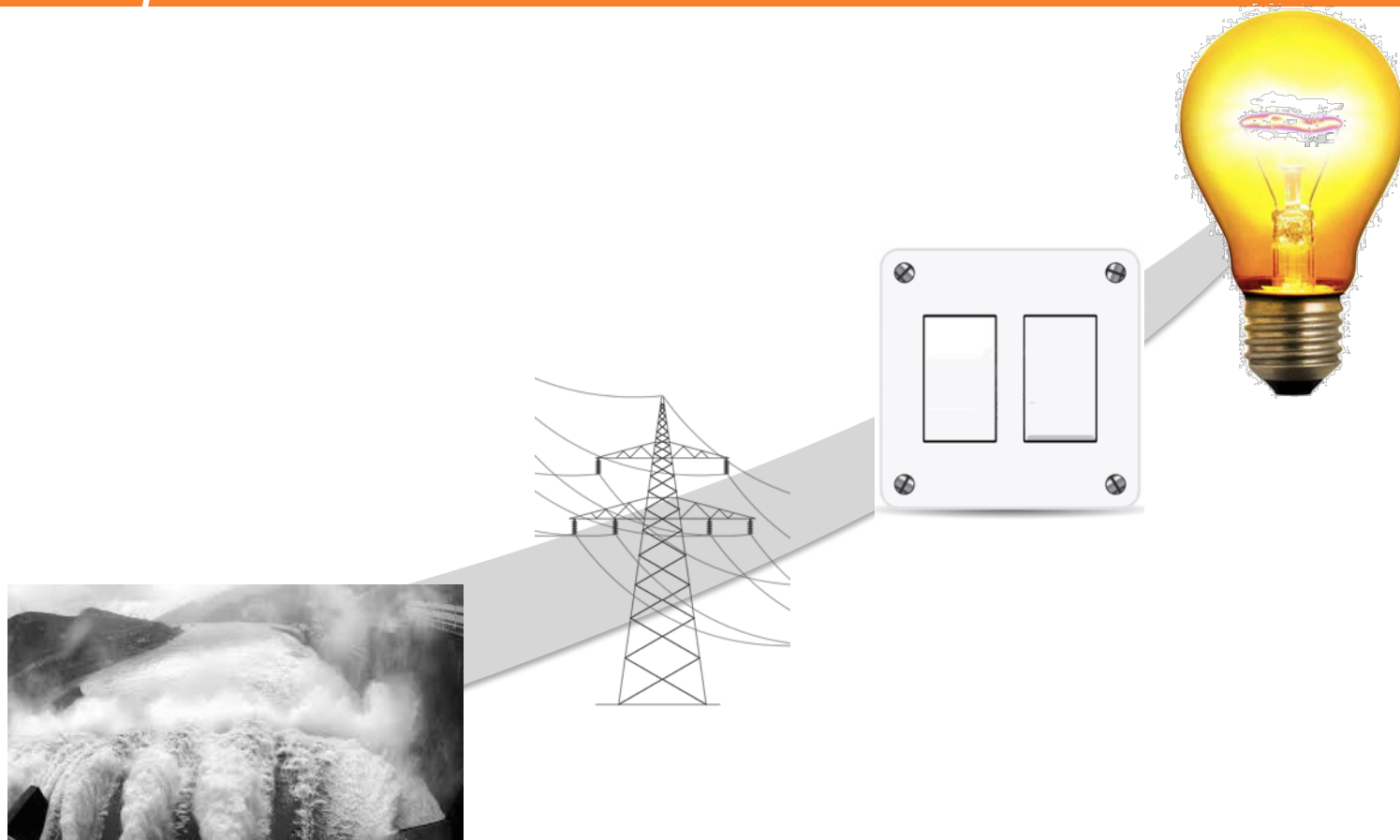# European life-science data infrastructure: Data, Computing and Services to Communities

*Niklas Blomberg*
*European Life Sciences Infrastructure for Biological Information*
*www.elixir-europe.org*

*What appears to be a simple, reliable user experience...*

*...is made possible by robust, non-trivial infrastructure that often goes unnoticed.*

# *Biomedical research is requiring increasingly sophisticated infrastructure.*
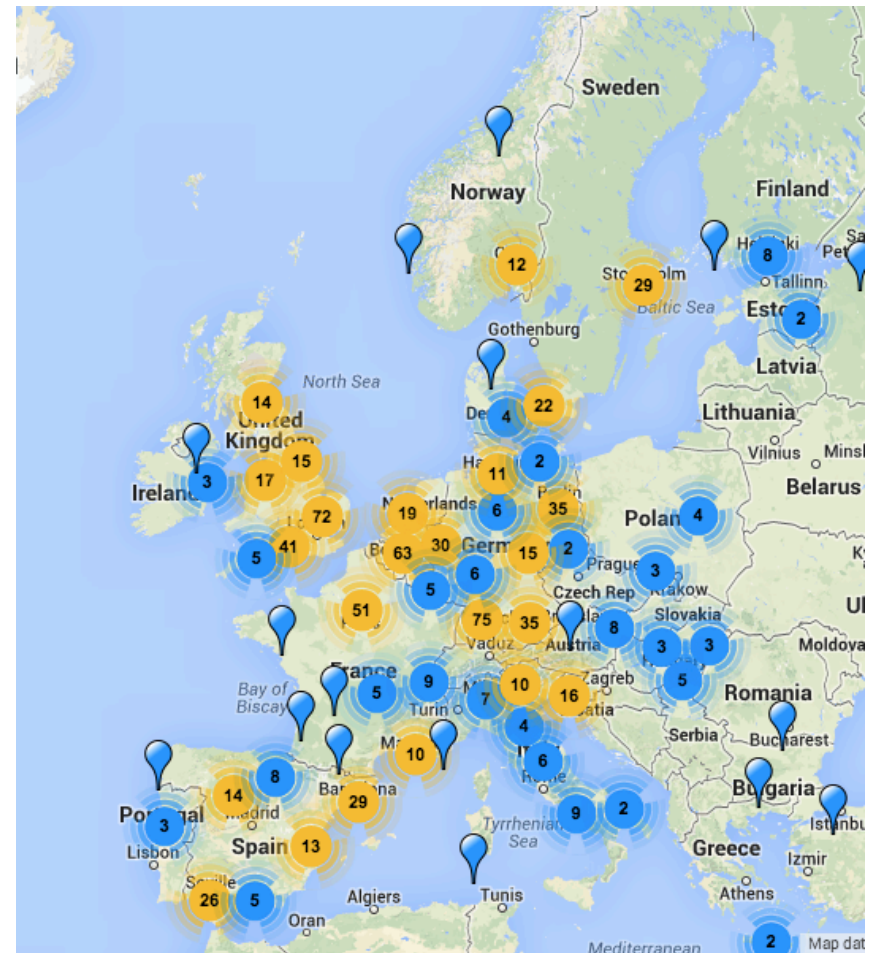


No formal infrastructure
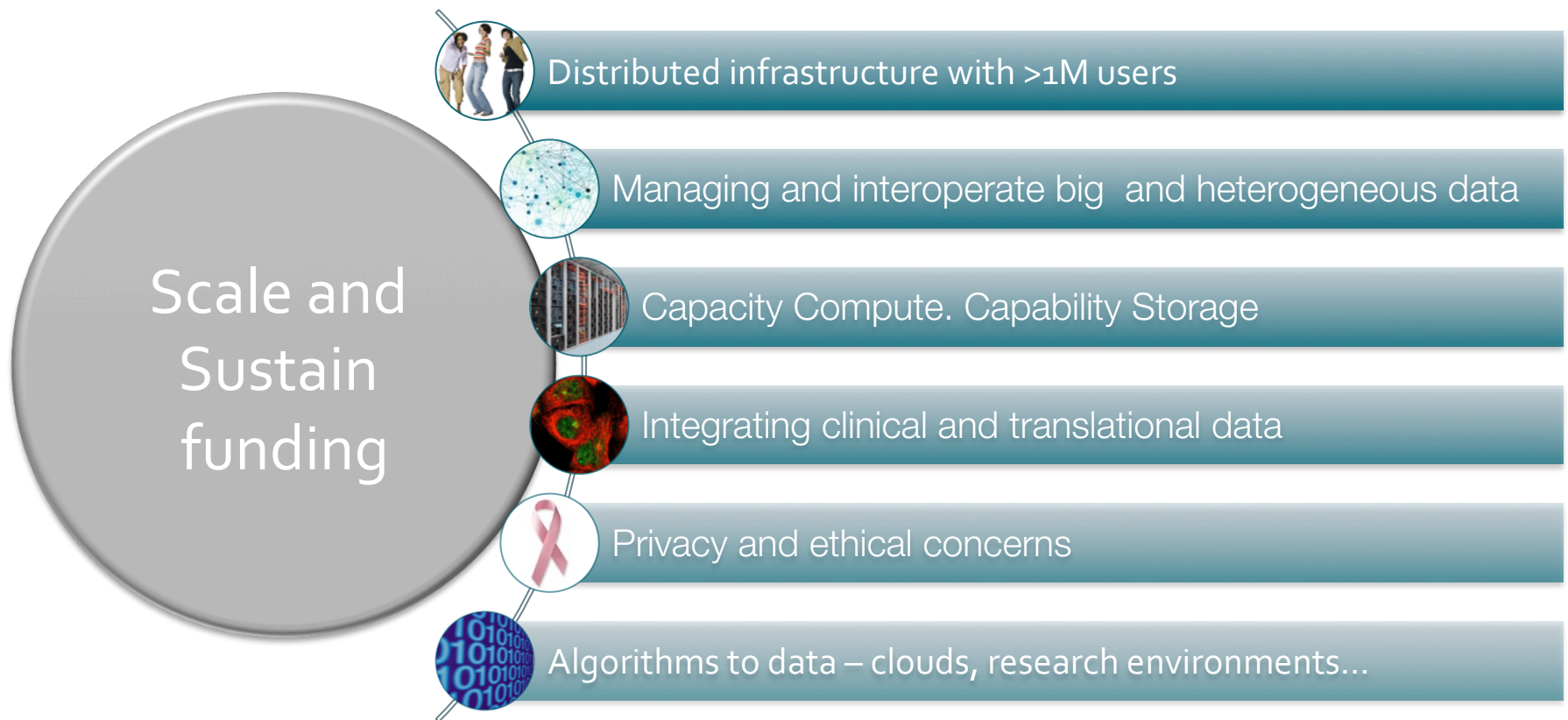(single use)

Basic infrastructure
(local)

Robust infrastructure
(large-scale, non-trivial,
interconnected)

# Life-science and data infrastructure

- Data production and using at a large number of sites across Europe
  - (European Illumina sales up 20% 2013)

- Human genomics projects but also plants, microbiota, environmental marker organisms

- Metabolomics & Proteomics coming of age
  - UK National Phenome facility

- Be scalable to 1000s of sites

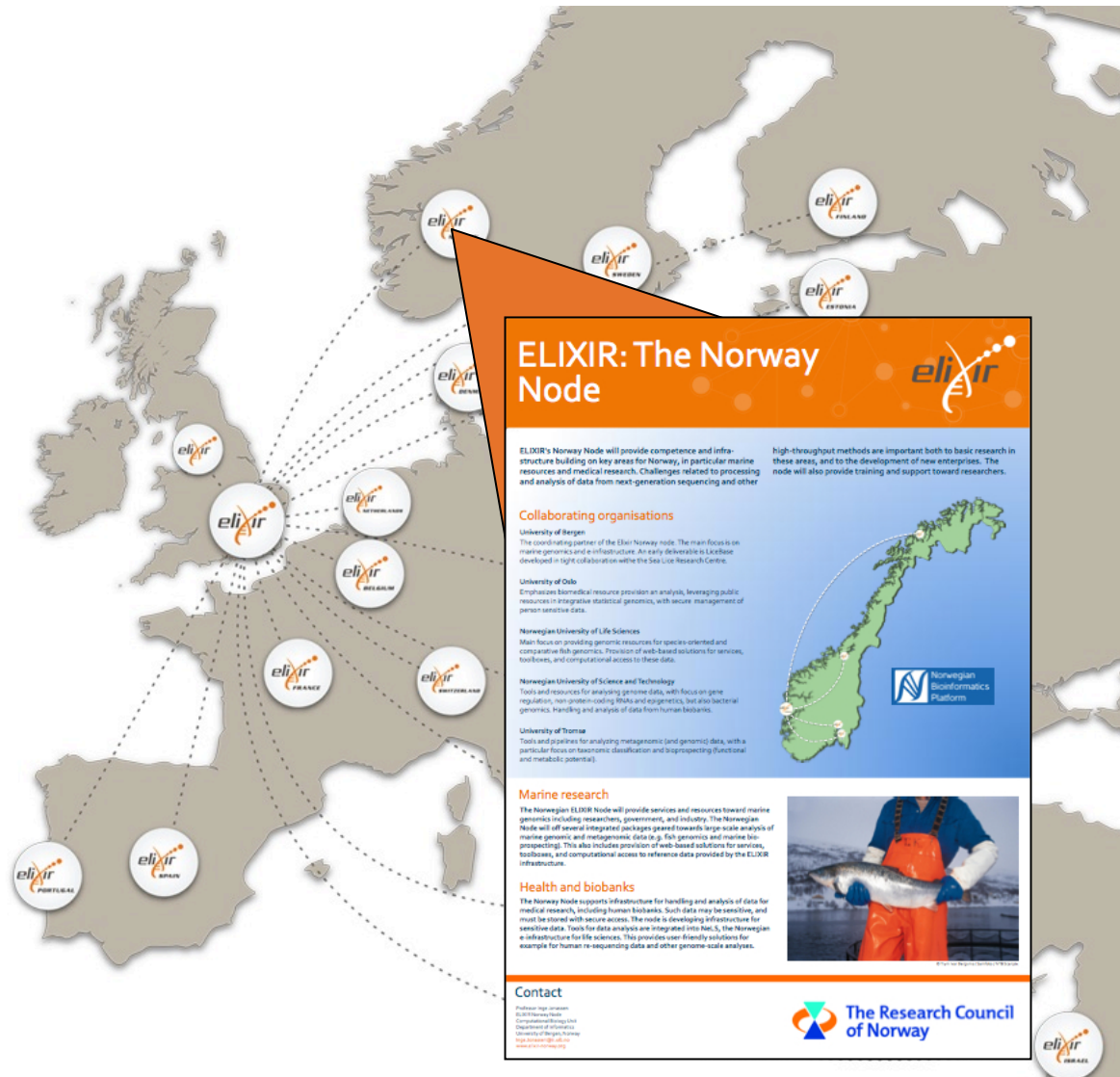- Deal with incomplete, conflicting, and incorrect data

# Challenges for life-science data services



**Scale and Sustain funding**

- Distributed infrastructure with >1M users
- Managing and interoperate big and heterogeneous data
- Capacity Compute. Capability Storage
- Integrating clinical and translational data
- Privacy and ethical concerns
- Algorithms to data – clouds, research environments…

elixir

# A distributed infrastructure to scale with the challenges

- ELIXIR deliver services through national ELIXIR Nodes

- ELIXIR Nodes build local bioinformatics capacity throughout Europe

- ELIXIR Nodes build on national strengths and priorities



http://www.elixir-europe.org/about/elixir-nodes

# *ELIXIR*

- Elixir Consortium Agreement (ECA) entered into legal force Jan 2014

- 8 countries signed to date
  - Czech Republic, Denmark, Estonia, Netherlands, Norway, Sweden, Switzerland and the UK

- Further 9 countries have signed MoU and are working towards national signatures

- Discussions on-going with additional prospective member states

# ELIXIR Infrastructure

- **Data**

  *Sustain core data resources*

- **Tools**

  *Services & connectors to drive access and exploitation*

- **Compute**

  *Access, Exchange & Compute on sensitive data*
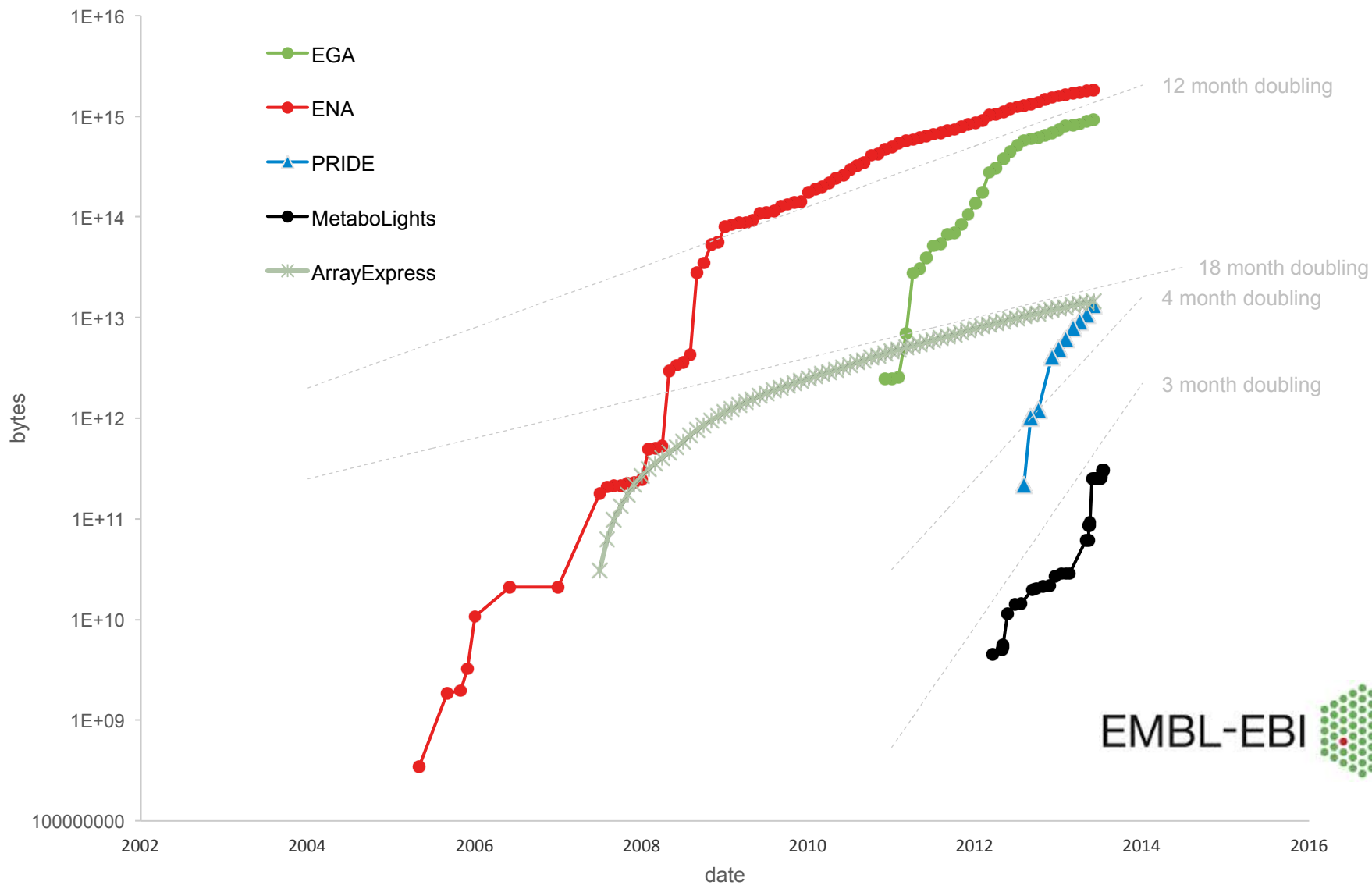
- **Standards**

  *Integration and interoperability of data and services.*

- **Training**

  *Professional skills for managing and exploiting data*

# Growing data

# *ELIXIR pilots to address key challenges in biomedical research:*



1. Cloud computing
   **"Embassy cloud"**: Access reference data in a virtual environment – work as though you are at EMBL-EBI or SIB, Switzerland

2. Authentication & Authorisation
   Improved methods and processes for access to clinical data

3. High-Performance Computing
   **"Lightpath"**: Connections for on-demand reference data to remote HPC centres at EMBL-EBI and CSC Finland

# Cross-site VM Operation - pilot

- Perform analysis via cloud infrastructures and VMs

- Transfer VMs between computing centers to allow researchers to perform analyses that they could not otherwise do locally

- Supported by 5 NRENs and in collaboration with EGI
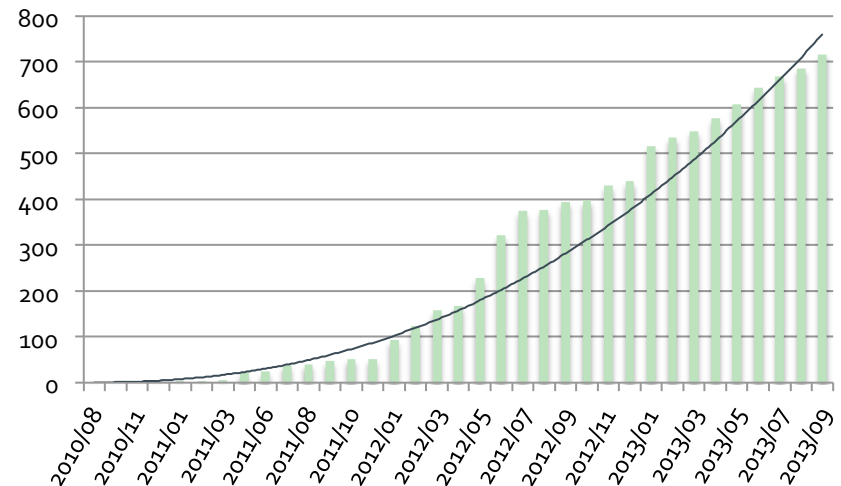
# ELIXIR Pilot: EGA as a joint venture

- Primary archive for any data consented for sharing in the context of research but not for fully public distribution

  - Secure storage, management and dissemination of data – raw or processed - from biomedical research projects.

  - Phenotypic data collected from the subjects.

  - Submissions must be de-identified and in accordance with the informed consent.

  - Data are packed into datasets that are governed by a Data Access Committee (DAC).

    - Authentication - each DAC approved individual will have a personal EGA account.

    - Authorization – DACs attach access permission(s) to the EGA account(s).

- EGA hosts more than 450 studies and discoverability to the 732 that are in both EGA and dbGaP

- EGA supports more than 400 user requests per month



Under the ELIXIR pilot project, the CRG and the EBI have agreed to "Explore ways in which the CRG's emerging Node could share responsibility for production of the EGA in future"…

… which translates into managing peer database representations of the EGA Project hosted jointly by the Hub and the Spanish node of ELIXIR
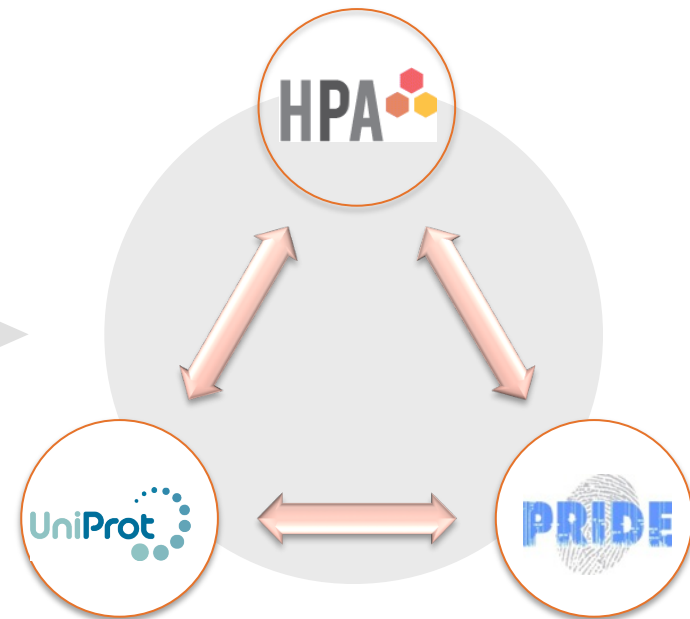
european genome-phenome archive

# Data interoperability – Human Protein Atlas
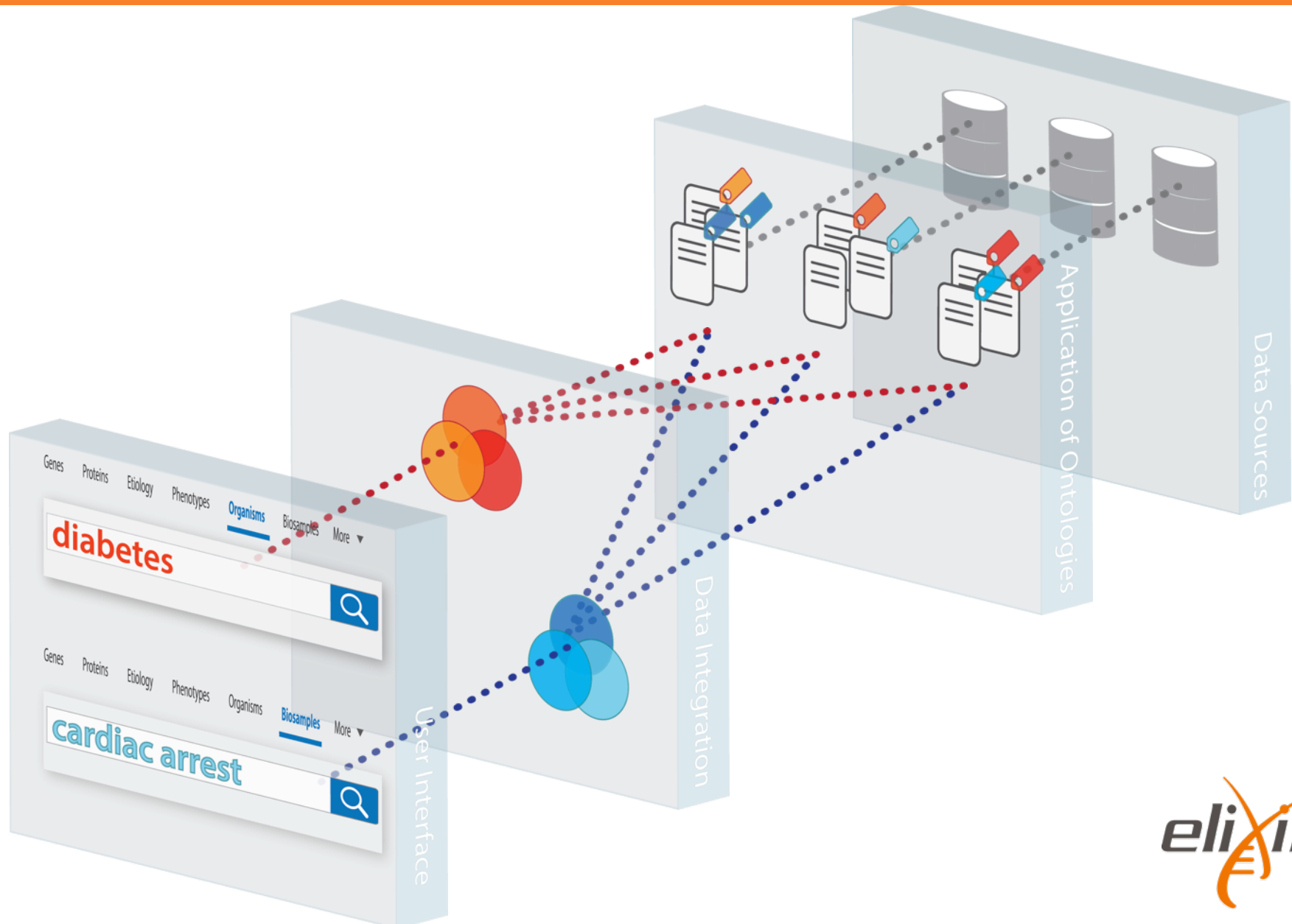


The Human Protein Atlas portal is a publicly available database with millions of high-resolution images showing the spatial distribution of proteins in 46 different normal human tissues and 20 different cancer types, as well as 47 different human cell lines.
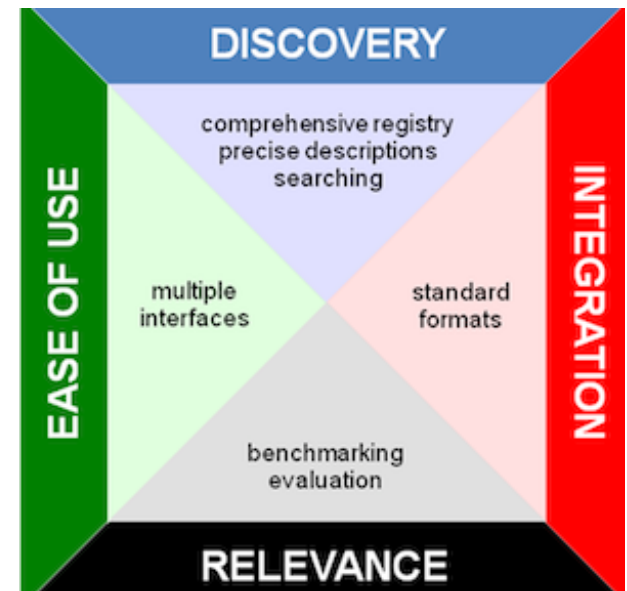
- Computational 'data and service' bridges between the BMS RIs

- Interoperability between data and services in the biological, medical, translational and clinical domains

- Link basic biological research data with clinical research and associated data



Data infrastructure
ELIXIR: the pan-European research infrastructure for biological information

Biobanking
BBMRI: The Biobanking and Biomolecular Resources Research Infrastructure

Drug discovery
EU-OPENSCREEN: the European Infrastructure of Open Screening Platforms for Chemical Biology

From mouse to human
Infrafrontier: The infrastructure for mouse disease models and phenotype data

Controlling epidemics
ERINHA: the European Research Infrastructure for Highly Pathogenic Organisms

Natural products from marine environments
The European Marine Biological Resource Centre (EMBRC)

Molecular structures
INSTRUCT: Integrated structural biology unlocking the secrets of life

Clinical trials
ECRIN: The European Clinical Research Infrastructure Network

Translational medicine
EATRIS: The research infrastructure for translational medicine.

Identifying biomarkers
Euro-BioImaging: the research infrastructure for imaging technologies

# Service and Resource Registry

- Provides a simple search interface

- Content: 1943 tools etc., 22,232 annotation
  - E.g. URL, text, ontology term: type, formats ..

- Classifies tools using an ontology
  - E.g. Sequence analysis tool

- Download complete content

- Supports a wide scope of tools

- Provides an interface to the literature

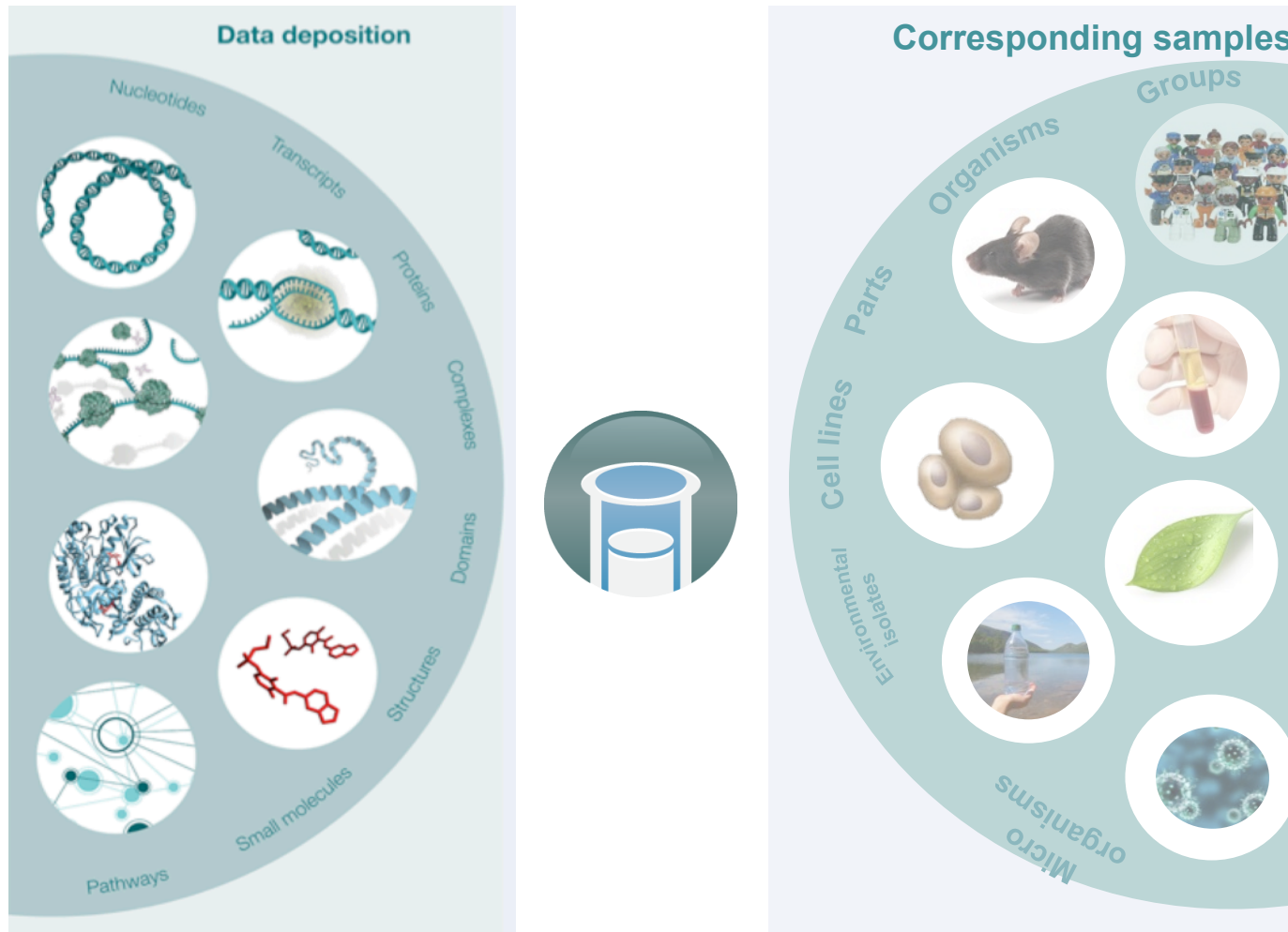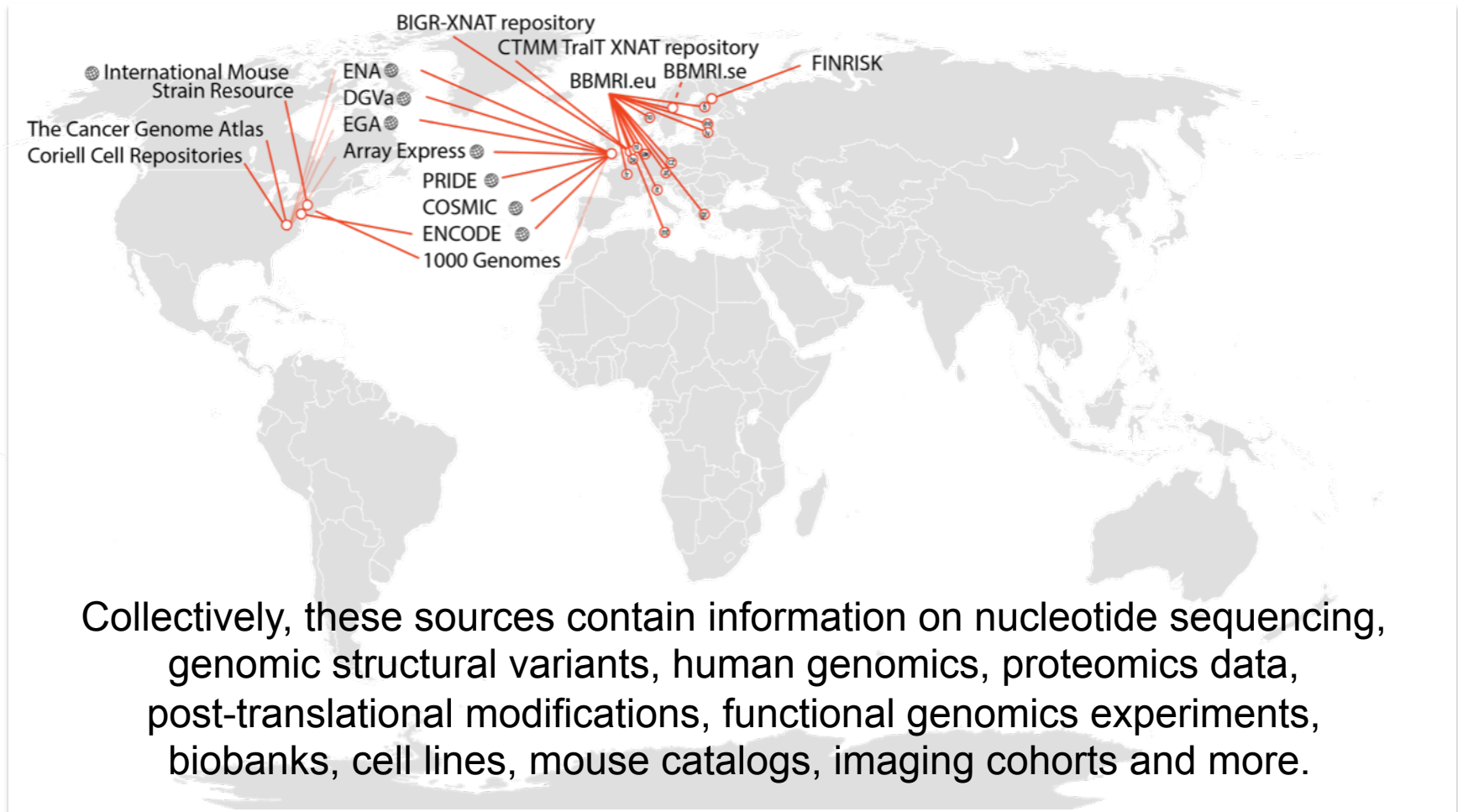- Simple spreadsheet population

- Domain neutral



http://bioregistry.cbs.dtu.dk/

# *Service and Resource Registry*

# BMB: Biosample information integration and discovery



Data deposition

Nucleotides, Transcripts, Proteins, Complexes, Domains, Structures, Small molecules, Pathways

Corresponding samples

Organisms, Groups, Parts, Cell lines, Environmental isolates, Micro organisms

http://www.ebi.ac.uk/biosamples/

Collectively, these sources contain information on nucleotide sequencing, genomic structural variants, human genomics, proteomics data, post-translational modifications, functional genomics experiments, biobanks, cell lines, mouse catalogs, imaging cohorts and more.

# Example query: "Cardiac Arrhythmia" yields over 450 samples from 10 groups (via BBMRI, Array Express, ENA)

## Search results

BBMRI.eu : Atrial Fibrillation Network Munich - M4-Cluster-Biobank

BBMRI.eu : Atrial Fibrillation Network Munich

BBMRI.eu : Atrial Fibrillation Network Munich

12 Homo sapiens samples from ENA SRA

Transcription profiling of mouse rapidly stimulated atrial myocytes: Conserv

Transcription profiling of mouse model of cardiac failure - particulate matter

Transcription profiling of human atrial and ventricular myocardium from pati
ventricular non-failing myocardium to identify the transcriptional basis for ul

Gender dependent differences in molecular electrophysiological targets in fa

NHLBI GO-ESP: Family Studies (Familial Atrial Fibrillation)

Using iPSC-derived neurons to uncover cellular phenotypes associated with

Circulating microRNAs to predict neurological outcome after sudden cardiac

Valvular heart disease and atrial fibrillation regulate microRNA expression profiles in left and right atria differently

LmnaN195K Mouse Model

Gene profiling of Hand2 target genes and transcriptional regulation of Hand2 expression in the postnatal myocardiu

Transcription profiling of rat heart transplants from Lewis to Lewis and Lewis to F344 strains with and without cold s
ischemia.

Gene expression analysis of cardiac left-ventricle tissue from hybrid mice harboring the Scn5a-1798insD/+ mutatio

Rac1-Induced Connective Tissue Growth Factor regulates Connexin 43 and N-Cadherin Expression in Atrial Fibrillati

Molecular Remodeling of Ion Channels in Human Atrial and Ventricular Myocytes Associated with Ischemic Cardiom

cardiac a

**cardiac a**rrhythmia
- Brugada syndrome
- ventricular fibrillation
- Familial short QT syndromw
- **atrial fibrillation**
- Catecholaminergic polymorphic ventricular
- Familial long QT syndrome
  - Timothy syndrome
  - Jervell and Lange-Nielsen syndrome
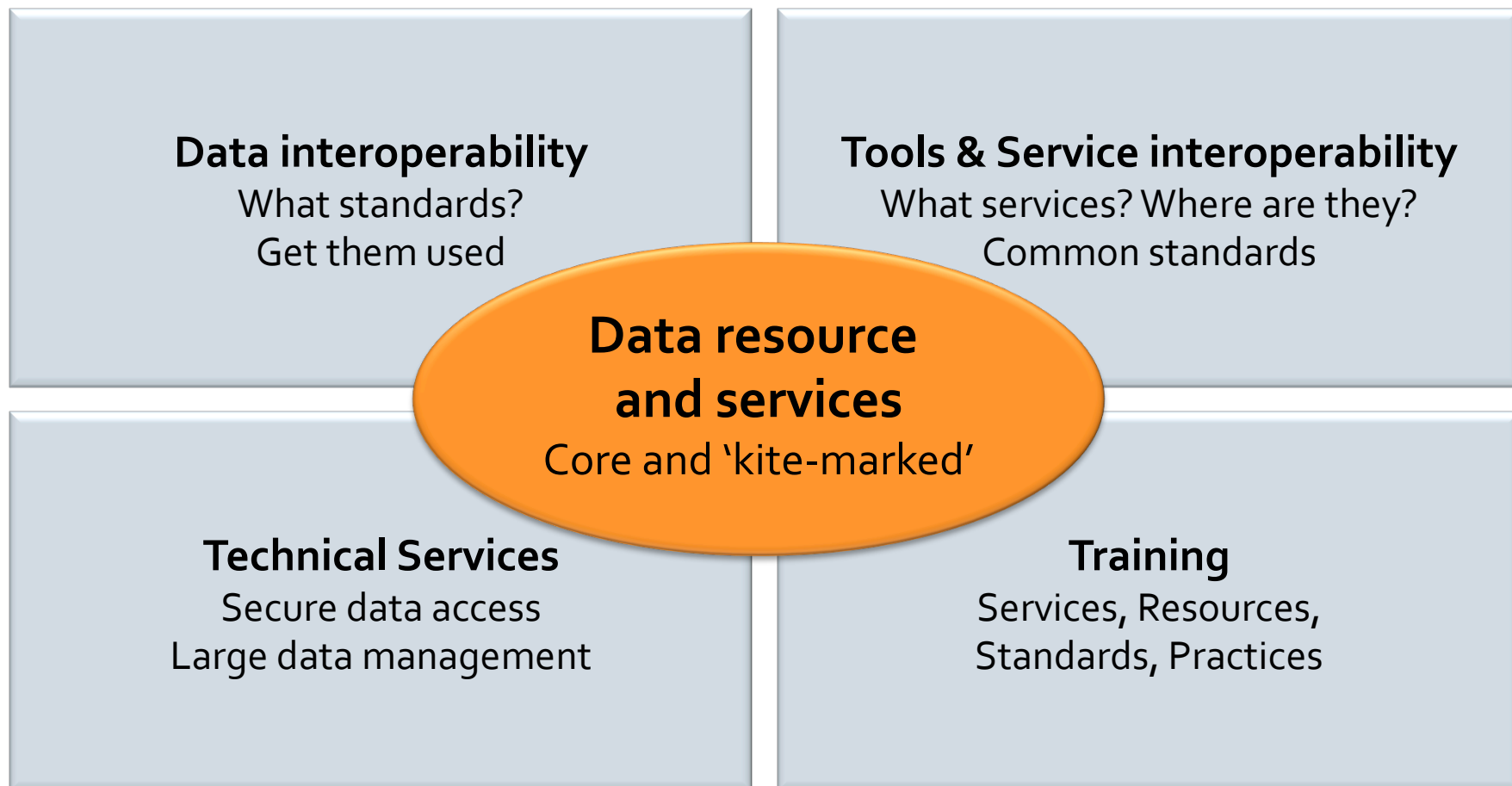  - Romano-Ward Syndrome
- acquired long QT syndrome
- sudden **cardiac a**rrest
- atrial flutter

**cardiac a**trium

# *Thank you*



Belgium    Czech Republic    Denmark    EMBL    Estonia

Finland    France    Greece    Israel    Italy

Netherlands    Norway    Portugal    Slovenia    Spain

Sweden    Switzerland    United Kingdom

The molecules of life

# Embassy Cloud logical view

# Data-focused Work Stream Headlines

**Data interoperability**
What standards?
Get them used

**Service interoperabilitry**
What services?
Common standards

**Data resource
and services**
Core and 'kite-marked'

**Technical Services**
Secure data access
Large data management

**Training**
Services, Resources,
Standards, Practices

Data infrastructure for Europe's life science research sector

elixir

# Data-focused Work Stream Headlines

**Data interoperability**
What standards?
Get them used

**Service interoperabilitry**
What services?
Common standards

**Data resource
and services**
Core and 'kite-marked'

**Technical Services**
Secure data access
Large data management

**Training**
Services, Resources,
Standards, Practices

Data infrastructure for Europe's life science research sector

elixir