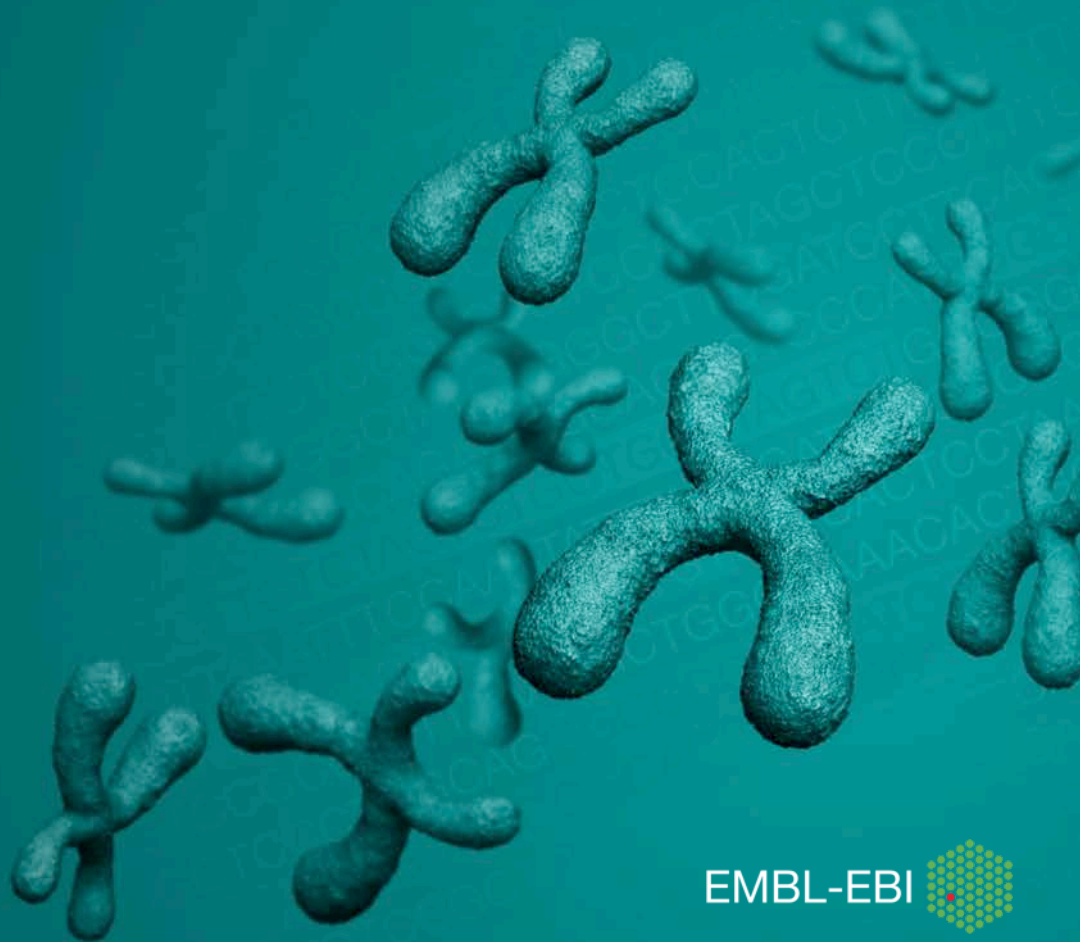# Sharing data to support Science

Laura Clarke

12th July 2016

EMBL-EBI

# OUR MISSION

To provide freely available data and bioinformatics services to all facets of the scientific community in ways that promote scientific progress

EMBL-EBI

# Interpreting human biology

How and why do we differ from one another?

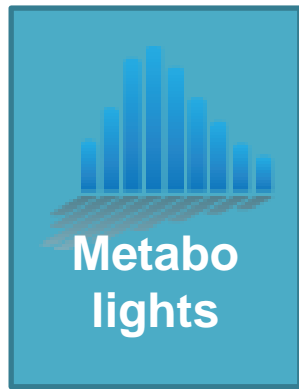How do you compare millions of genomes?

How important are lifestyle choices?

What causes susceptibility to disease?

What makes some people more sensitive to drugs?

EMBL-EBI

# Data Archives

BioStudies

EVA
European Variation Archive

Metabo lights

PRIDE

DGVarchive

Array Express

EUROPEAN GENOME-PHENOME ARCHIVE

ENA
European Nucleotide Archive

BioSD

EMBL-EBI

# Added Value Services

 e!Ensembl

 Gene Expression atlas

 Ve!P

 GWAS Catalog

 REACTOME

 Open Targets

 Locus·Reference·Genomic

 Europe PMC

 UniProt

 EMBL-EBI

# Sharing Good Data

EMBL-EBI

# Good data is well described data

- Needs
  - Well structured
  - Consistent naming
  - Specific descriptions
- Enables
  - Aggregation
  - Integration
  - Tracking

# Motivations



- Make your data usable
  - Reduce ambiguity
  - Facilitate reproduction of results
  - Improve integration across labs, projects and data modalities
- Make your data discoverable
  - Other researchers
  - Informatics services (Ensembl, Gene Expression Atlas)
- Improve your analysis
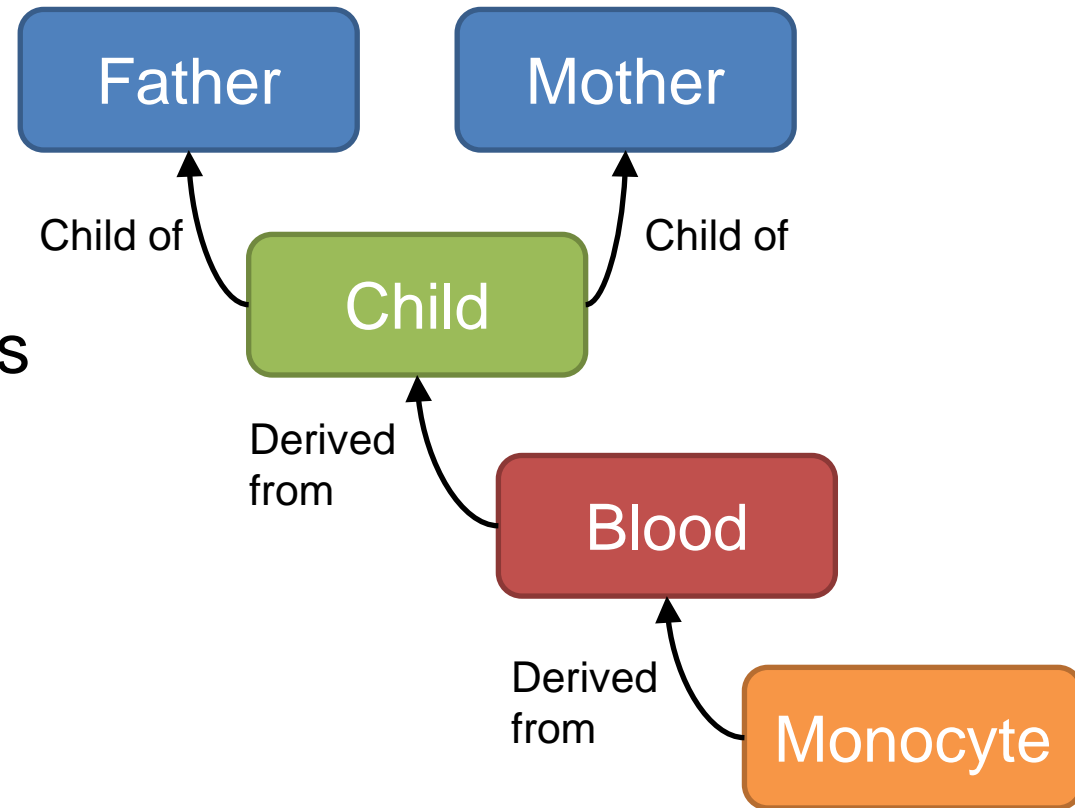  - Easier to find batch effects and confounding factors

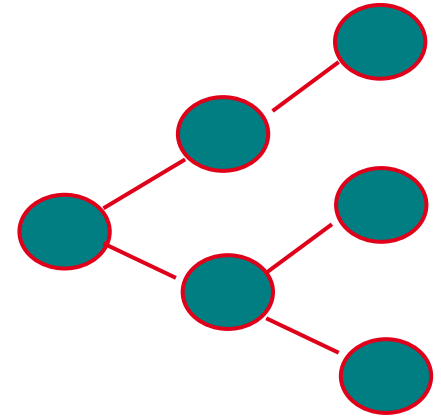# BioSamples

# BioSamples

- Provides
  - Unique Identifier
    - SAMEA3105765
  - User defined attributes
  - Ontology Support
  - Relationship support
    - Child of
    - Derived from
    - Same as

Father

Mother

Child of

Child of

Child

Derived from

Blood

Derived from

Monocyte

EMBL-EBI

# Ontology Support

- An ontology is

  - A classification of the kinds of entities that exist.

  - A specification of the meanings of terms in a conceptual vocabulary.

- Ontologies make it easier to

  - Understand what the description means

  - Enable more intuitive searching of the data

# Ontology driven search



https://www.targetvalidation.org

# Ontology Annotation
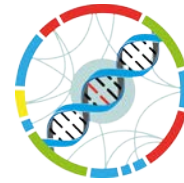
- Finding the right term can be challenging

- EMBL-EBI has tools to help

- Zooma

  - Using past knowledge to inform new annotation

  - http://www.ebi.ac.uk/spot/zooma/

- Ontology Lookup service (OLS)

  - Indexes 150 biomedical ontologies (4.5 million terms, 11 million relations)

  - Web and programmatic interfaces

  - http://www.ebi.ac.uk/ols

OLS

$ZOOMA^2$

EMBL-EBI

# Sharing Data

EMBL-EBI

# Data Archives

# Data Hubs

- ENA tool available to collaborators

- Allows consortium wide access to targeted prepublication and public data

- Data and metadata

- Supported data types

  - Sequence

  - Alignment

  - Variants

  - Assemblies

# Sharing controlled access data

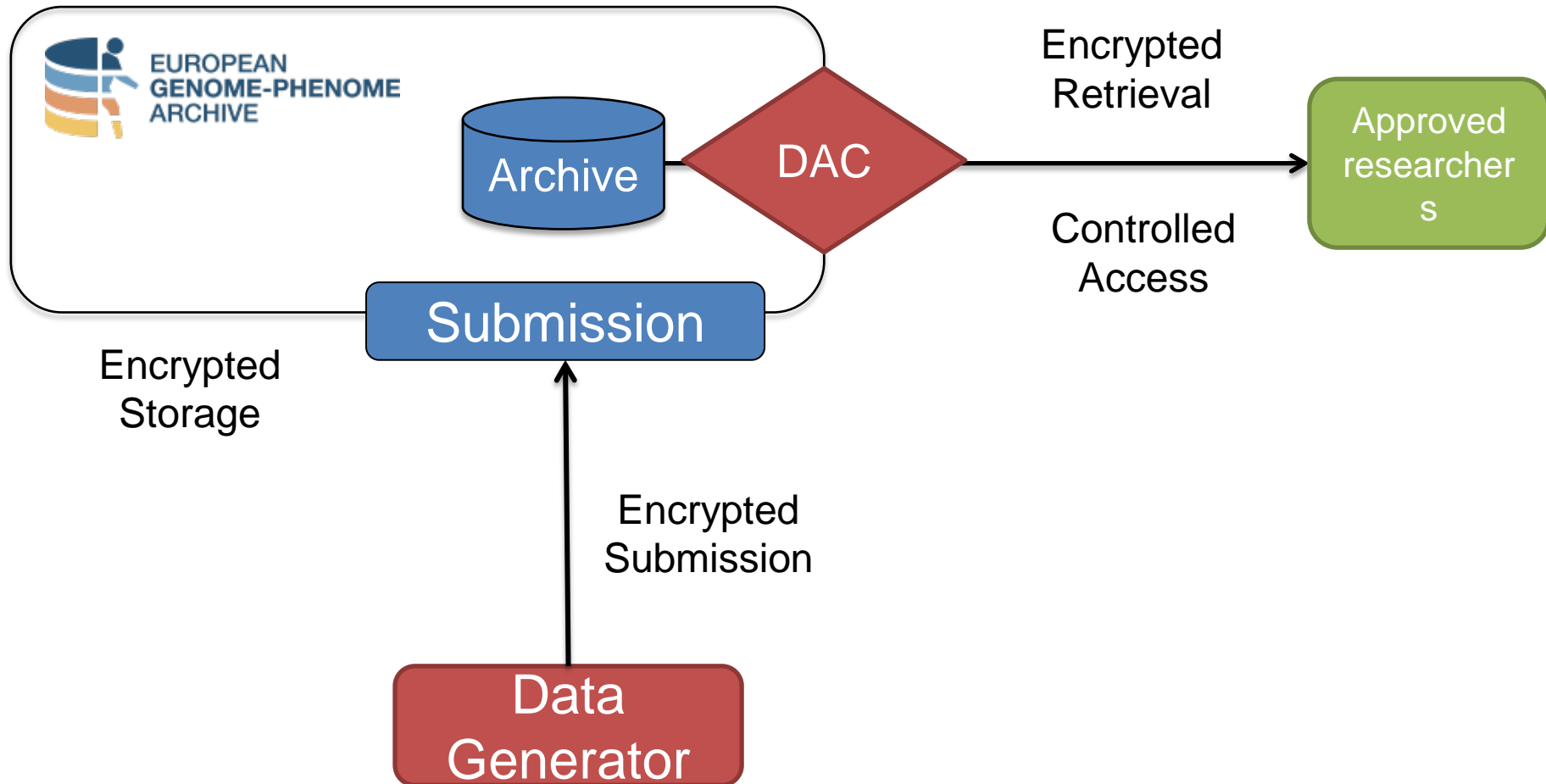EMBL-EBI

# What is controlled access data

- Driven by participant consent

- Some consents specify data is released

  - In a controlled manner

  - To bone fide researchers

- This data tends to be

  - Human, personally identifiable data types ('raw', processed, phenotypic)

  - Affiliated to bio-medical research or consortium projects
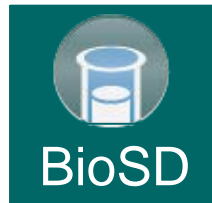
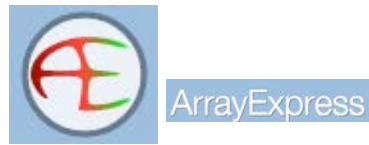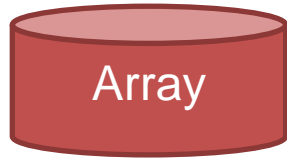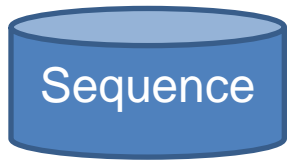EMBL-EBI

# What is the EGA?

- Launched 14th July 2008

- Role: **secure** archive for **controlled** distribution of **consented** genetic & phenotypic data

- **8000+** data access accounts; **395** submission accounts

- **~2.5PB** available for download; **~1900** datasets; **1.6PB** distributed last year

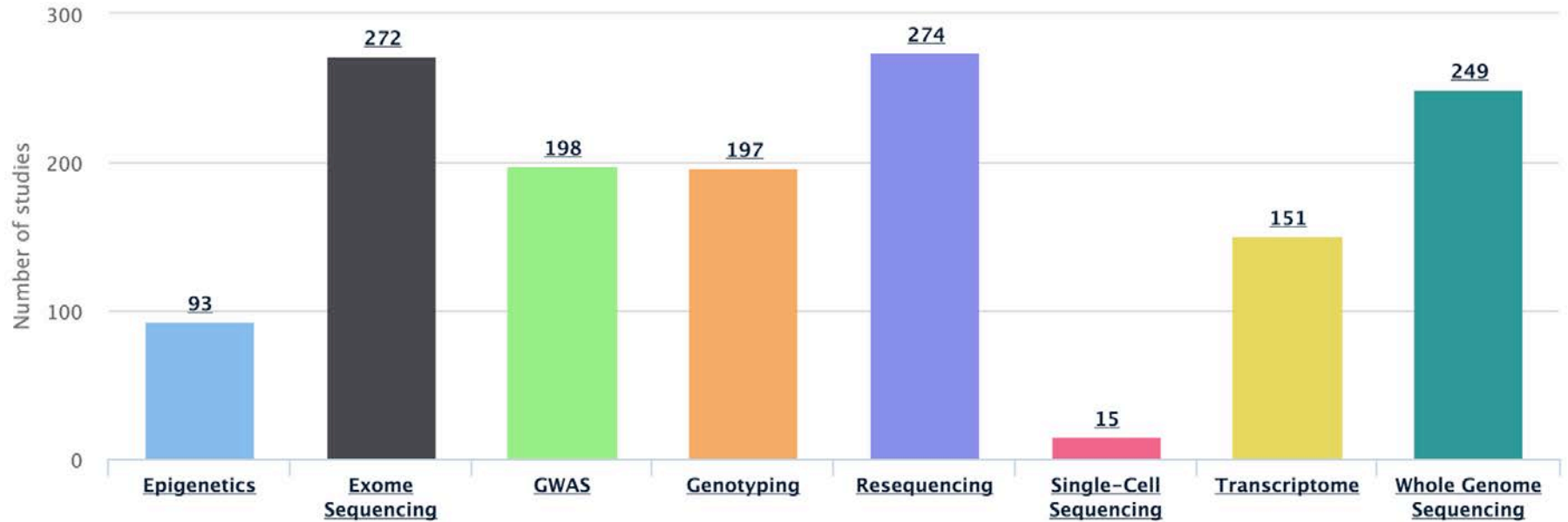- **+200** contacts to Helpdesk/month

# EGA - Architecture

# Where do I submit data?



Open & public archives

Controlled archives

EMBL-EBI

# EGA Studies by Technology

# Added Value Services

 e!Ensembl

 Gene Expression atlas

 Ve!P

 GWAS Catalog

 REACTOME

 Open Targets

 Locus·Reference·Genomic

 Europe PMC

 UniProt

# Putting you data in context

EMBL-EBI

# Ensembl

- Genome Browser
  - 65 vertebrate species
- Open Data
- Open Code
- High quality annotation
  - Genes
  - Regulatory regions
- Comparative Analysis
- Variation Data
- Data mining tools

# Trackhubs

- Publicly hosted text file

- Supports indexed "Big" formats

- See your data in context of broad genomic annotation

- Attach many files to Ensembl (or UCSC) at once

- Provide community with clearly defined data collections

# OUR MISSION

To provide freely available data and bioinformatics services to all facets of the scientific community in ways that promote scientific progress

EMBL-EBI

# Collaboration opportunities

- Data Coordination Services

- Specialized support for archive use

- Increased data integration with value added services

- New technology development



EMBL-EBI

# Thanks



EMBL-EBI

# Questions?

EMBL-EBI