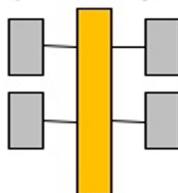
SQL - it was good while it lasted but it's time to move on

Six Types of Databases

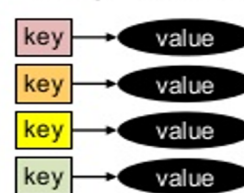
Relational



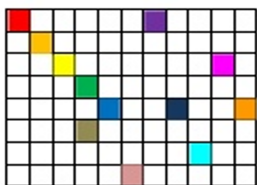
Analytical (OLAP)



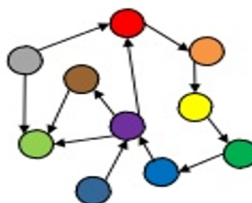
Key-Value



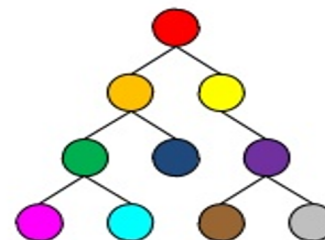
Column-Family



Graph



Document



Copyright Kelly-McCreary & Associates, LLC

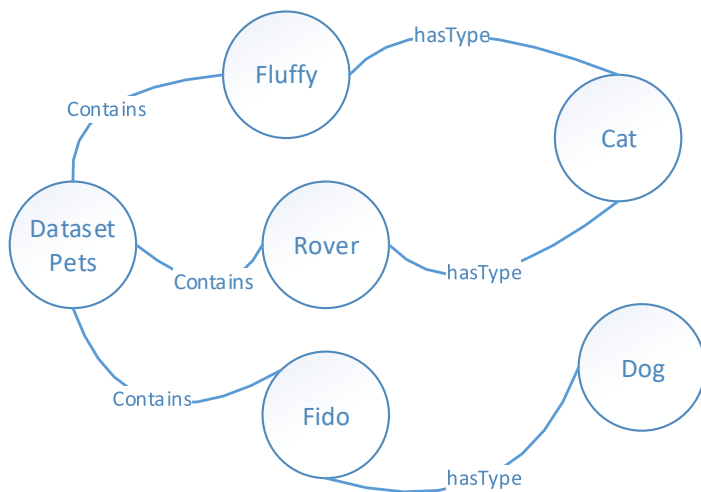
11

What and why of HBase/JanusGraph

- Traditional relational data stores will not scale and it's not always easy to alter the schema
- Instead of:

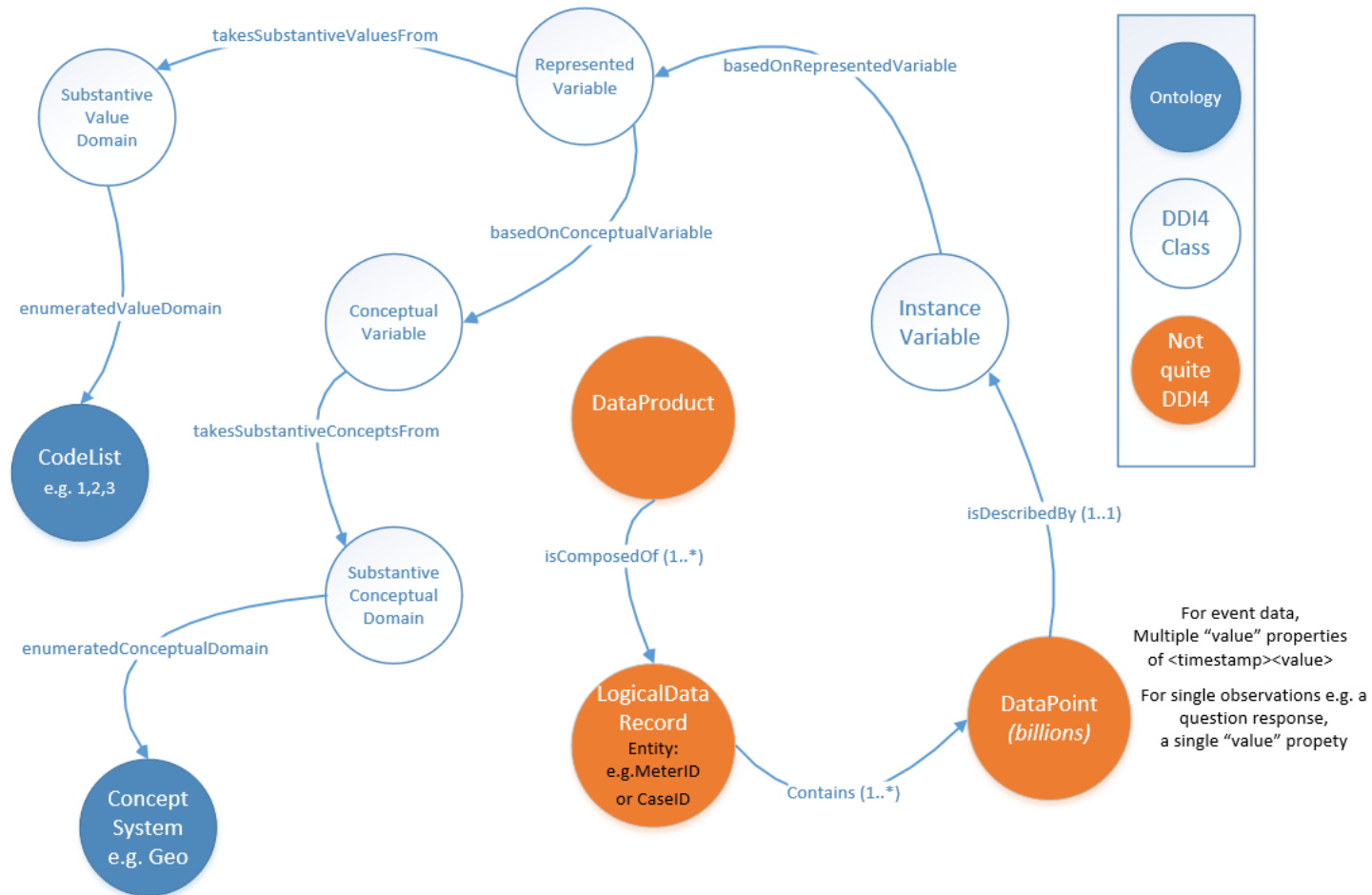
	Name	Type
Row1	Fluffy	Cat
Row2	Rover	Dog
Row3	Fido	Dog

- A property graph would be modelled like:



<i>6 statements in the database</i>		
PetsDataset	Contains	Fluffy
PetsDataset	Contains	Rover
PetsDataset	Contains	Fido
Fluffy	hasType	Cat
Rover	hasType	Dog
Fido	hasType	Dog

SERL Logical Data Model (DDI4)



The access bit

To do FAIR properly, we must look to promote more machine-actionable access and rights models:

Unify:

- Consents
- Rights
- Licensing
- Access Mediation

in a single model

ODRL (open digital rights language) provides a ready-made machine-actionable “vocabulary” to describe these.

Assets *have*

Policies *consisting of*

Rules (Permissions, Obligations and Prohibitions)

which apply to **Parties**

and which determine **Actions**

which may have **Constraints**



Simple ODRL example

```
{
  "@context": {
    "odrl": "http://www.w3.org/ns/odrl/2/"
  },
  "@type": "odrl:Agreement",
  "@id": "http://ukdataservice.ac.uk/policy:12",
  "target": "http://ukdataservice.ac.uk/asset:2000",
  "assigner": "http://ukdataservice.ac.uk/organisation:55",
  "permission": [{
    "assignee": "http://ukdataservice.ac.uk/guest:0001",
    "action": "odrl:viewmetadata"
  }],
  "permission": [{
    "assignee": "http://ukdataservice.ac.uk/group:122",
    "action": "odrl:download"
  }]
}
```

=>

For Study 2000, ONS (*organisation #55*) have declared that guest users can view the metadata and UK users (*group #122*) can download the study

Front-end UX for SERL

- (1) Single entry point. We allow researcher to search for the variables, time-ranges and geography they are interested in (a “Universal Query”)
- (2) Next step is to identify possible linkages and access criteria while incrementally filtering and/or aggregating.
- (3) Once data product is defined, we then execute the access conditions and backhaul the data into a virtual environment, with the analytic tool they have chosen to use



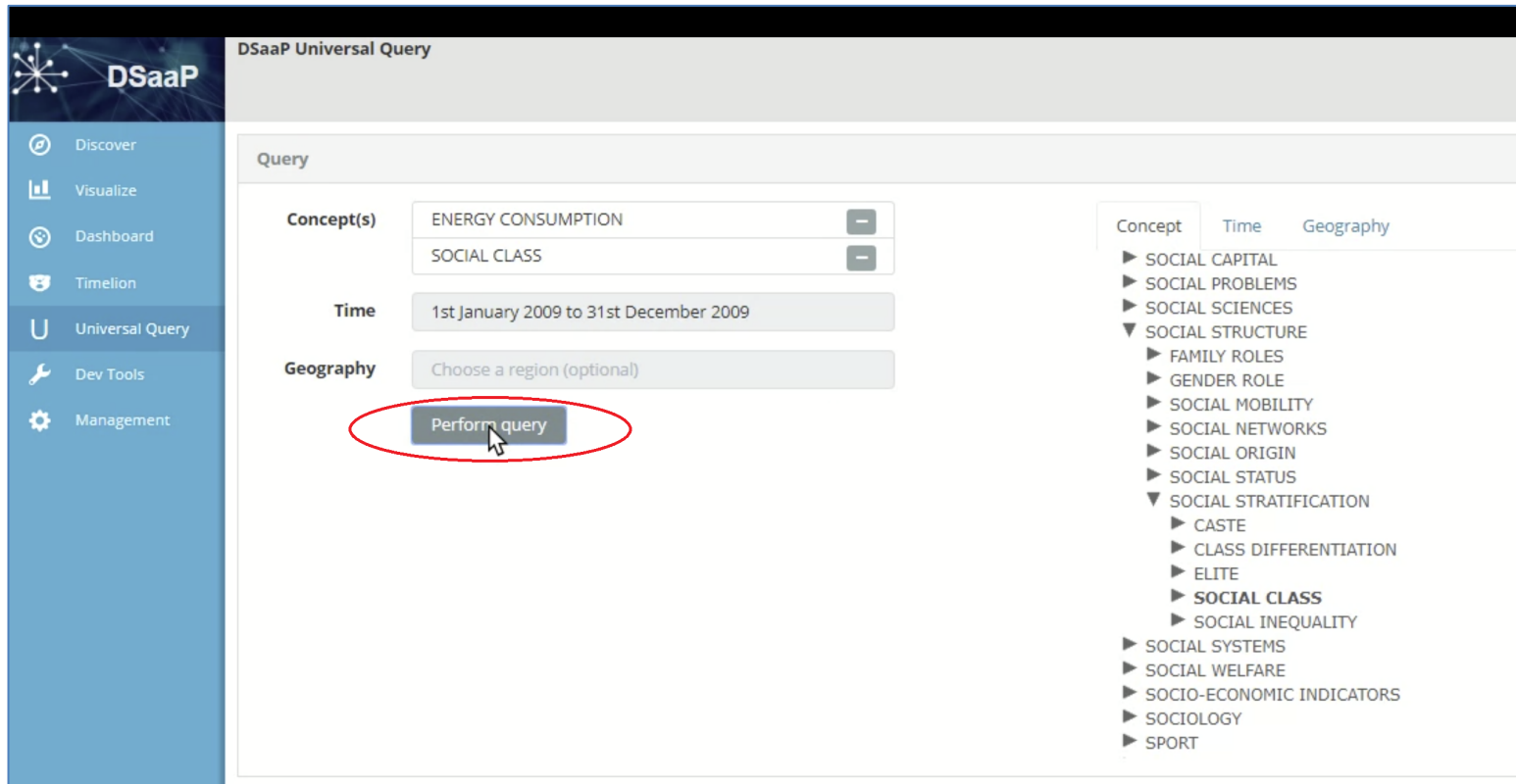
Example use case on PoC platform

- Research question: Average energy consumption by social class
- Two demographics : affluent and poor
- Geospatial visualisation by local authority district
- Uses smart meter energy data and household survey data
- Uses Spark to quickly process millions of records



Select concepts and time

In this case: “ENERGY CONSUMPTION”
and “SOCIAL CLASS” in 2009



The screenshot shows the DSaaP Universal Query interface. On the left is a navigation menu with options: Discover, Visualize, Dashboard, Timelion, Universal Query (selected), Dev Tools, and Management. The main panel is titled "DSaaP Universal Query" and contains a "Query" section with the following fields:

- Concept(s)**: Two input boxes containing "ENERGY CONSUMPTION" and "SOCIAL CLASS", each with a minus icon to its right.
- Time**: A date range input box containing "1st January 2009 to 31st December 2009".
- Geography**: A dropdown menu with the text "Choose a region (optional)".
- Perform query**: A button at the bottom of the query section, which is circled in red and has a mouse cursor pointing at it.

On the right side of the interface is a list of concepts under three tabs: Concept, Time, and Geography. The "Concept" tab is active, showing a hierarchical list of concepts:

- ▶ SOCIAL CAPITAL
- ▶ SOCIAL PROBLEMS
- ▶ SOCIAL SCIENCES
- ▼ SOCIAL STRUCTURE
 - ▶ FAMILY ROLES
 - ▶ GENDER ROLE
 - ▶ SOCIAL MOBILITY
 - ▶ SOCIAL NETWORKS
 - ▶ SOCIAL ORIGIN
 - ▶ SOCIAL STATUS
- ▼ SOCIAL STRATIFICATION
 - ▶ CASTE
 - ▶ CLASS DIFFERENTIATION
 - ▶ ELITE
 - ▶ **SOCIAL CLASS**
 - ▶ SOCIAL INEQUALITY
- ▶ SOCIAL SYSTEMS
- ▶ SOCIAL WELFARE
- ▶ SOCIO-ECONOMIC INDICATORS
- ▶ SOCIOLOGY
- ▶ SPORT

Filter/aggregate by variables

The screenshot displays the DSaaS interface with a sidebar on the left containing navigation links: Discover, Visualize, Dashboard, Timelion, Universal Query (selected), Dev Tools, and Management. The main panel shows a list of variables grouped by dataset:

- 7591_e edrp_elec**

Variable	Question	Concept
ELECKWH		ENERGY CONSUMPTION
- 7591_eg_y edrp_annual_total_energy**

Variable	Question	Concept
TOTAL_y		ENERGY CONSUMPTION
- 7591_g edrp_gas**

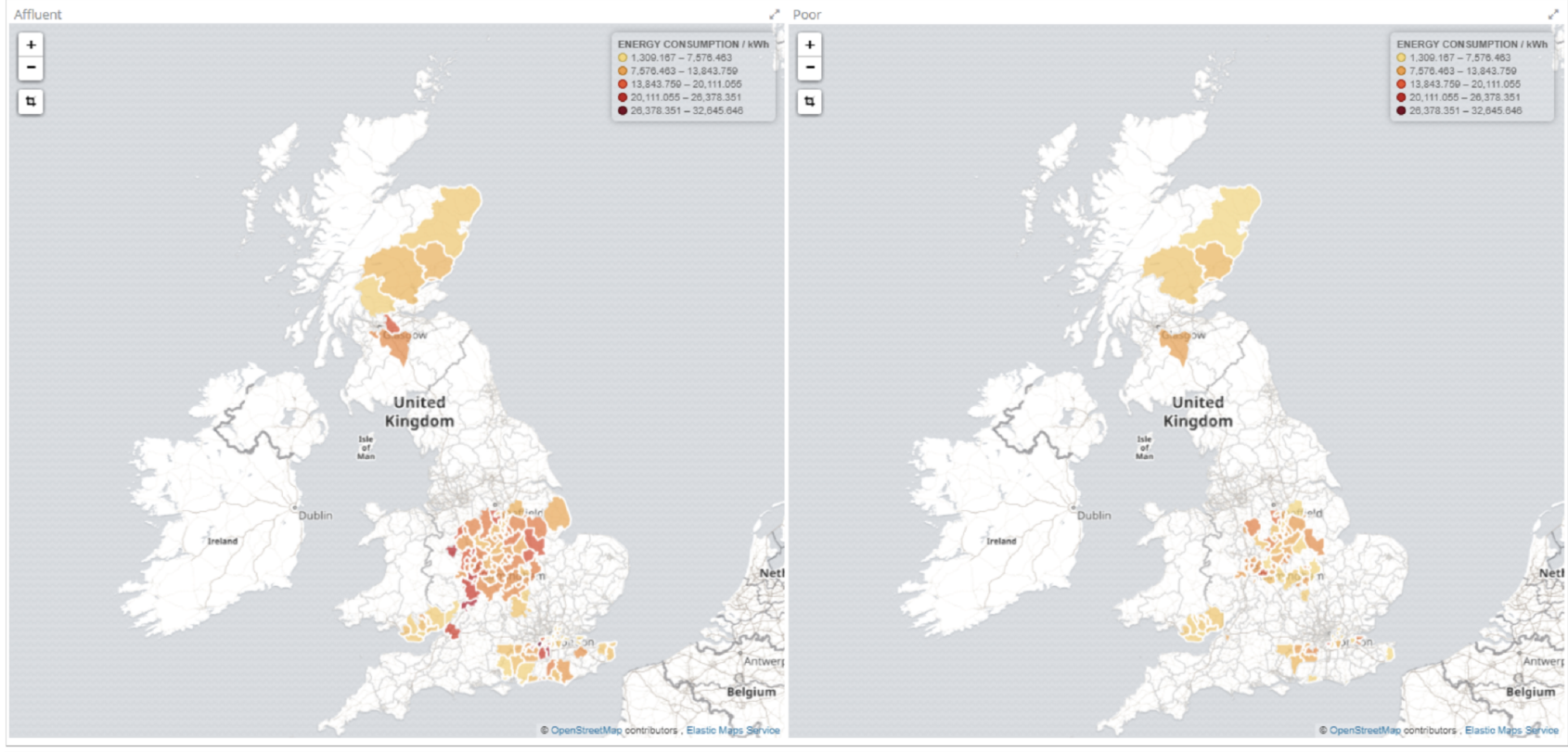
Variable	Question	Concept
GASKWH		ENERGY CONSUMPTION
- 7591_hh edrp_household.sav**

Variable	Question	Concept
ACORN_Category		SOCIAL CLASS
LAD		Geography (Local Authority)

Below the variable list are buttons for 'Visualise' (with chart icons) and 'Advanced' (with a menu icon). The 'Advanced' section is expanded, showing two configuration areas:

- Variables to visualise:** A search box containing 'TOTAL_y' and a dropdown menu showing 'ENERGY CONSUMPTION'.
- Concepts to aggregate on:** A search box containing 'SOCIAL CLASS' and a list of aggregation ranges with checkboxes:
 - 'Affluent Achievers' to 'Comfortable Communities': ☒
 - 'Comfortable Communities' to 'Financially Stretched': ☐
 - 'Financially Stretched' to 'Not Private Households': ☒

Result



Core messages

- Unification of metadata and data at lifecycle, function and process level is now possible and moreover, essential.
- Concept driven data discovery at the variable level and lower allows for powerful and flexible generation of bespoke data products.
- Standards based around semantic web and DDI(based on GSIM – Generic Statistic Information Model).
- Datum based approach = *domain-agnostic research*
- Unified access model
- Derived and reproducible *data products*

