

# Data platforms to support research, evaluation & practice

David V Ford  
Professor of Health Informatics  
School of Medicine, Swansea University





# Outline

1. Swift overview of SAIL Databank as used in Wales
2. Everything we know in a box



Research Data  
Appliances

# The SAIL Databank

A national e-research  
platform for Wales

**David Ford**

Professor of Health Informatics  
College of Medicine  
Swansea University



## SAIL: 5 Minute Overview

- SAIL= Secure **Anonymised** Information Linkage
- Databank of >9 billion recordings on the people of Wales
- Data on >5 million people
- 500+ feeder systems from Wales, inc >350 GP (75%)
- Much data goes back 10-20 years
- All pre-linked
- £5m+ investment in high performance IT
- Strong privacy protection & IG
- Remote access from anywhere, given approvals
- Industrial strength, reusable infrastructure.



# SAIL: Data resources

## Core holdings:

- Annual District Birth Extract (ADBE)
- Annual District Death Extract (ADDE)
- Bowel Screening Wales (BSW)
- Breast Test Wales (BTW)
- Cervical Screening Wales (CSW)
- Congenital Anomaly Register and Information Service (CARIS)
- Emergency department Data Set (EDDS)
- National Community Child Health Database (NCCHD)
- Outpatient Dataset (OPD)
- Patient Episode Database for Wales (PEDW)
- Primary Care GP dataset (360 practices, 2m people, 75% of practices)
- Welsh Cancer Intelligence and Surveillance Unit (WCISU)
- Welsh Demographic Service (WDS)


## Project-specific holdings:

- Clinical trials participants
- Conventional cohort participants
- Cross sectional survey participants
- Many others

## Reference data:

- Data quality reports
- Extract histories
- Coding and mapping information
- Metadata
- Organisation codes

# SAIL features

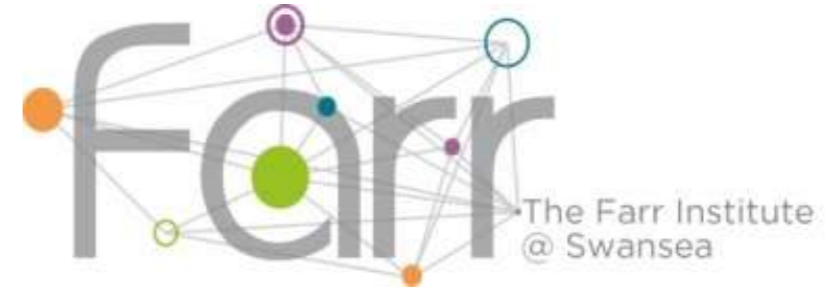
1. Built to consume (new data) swiftly 
2. Secure sharing for projects, based on views and a UKSERP instance (more later).
3. Total population coverage
4. Used for observational research, trial feasibility, outcome data for trials, extending cohorts, natural experiments . .
5. Lots of trial and cohort study participants embedded in SAIL

## SAIL: 5 Minute Overview

1. Currently supporting externally funded projects with value >£90m
2. >100 Peer-reviewed publications to date
3. 140+ approved and active SAIL projects.
4. Over 300 registered users, with 150 active today from across the world
5. 100 staff in Swansea working on Health Informatics-related projects
6. Average 35 day turnaround from application to data
7. Applications open to all

# SAIL: 5 Minute Overview

1. Now new UK Centre of Excellence in HI
  - CIPHER - Centre for Improvement in Population Health through E-records Research (£4.4m). Partners: Cardiff Bristol, Sussex, Australia, Canada
2. CIPHER now one of 4 UK-wide Farr Institute Centres (+£5m)
3. New award: CADRE - Centre for Administrative Research & Evaluation (£8.0m) ESRC, as part of the ESRC ADRN and its “Big data” initiative. Partners: Cardiff





# SAIL: 5 Minute Overview

## New Data Science building @ Swansea

- Security and privacy protection from the ground up
- Funded by MRC (Farr); ESRC (ADRN) and Welsh Government
- High security home to Farr@Swansea and CADRE
- Office space for NHS and other public sector staff and industry to work alongside university staff



# Research Data Appliances & UKSERP

Everything we know, in a box



## 5 Minute Overview

### Now in beta test:

1. **UK Secure E-Research Platform (UKSeRP)** - based on the SAIL Gateway - massively extended to provide a secure platform for data sharing across the UK - not just SAIL data (for Farr and ADRN). 270001 certified and IG Toolkit compliant
2. **National Research Data Appliances:** New concentrator technology to gift to NHS and other organisations, including automated matching, anonymisation, data management, metadata capture, data quality assessment, etc.
3. New focus on the capture and analysis of electronic free text data (on-board NLP in the Appliances)
4. Funded by MRC Farr institute and ESRC ADRN grants

# Research Data Appliance

- Brings SAIL's capabilities onto combined hardware and software
- Shrink wrapped, ready to go.
- Any size from data stick to massive!
- Easy to use, low expertise barrier
- Multiple Appliances work together as a larger whole
- Purposes: concentrate data, make it research ready
- Provide utility to data owners and partners
- Provided free (limited numbers) to our collaborators

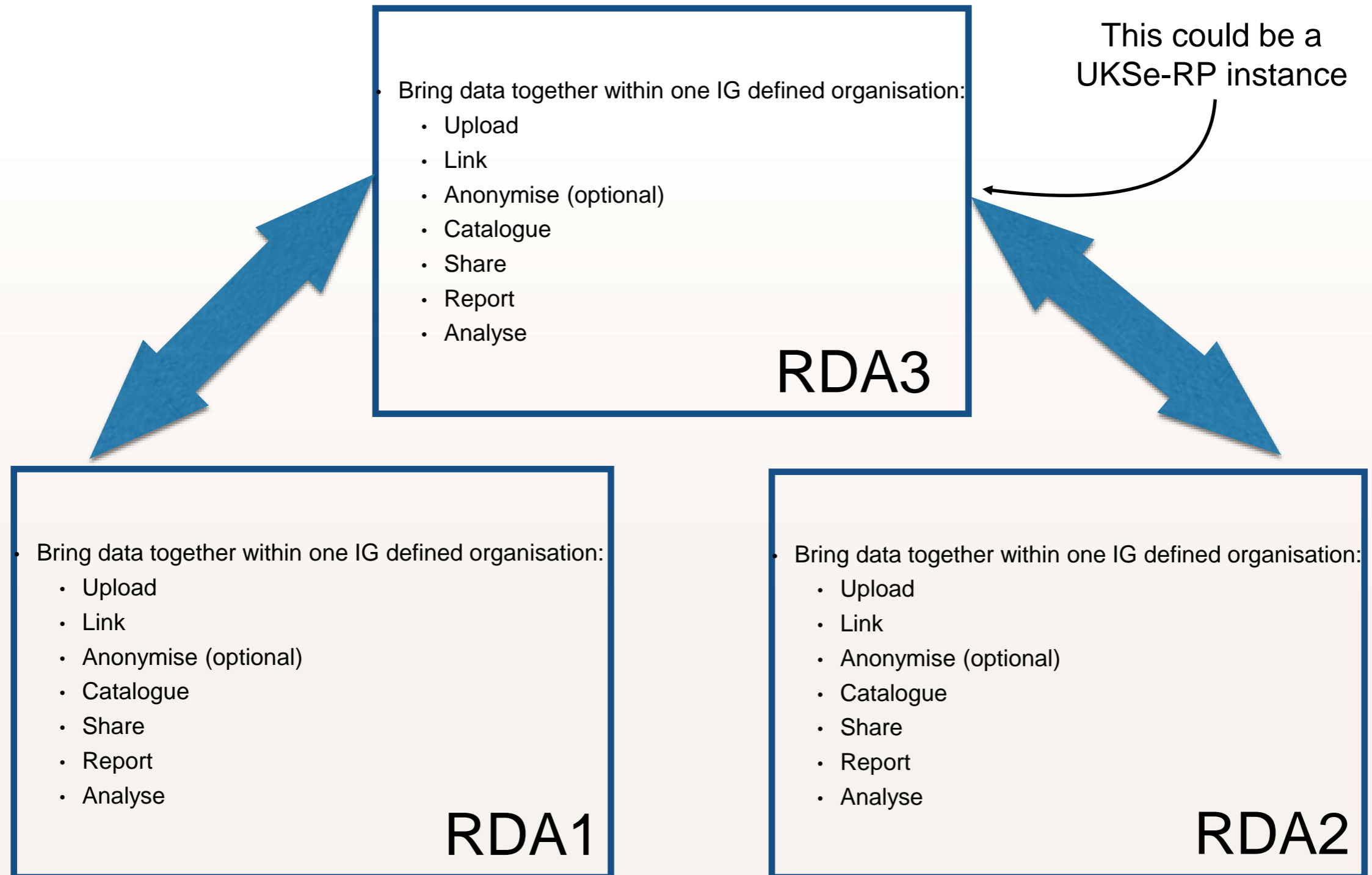
# RDA Usecase

Bring data together within one IG defined organisation:

- Upload
- Link
- Anonymise (optional)
- Measure
- Catalogue
- Share
- Pool
- Manage & Report
- Analyse

All under strong IG control

# A Simple RDA deployment model

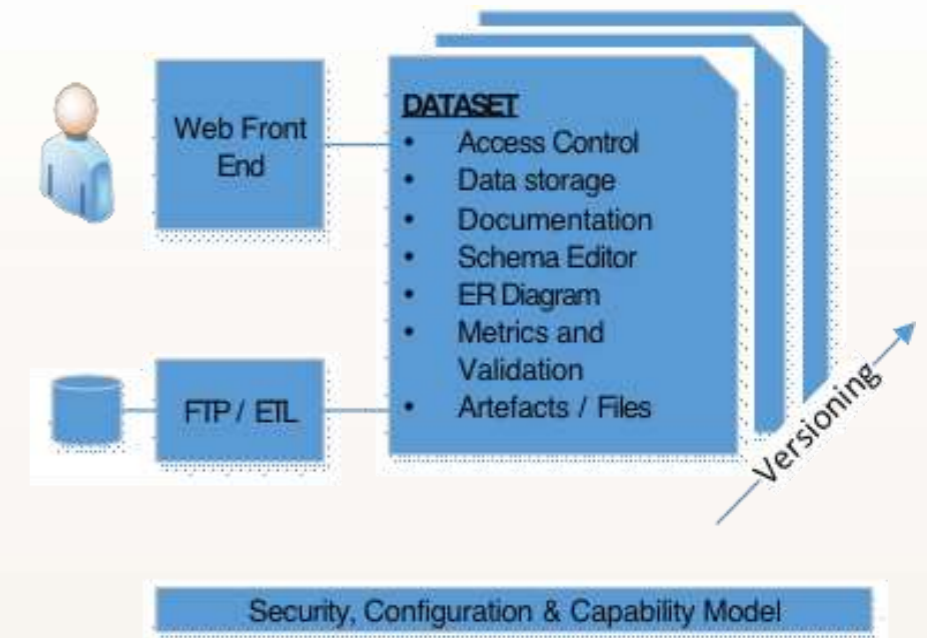


# Research Data Appliances

Concepts  
(Presented at SHIP 2013)

## Benefits to health partners:

- Data concentrator with file store and database capabilities
- On board ETL
- Very easy to use (low technical barriers)
- Built on best privacy practice principles
- On-board high quality probabilistic identity matching system
- Integrated automated Metadata Catalogue and data quality reporting
- Full IG control locally - local data for local reuse.
- Robust standardised anonymisation (as required)
- NLP facilities option for unstructured (free text) data stores



# Self management

Home Programme 1 Programme 2 External Data Catalogues

Projects & Datasets Local Data Catalogue Users & Permissions Data Out

Projects Create New Project Users Create New User

## Users & Permissions

+ Create New User

First Name	Last Name	Username
<input type="text"/>	<input type="text"/>	<input type="text"/>
Username	Telephone	Signed Agreement
<input type="text"/>	<input type="text"/>	<input type="text" value="DD/MM/YYYY"/>
Email Address	Manager Email Address	
<input type="text"/>	<input type="text"/>	

This user has no effective permissions yet.

Create New User Cancel

## Users

Sort by A-Z  Search

## Permissions

Advanced Mode  ON  OFF

### + Create New Permission

Permission Description	Aspect	Type
<input type="text" value="Enter Description..."/>	<input type="text" value="Select..."/>	<input type="text" value="Select..."/>
CalcGroup	String	Inherit down?
<input type="text" value="Enter CalcGroup..."/>	<input type="text" value="Enter String..."/>	<input checked="" type="radio"/> Yes <input type="radio"/> No

Create [Delete](#)

Permission Description DB Read	Aspect DB2	Type Read
CalcGroup [@prog].[@asp].{@perm}	String —	Inherit down? No
<a href="#">Edit</a>	<a href="#">Delete</a>	



# Congenital Anomaly Register and Information Service (CARIS)

Unpublished - Version 5 (Draft) View Log

Dataset Description Data Files Supporting Files Entity Relationship Diagram Share Settings Publish to Catalogue

## Dataset Description

### Overview

Give your Dataset a unique name, so others can understand what the data is about.

Dataset Name (Mandatory)

8

i

Theme

11

i

Congenital Anomaly Register and Information Service (CARIS)

Deprivation

Data Providing Organisation (Mandatory)

9

i

Data Type (Mandatory)

12

i

Congenital Anomaly Register and Information Service (CARIS)

Research Data (survey, questionnaire etc.)

Description (Mandatory)

10

i

Dataset Level (Mandatory)

13

i

Sed ut perspiciatis unde omnis iste natus error sit voluptatem accusantium doloremque laudantium, totam rem aperiam, eaque ipsa quae ab illo inventore veritatis et quasi architecto beatae vitae dicta sunt explicabo.

Other, please specify...

Some other text here

Artifact Name Artifact Name Artifact Name

Link a file

Purpose

14

i

Tags

15

i

Sed ut perspiciatis unde omnis iste natus error sit voluptatem accusantium doloremque laudantium, totam rem aperiam, eaque ipsa quae ab illo inventore veritatis et quasi architecto beatae vitae dicta sunt explicabo.

tag, another tag, a really long tag.

Link a file

# Congenital Anomaly Register and Information Service (CARIS)

Unpublished - Version 5 (Draft) View Log

Dataset Description Data Files Supporting Files Entity Relationship Diagram Share Settings Publish to Catalogue

## Data Files (3 files)

Upload new data file

Remote upload data file

### Friendly Name of Data Table (NEW)

CSV

Original-file-name-01.csv

Sed ut perspiciatis unde omnis iste natus error sit voluptatem accusantium doloremque laudantium, totam rem aperiam, eaque ipsa quae ab illo inventore veritatis et quasi architecto beatae vitae dicta sunt explicabo.

Data Details Data Import Settings

Options

Include in Data Quality Report

Status: Errors - Not ready to publish

## Data Details

Name (Mandatory)

2

i

Date Format

5

i

Friendly Name of Data Table

DD/MM/YYYY

Table Name

3

i

Date / Time Format

6

i

friendly-name-of-data-table

DD/MM/YYYY 00:00:00

Description (Mandatory)

4

i

Distribution Column

7

i

Sed ut perspiciatis unde omnis iste natus error sit voluptatem accusantium doloremque laudantium, totam rem aperiam, eaque ipsa quae ab illo inventore veritatis et quasi architecto beatae vitae dicta sunt explicabo.

PROV\_UNIT\_CD

Save and Close

## Data Files (3 files)

Upload new data file

Remote upload data file



### Friendly Name of Data Table (NEW)

Original-file-name-01.csv

Sed ut perspiciatis unde omnis iste natus error sit voluptatem accusantium doloremque laudantium, totam rem aperiam, eaque ipsa quae ab illo inventore veritatis et quasi architecto beatae vitae dicta sunt explicabo.

Data Details

Data Import Settings

Origin: My Dataset

Options

Include in Data Quality Report

Status: **Errors - Not ready to publish**

## Data Import Settings

Personal Identifiable Data (PID) 29

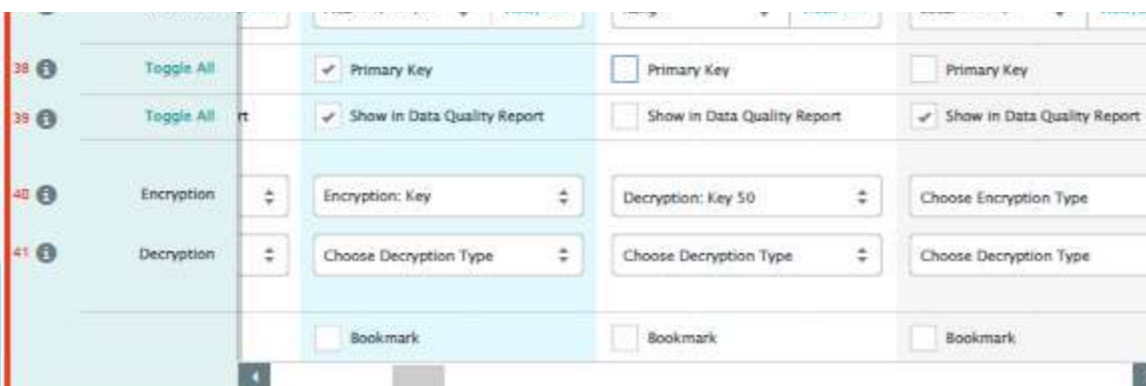
Choose template

Import Settings 30


Choose existing settings

Import Settings from Local Data Catalogue 31

View list of settings



Data “schema” automatically computed based on data contained in uploaded file

	All Fields	A-Z	Errors (3)	Bookmarks (5)	PID (3)
Field Name	PROV_UNIT_CD 	SPELL_NUM_E	EPI_NUM		
Friendly Name	Organisation Code (Unit Code of Provider)	Hospital Provider Spell Number	Episode Number		
Field Description	This is the organisation code of the health care provider. The provider code identifies the ...	This is the organisation code of the health care provider. The provider code identifies the ...	This is the organisation code of the health care provider. The provider code identifies the ...		
Personal Identifiable Data (PID) Type	N/A	NHS Number	N/A		
Field Type	Char <input type="text" value="5"/> Size <input type="text" value="3"/>	Float <input type="text" value="15"/> Precision <input type="text" value="5"/> Scale	Char <input type="text" value="2"/> Size		
Validation	Please Specify... <input type="text"/> View/Edit	Range <input type="text"/> View/Edit	Local Lookup <input type="text"/> View/Edit		
Toggle All	<input checked="" type="checkbox"/> Primary Key	<input type="checkbox"/> Primary Key	<input type="checkbox"/> Primary Key		
Toggle All	<input checked="" type="checkbox"/> Show in Data Quality Report	<input type="checkbox"/> Show in Data Quality Report	<input checked="" type="checkbox"/> Show in Data Quality Report		
Encryption	Encryption: Key <input type="text"/>	Decryption: Key 50 <input type="text"/>	Choose Encryption Type <input type="text"/>		
Decryption	Choose Decryption Type <input type="text"/>	Choose Decryption Type <input type="text"/>	Choose Decryption Type <input type="text"/>		
	<input type="checkbox"/> Bookmark	<input type="checkbox"/> Bookmark	<input type="checkbox"/> Bookmark		

# Publish Dataset

- Depend on Configuration/Capabilities. Data will now be available

Entity Relationship Diagram   Share Settings   **Publish to Catalogue**

## Publishing Checklist

All Data Files have been validated	✓
Remote Files uploaded	✓
All mandatory information is complete	⚠ (5) Errors, fix these
Entity Relationship Diagram confirmed	✓

**Publish to Catalogue**

**When you publish some things will happen...**

- A list of things that will happen to the dataset.
- These things have been previously defined.
- This is just a static statement.

Name of Dataset

Version 2 (Draft)   **Progress**

Theme Health  
Date Type Clinical System

## Publishing Progress

Elapsed Time: 10 Weeks, 30 Days and 23 Hours

Data File 1	Complete	✓
Data File 2	Error	View Log
Really long names will get...	Metrics	
Data File 4	Validating	
Data File 5	Cleaning Up	
Generate Data Quality Report	Waiting	
Converting to Data Catalogue	Waiting	
Seal & Create New Version	Waiting	

**Close**

# Data Catalogue - Key Component

The screenshot displays the SAIL DataBank interface. At the top, there are navigation tabs for 'Projects & Datasets', 'Local Data Catalogue', and 'Users & Permissions'. Below these, there are sections for 'My Datasets (5)' and 'Favourite Datasets (10)'. A search bar contains the text 'Wales congenital' with a dropdown menu showing suggestions: 'Wales congenital', 'wales congenital anomaly', 'wales congenital', and 'wales'. The search results section is titled 'Wales congenital anomaly (107 Datasets)'. Below this, there are filters for 'Sorted by Relevance', 'All Themes (107)', 'All Data Types (107)', and 'All Dataset Levels (107)'. The main result is for 'Congenital Anomaly Register and Information Service (Caris)' with a star icon and a version dropdown set to 'Version 5 - Published 21/03/2013'. The description area is redacted with black bars. To the right, there are details for 'Theme: Health', 'Date Type: Clinical Sy:', 'Dataset Level: Individual Person', and 'Tags: Health, Wales, CARIS, ABMU, Swansea, Congenital Anomaly Register'. A 'Data Providing Organisation' section lists 'Congenital Anomaly Register and Information Service for Wales'.

# A Dataset

The screenshot shows a web interface for a dataset. At the top, there's a navigation bar with 'Home', 'Programme 1', 'Programme 2', and 'External Data Catalogues'. Below that, there are tabs for 'Projects & Datasets', 'Local Data Catalogue', and 'Users & Permissions'. A search bar is present with 'My Datasets (5)' and 'Favourite Datasets (10)' filters. The main content area features a dark blue header for the dataset: 'Congenital Anomaly Register and Information Service (CARIS)'. Below this, it shows 'Version 5' and 'Published 21/03/2013'. To the right, there's an 'Administrative Contact' section for 'Dr. Jones' with contact details and a 'Request Subscription to Data' button. Below the contact info are 'Request Subscription to Data' and 'Print Dataset Page' buttons. The 'Overview' section contains a 'Description' with placeholder text and a list of 'Data File' attachments. To the right of the description is a 'Data Quality Report' (2.6MB) and a 'Theme' section listing 'Health' and 'Clinical System Data'. Below the theme is the 'Dataset Level' 'Individual Person' and a 'Tags' section with 'Health, Wales, CARIS, ABMU, Swansea, Congenital Anomaly Register'.

Specific version & Date

Contact

Request

All section attach files

VIMO

Theme / Type / Level

Tags

# A Dataset (cont.)

## Coverage

### Data Coverage

Sed ut perspiciatis unde omnis iste natus error sit voluptatem accusantium doloremque laudantium, totam rem aperiam, eaque ipsa quae ab illo inventore veritatis et quasi architecto beatae vitae dicta sunt explicabo.

Nemo enim ipsam voluptatem quia voluptas sit aspernatur aut odit aut fugit, sed quia consequuntur magni dolores eos qui ratione sequi nesciunt.

- [Data File](#)
- [Another Data File](#)
- [Another Data File](#)
- [Another Data File](#)

[View all Files](#)

### Inclusion / Exclusion Criteria

Nemo enim ipsam voluptatem quia voluptas sit aspernatur aut odit aut fugit, sed quia consequuntur magni qui ratione voluptatem sequi nesciunt.

- [Data File](#)

[View all Files](#)

## Collection

### Data Collection Method

This information is not available.

- [Data File](#)

[View all Files](#)

### Dataset Period

March 1989 - December 2001

- [Data File](#)

[View all Files](#)

### Refresh Frequency

Data is refreshed every month w processing.

- [Data File](#)

[View all Files](#)

## Highlights & Know Issues

### Data Highlights

Sed ut perspiciatis unde omnis iste natus error sit voluptatem accusantium doloremque laudantium, totam rem aperiam, eaque ipsa quae ab illo inventore veritatis et quasi architecto beatae vitae dicta sunt explicabo.

Nemo enim ipsam voluptatem quia voluptas sit aspernatur aut odit aut fugit, sed quia consequuntur magni dolores eos qui sequi nesciunt.

- [Data File](#)

[View all Files](#)

### Known Issues

Nemo enim ipsam voluptatem quia voluptas sit aspernatur aut odit aut fugit, sed quia consequuntur magni dolores eos qui ratione voluptatem sequi nesciunt.

- [Data File](#)

[View all Files](#)

[Back to top](#)

## Data Files

### Data Tables (2/2)

[View all](#)



#### Name of Data Table

At vero eos et accusamus et iusto odio dignissimos ducimus qui blanditiis praesentium voluptatum deleniti atque corrupti quos dolores et quas molestias.



#### Name of Data Table

At vero eos et accusamus et iusto odio dignissimos ducimus qui blanditiis praesentium voluptatum deleniti atque corrupti quos dolores et quas molestias excepturi sint occaecati cupiditate non provident, similique sunt in culpa qui officia deserunt mollitia animi, id est laborum et dolorum fuga.



DDI, SPSS, SAS, STATA

### SAIL Specific

Dataset Category  
Core

#### Access Requirements

As a core SAIL dataset available in accordance with our standard Information Governance procedure.

## Supporting Files

### Supporting Files (5/20)

[View all](#)



#### Entity Relationship Diagram

orig-file-name.png (345 kb)

At vero eos et accusamus et iusto odio dignissimos ducimus qui voluptatum...



#### Friendly Name Here

orig-file-name.doc (4.2 mb)

At vero eos et accusamus et iusto odio dignissimos ducimus qui voluptatum...



#### Friendly Name Here

orig-file-name.pdf (1.6 mb)

At vero eos et accusamus et iusto odio dignissimos ducimus qui voluptatum...



#### DDI File

orig-file-name.pdf (3.6 mb)

At vero eos et accusamus et iusto odio dignissimos ducimus qui voluptatum...



#### Another File

orig-file-name.pdf (3.6 mb)

At vero eos et accusamus et iusto odio dignissimos ducimus qui voluptatum...

# Data Catalogue - a specific table

**Friendly Name of Data Table**

Download Modified 3 months ago
Records: 55,160,095 | Fields: 9

**Description**

Sed ut perspiciatis unde omnis iste natus error sit voluptatem accusantium doloremque laudantium, totam rem aperiam, eaque ipsa quae ab illo inventore veritatis et quasi architecto beatae vitae dicta sunt explicabo.

Nemo enim ipsam voluptatem quia voluptas sit aspernatur aut odit aut fugit, sed quia consequuntur magni dolores eos qui ratione voluptatem sequi nesciunt.

Nemo enim ipsam voluptatem quia voluptas sit aspernatur aut odit aut fugit, sed quia consequuntur magni dolores eos qui ratione voluptatem sequi nesciunt.

**Validation (VIMO)**

Category	Percentage
Valid	80%
Invalid	10%
Missing	7%
Outliers	3%

**Connection String**

```
server.database.table.12345
```

**Heading?**

All Fields | A-Z

Code Name	PROV_UNIT_CD	SPELL_NUM_E	EPI_NUM	OPER_DT
Friendly Name	Organisation Code (Unit Code of Provider)	Hospital Provider Spell Number	Episode Number	Operatio date
Description	This is the organisation code of the health care provider. The provider code identifies...	A number (alphanumeric) to provide a unique identifier for each hospital provider...	A number used to identify episodes uniquely, and is a sequence number for each...	A numbe position to a pati
Field Type	char Size 3	int Size 4	char Size 2	smallint
VIMO	100% Valid	100% Valid	91.43% Valid	100%
Metrics	Top 10 Values	Top 10 Values	Top 10 Values	Date

Records: 55,160,095
Fields: 9

**VIMO**

Valid | Invalid | Missing | Outliers

**Top 10 Values**

Rank	Percentage
1	93.98%
2	4.39%
3	1.02%
4	0.32%
5	0.11%
6	0.05%
7	0.02%
8	0.01%
9	0%
10	0%

**Datetime Range**

Start Datetime: 1900-01-01 00:00:00

End Datetime: 2020-12-28 00:00:00

**View Min, Max, Mean**

Min: 1 | Max: 12

Mean: 2.18225767767603



# Key features of relevance

- Appliances federate with each other to create a sharing network
- Network can be hierarchical or peer-to-peer
- Full IG controls available at every point
- Metadata builds automatically and publishes to a global catalogue
- Data quality measurement automated
- NLP address rich, free text datasets, converting them to SNOMEDCT
- High quality identity reconciliation automatic, de-identification optional
- UKSERP provides scalable, performant analytics platform, with full IG controls

# Questions?

