# European Genome phenome Archive at the European Bioinformatics Institute

Helen Parkinson
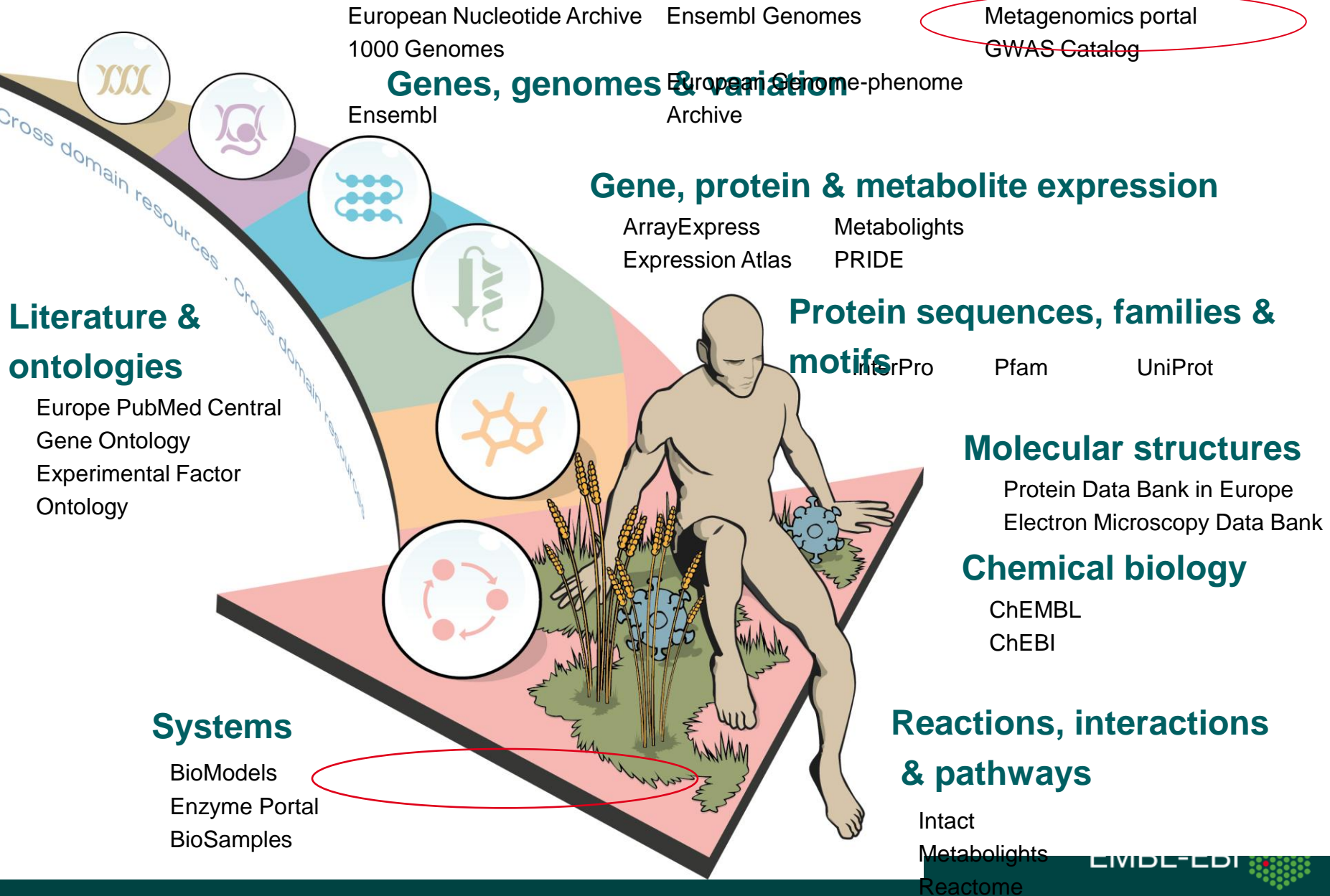
Head of Molecular Archives

EMBL-EBI

# What is EMBL-EBI?

- International, non-profit research institute

- Part of the European Molecular Biology Laboratory

- Europe's hub for biological data services and research

- ~600 members of staff from 53 nations.



*70 Petabytes of storage*
*>40,000 CPU Cores*
*20 Gbit Internet connection*

EMBL-EBI

# Data resources at EMBL-EBI

European Nucleotide Archive    Ensembl Genomes    Metagenomics portal
1000 Genomes                                       GWAS Catalog

## Genes, genomes & variation
European Genome-phenome

Ensembl    Archive

## Gene, protein & metabolite expression

ArrayExpress    Metabolights
Expression Atlas    PRIDE

## Literature & ontologies

Europe PubMed Central
Gene Ontology
Experimental Factor
Ontology

## Protein sequences, families & motifs

InterPro    Pfam    UniProt

## Molecular structures

Protein Data Bank in Europe
Electron Microscopy Data Bank

## Chemical biology

ChEMBL
ChEBI

## Systems

BioModels
Enzyme Portal
BioSamples

## Reactions, interactions & pathways

Intact
Metabolights
Reactome

Cross domain resources · Cross domain resources

EMBL-EBI

# EBI as a use case for 'data linkage'

… there are difficulties in defining boundaries around what is technically possible, what is legally permitted, and what should be done ethically. Indeed, linking data can also compromise the privacy of an individual's personal information, making them identifiable …
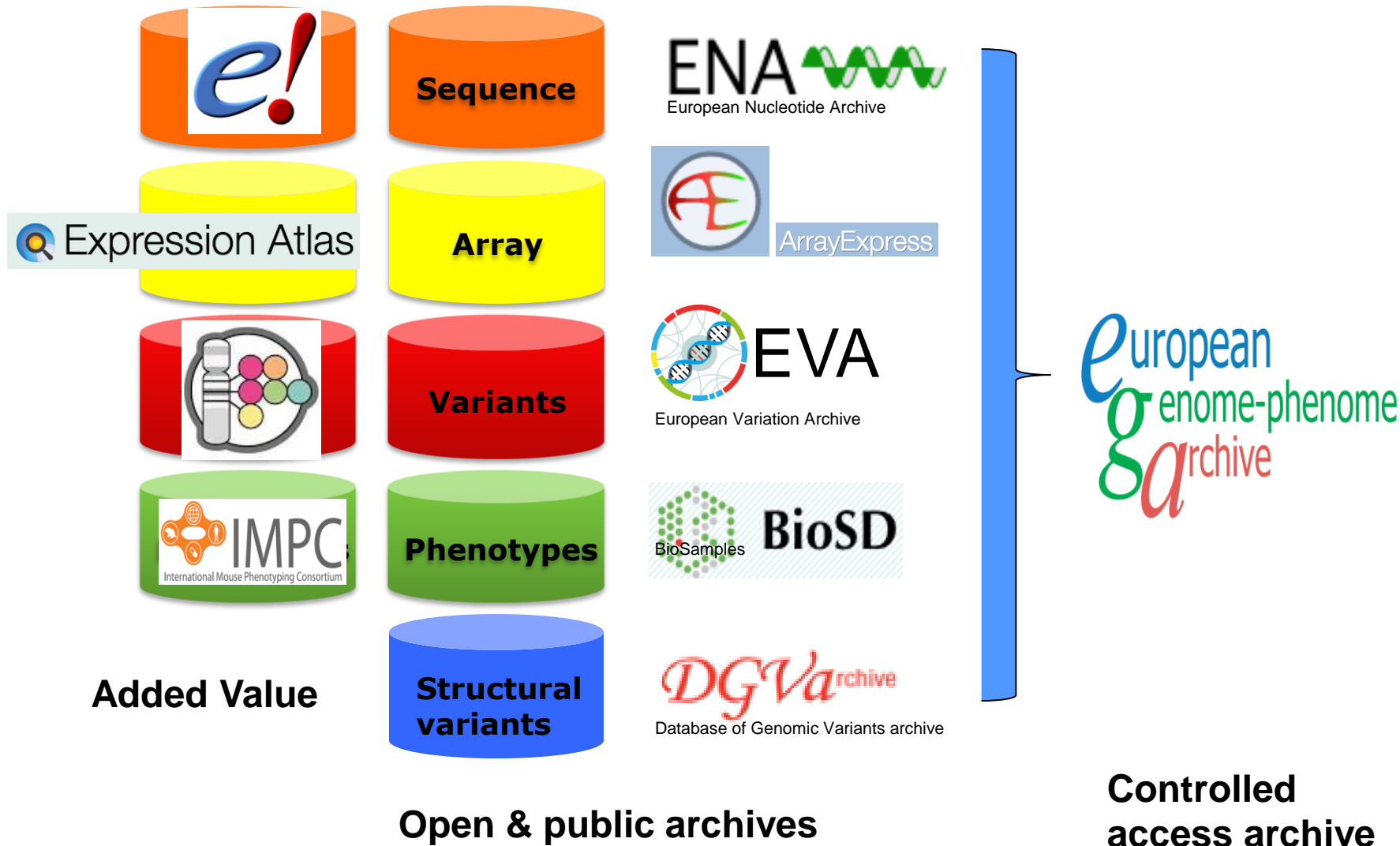
Use case
Ethical issues
Technical requirements
*balancing the drive for automation against the need for human judgement*

# EBI's Data Sources



**Sequence**

**Array**

**Variants**

**Phenotypes**

**Structural variants**

**Added Value**

ENA — European Nucleotide Archive

ArrayExpress

EVA — European Variation Archive

BioSamples BioSD

DGVarchive — Database of Genomic Variants archive

european genome-phenome archive

**Open & public archives**

**Controlled access archive**

EMBL-EBI

# BioSamples Database at the EBI

- Repository of information about biological materials, or "samples" – 4.3 million samples, 2 million human

- Aggregates reference samples – such as HapMap, Coriell cell lines or samples from 1000 genomes – with sample data from EBI assay databases

- Enables cross linking of assay data to metadata about the samples from which they are derived

- Accepts direct submission of sample data

BioSamples results for *1000 genomes*

Show more data from EMBL-EBI

| Page 1 .. 3 4 **5** 6 7 .. 34 | | Showing **101 - 125** of **842** Samples | | | Page size **25** 50 100 250 500 |
|---|---|---|---|---|---|
| Accession | Organism | Name | Description | Groups | Database |
| SAMN00779998 | Homo sapiens | Human 1000 genomes individual HG02605 | | | |
| SAMN00779989 | Homo sapiens | Human 1000 genomes individual HG02723 | | | |
| | | Human 1000 genomes | | | |

EMBL-EBI

# Summary vs. consented data

# What is the EGA?

The EGA is a resource for permanent secure archiving and sharing of all types of potentially identifiable genetic and phenotypic data resulting from biomedical research projects..



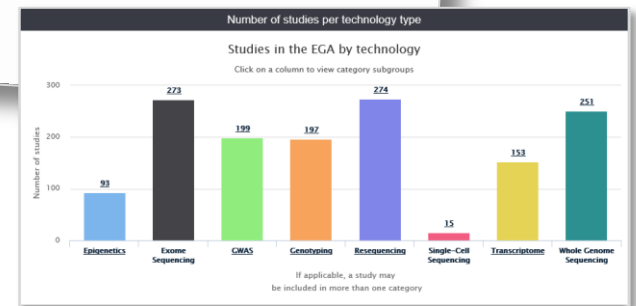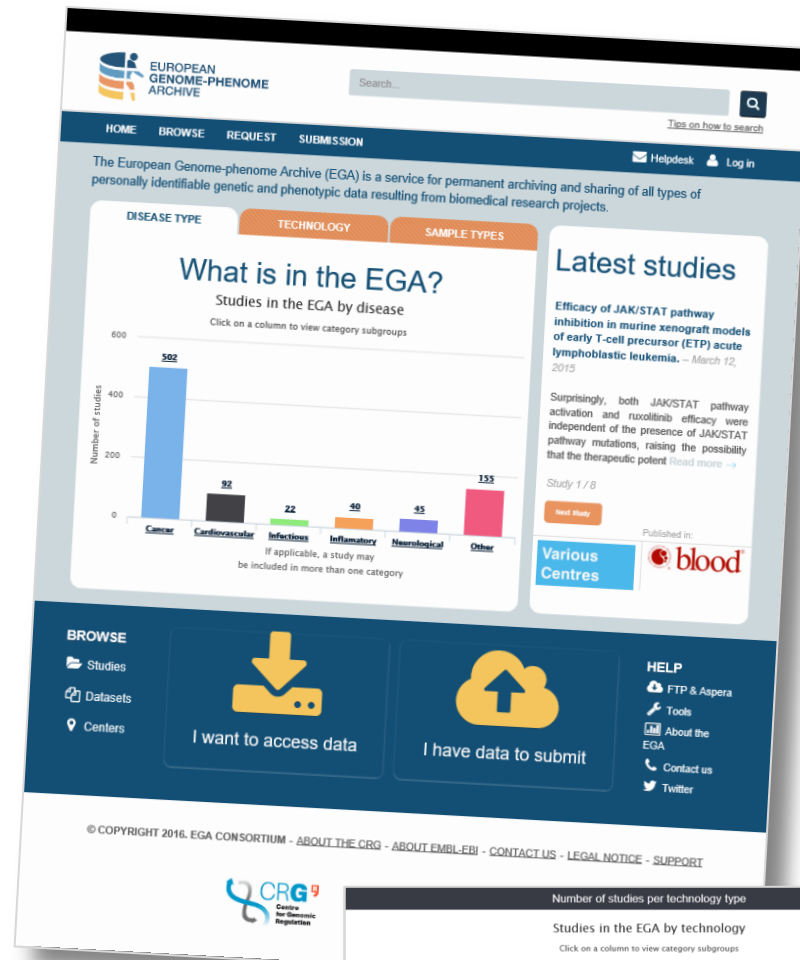Data is provided by research centers and health care institutions.

Access is controlled by Data Access Committees.

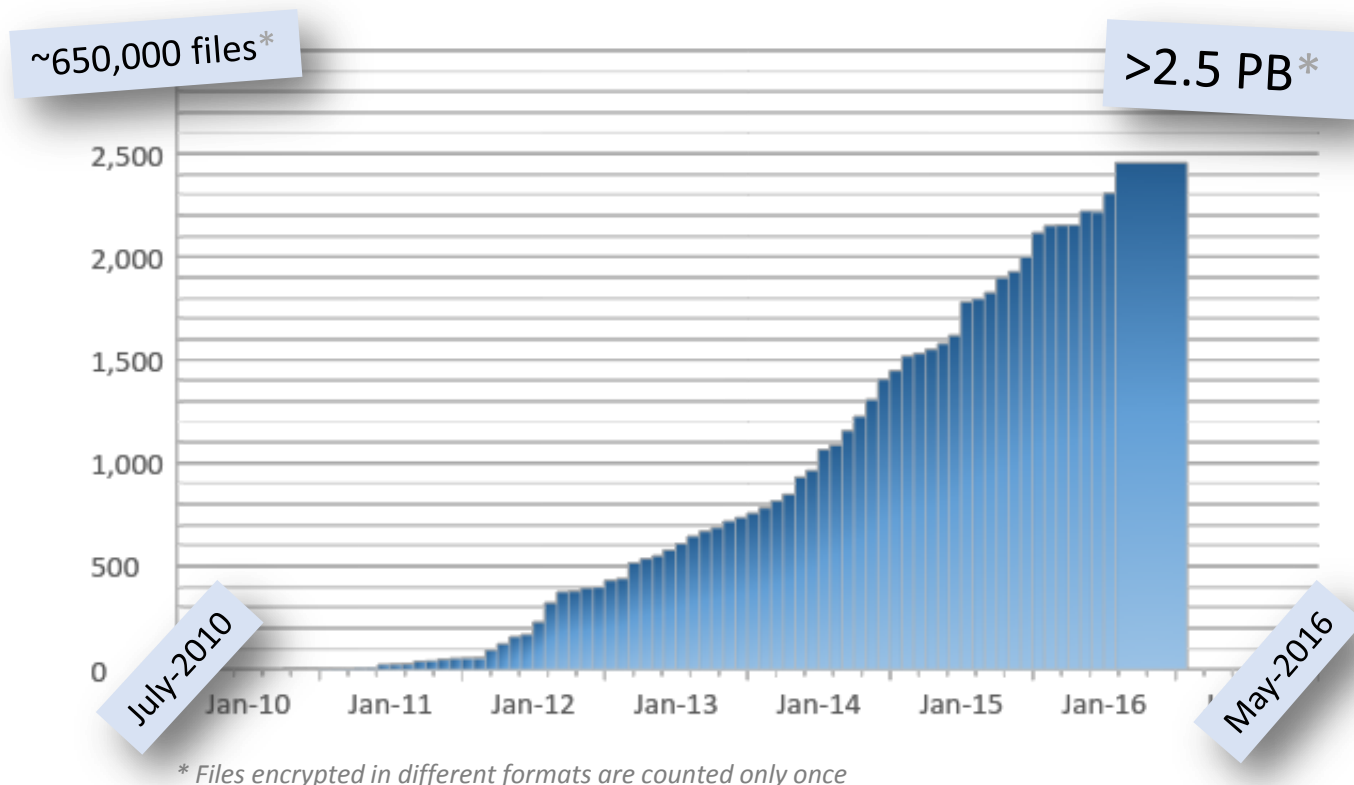Data requesters are researchers from other research or health care institutions.

Data and services and Federated between CRG and EBI

# EGA Numbers

- 911 Studies (31% growth in 2015)

- 1,966 Datasets (46% growth)

- 492 Data providers (92% growth)

- >8,000 Data requesters (20% growth year on year)

- 200 helpdesk tickets per month

# The EGA contains a growing amount of data



~650,000 files*

>2.5 PB*

July-2010

May-2016

* Files encrypted in different formats are counted only once

EMBL-EBI

# Data Sharing and Access

- EGA operates at two sites  EBI and CRG, Barcelona

- All data hosted at EBI – highly accessed data stored at CRG

- Study level meta data is shared between dbGAP (USA) and in future DDBJ

- Study level data is shared in the BD2K funded Data Discovery Index

- GA4GH Beacon provides variant le

Since January 1, 2015

- 1921245 files distributed

- 4.5 PB in Data Volume distribu

# Submitter accounts

**Data encryption & uploading**

**Metadata submission**

**Submitter account**

**Datasets**

Requirements:
- Appropriate affiliated consent agreement
- Submission statement
    - List of valid and institutional mail accounts
    - PI physical signature

Formed by:
- FTP box (10 Tb)
- Metadata submission credentials

# DAC Demographics

- 322 DACs
- Oversight of some 1966 datasets

- Created 8056 individual user accounts

- Our biggest DAC is Wellcome Trust Sanger Institute – which has 664 datasets

-  WTSI has 3 DACs in top ten with4 K users accessing them

- DACs creating >5 accounts per month provided with
    secureID key

EMBL-EBI

# Consent Codes

**Table 1. Data use categories and requirements (Consent Codes): definition and abbreviation.**

**Consent Codes**

| Name | Abbreviation | Description |
| --- | --- | --- |
| Primary Categories (I$^{ry}$) | | |
| No restrictions | NRES | No restrictions on data use. |
| General research use and clinical care | GRU(CC) | For health/medical/biomedical purposes, including the study of population origins or a |
| Health/medical/biomedical research and clinical care | HMB(CC) | Use of the data is limited to health/medical/biomedical purposes; does not include the population origins or ancestry. |
| Disease-specific research and clinical care | DS-[XX](CC) | Use of the data must be related to [disease]. |
| Population origins/ancestry research | POA | Use of the data is limited to the study of population origins or ancestry. |
| Secondary Categories (II$^{ry}$) (can be one or more extra conditions, in addition to I$^{ry}$ category) | | |
| Oher research-specific restrictions | RS-[XX] | Use of the data is limited to studies of [research type] (e.g., pediatric research). |
| Research use only | RUO | Use of data is limited to research purposes (e.g., does not include its use in clinical ca |
| No "general methods" research | NMDS | Use of the data includes methods development research (e.g., development of softwa algorithms) ONLY within the bounds of other data use limitations. |
| Genetic studies only | GSO | Use of the data is limited to genetic studies only (i.e., no "phenotype-only" research). |
| Requirements | | |
| Not-for-profit use only | NPU | Use of the data is limited to not-for-profit organizations. |
| Publication required | PUB | Requestor agrees to make results of studies using the data available to the larger scie community. |
| Collaboration required | COL-[XX] | Requestor must agree to collaboration with the primary study investigator(s). |
| Ethics approval required | IRB | Requestor must provide documentation of local IRB/REC approval. |
| Geographical restrictions | GS-[XX] | Use of the data is limited to within [geographic region]. |

# Federated 'Local' EGA



*"Offers solutions if data cannot leave the submitter facilities"*

# Local EGA

- Docker based Virtual Machine (infrastructure independent)
  - Database container
  - Data Archive container
  - API Access container
  - FTP server container
- Allows FTP upload of files, re-encryption, archiving, and distribution via downloader
- https://github.com/elixir-europe/human-data-local-ega

EMBL-EBI

# EGA Future

- Consent codes at data acquisition – transparency of consent

- EGA Local – for use as a country level EGA

- Shared meta data back to core for data discovery

- AAI Access – authentication, access, identification provides access to multiple resources with a common user id and AAI profile

- EGA in the cloud

- Addressing DAC scaling issues

- Improved usability of data – visualisation/slicing/query/input to analysis tools

- Improved streaming technology via GA4GH

- Improved submission tools  - efficiency and helpdesk, content

EMBL-EBI

# The EGA Beacon

- GA4GH Beacons are a discovery service:

  - which datasets include genomes with allele of interest?

  - 60 Beacons, 160 datasets, 25 organisation

- ELIXIR Objectives:

  - Provide ELIXIR reference implementation

  - Provide an example on capacity build across ELIXIR Nodes

- ELIXIR pilot project with partners from the Netherlands, Sweden, Finland, France, Belgium



https://ega.crg.eu/beacon_web/#/

# Beacon Challenges

Privacy Risks from Genomic Data-Sharing Beacons

Suyash S. Shringarpure ✉, Carlos D. Bustamante ✉

**IDASH PRIVACY & SECURITY WORKSHOP 2016**

HOME    ABOUT    **COMPETITION TASKS**

**SUBMIT HERE**

## Three tracks of competition tasks

Track 1: Practical Protection of Genomic Data Sharing through Beacon Services (privacy-preserving output rel...

Given a sample Beacon database, we will ask participating team to develop solutions to mitigate the Bustamante a...
each algorithm based on the maximum number of correct queries that it can respond before any individual can be
Bustamante attack. (data link)

EMBL-EBI

# Summary

EBI has data linkage use cases for both open and controlled access data

Likely that secure and summary data access mode will persist

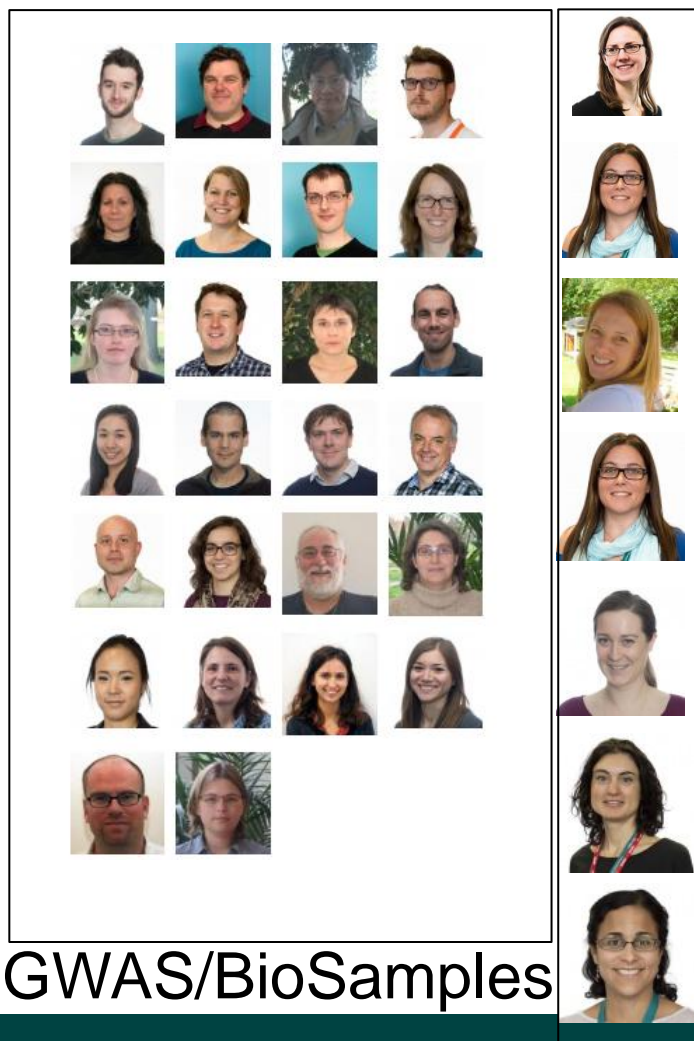Ethical issues  - largely determine the shape of our implementation

Technical requirements – federation, automation of access control, AAI, scaling

Ethics and policy around secure cloud sharing of EGA data

*Balancing the drive for automation against the need for human judgement*

EMBL-EBI

# *Acknowledgements*

## EGA @ EBI and CRG



**GWAS/BioSamples**

EMBL-EBI