# WHAT ARE OTHER EUROPEAN NETWORKS OFFERING? WHAT IS THE BENEFIT OF SHARING DATA AND SAMPLES THROUGH EXISTING STRUCTURES E.G. RD-CONNECT?

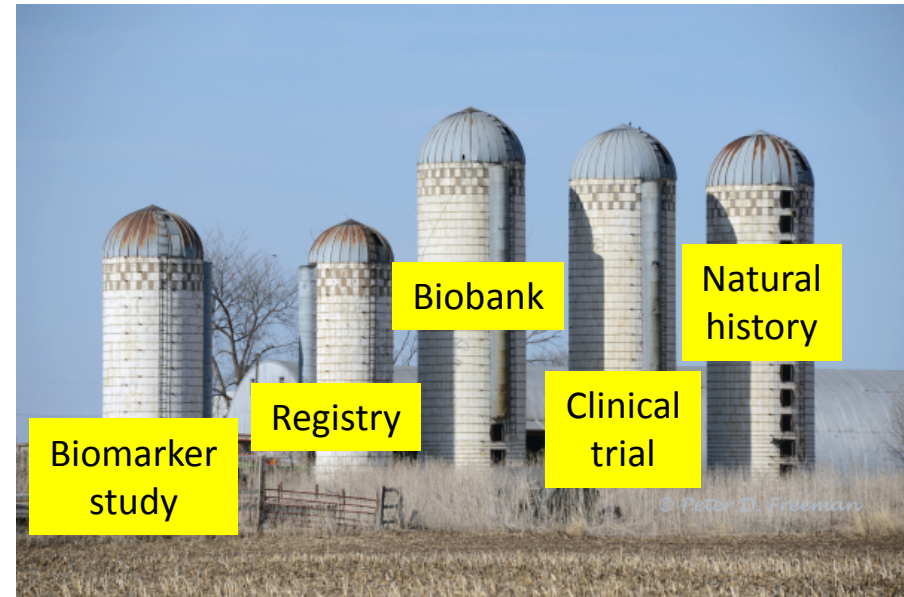Rachel Thompson • Newcastle University

**A: quite a lot!**

# The guiding principle

## Data sharing for research and better data analysis

- ☐ Gene and modifier discovery
- ☐ Samples for further research
- ☐ Genotype-phenotype correlation
- ☐ Patient recruitment
- ☐ Global natural history comparisons
- ☐ Biomarkers, therapeutic targets…



Biobank
Natural history
Registry
Clinical trial
Biomarker study

**Overcoming silos!**

**An integrated platform connecting databases, registries, biobanks and clinical bioinformatics for rare disease research**

**Overarching objectives:**

- Contribution to the IRDiRC objectives of delivering 200 new therapies for rare diseases and means to diagnose most rare diseases by the year 2020

- Development of an integrated, quality-assured and comprehensive platform in which complete clinical profiles are combined with -omics data and sample availability for rare disease research, in particular IRDiRC-funded research.

# RD-Connect's main aims

- Creation of central system and repository for **reprocessing**, **storing** and **analysing** omics data
  - Raw data hosted at European Genome-phenome Archive (EGA)
  - Raw data reprocessed through standard analysis pipeline for consistency
  - Reprocessed data accessible via Barcelona platform with user-friendly online analysis interface
- Integration of phenotypic data
- Integration of biosample data
- Development of new bioinformatic tools
- Ethical and legal considerations for data sharing
- Patient input
- Outreach and impact: interaction with rare disease community

# Sharing: What?

- Raw data from all types of studies
- Genomic data
- Phenotypic data
- Natural history data
- Clinical trial data
- Biosamples (blood, DNA, tissue samples, cell lines...)
- Linked data and samples
- Access to patients
- ...

# Sharing: Barriers

- **General**
  - Privacy protection issues: "do I have the patient's permission?"
  - Lack of infrastructure: "I want to share data but where do I put it?"
  - Lack of standards and interoperability
- **Academia**
  - Culture of protecting research results: "someone else might scoop my publication!"
  - Lack of incentives for sharing
- **Industry**
  - IP issues/competition (when pharma is asked to share its own data)
  - Concerns over data quality, regulatory compliance (when pharma wants to reuse data from academia)

# Sharing: Benefits

- Overcoming the "rare disease problem"
  - Cohort size
  - Powering trials
  - Finding confirmatory cases
- Reducing costs
- Reducing duplication of effort
- Facilitating validation of results
- Enabling engagement with experts and the patient community
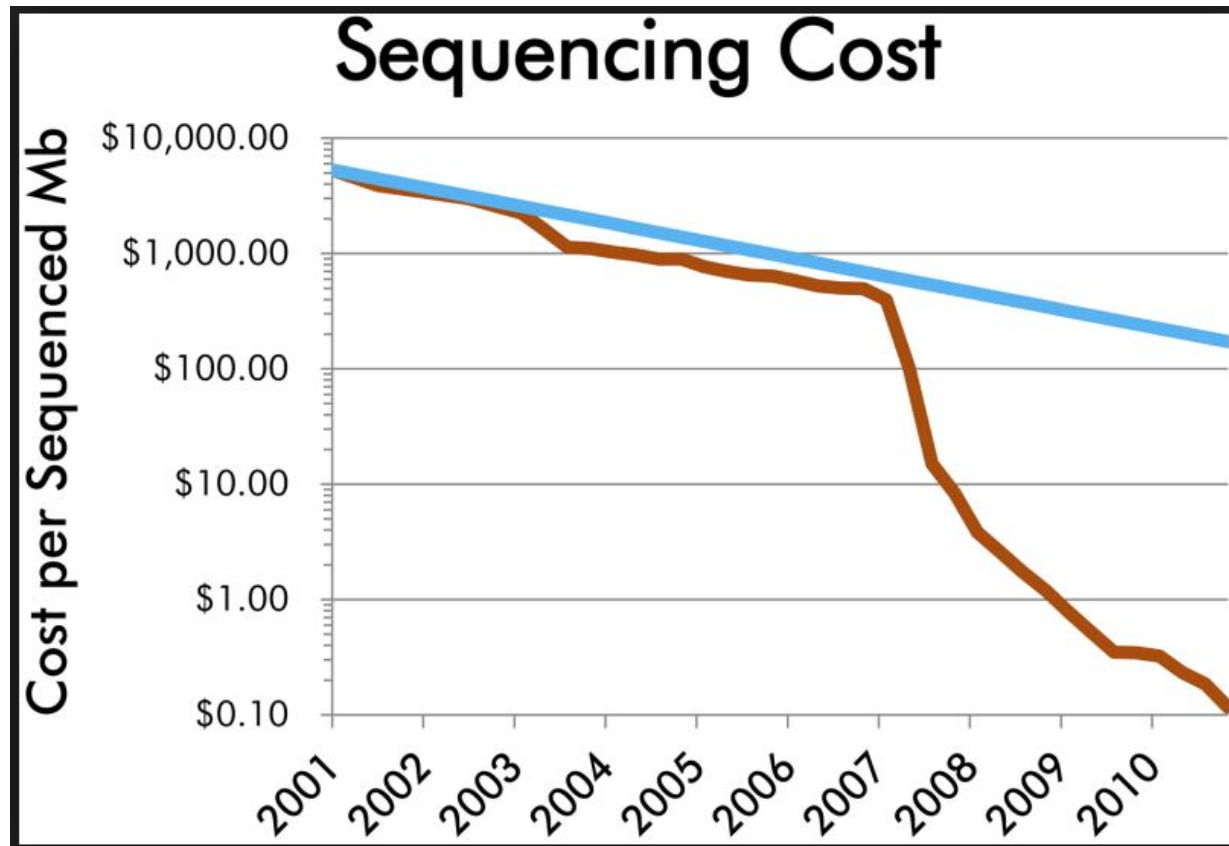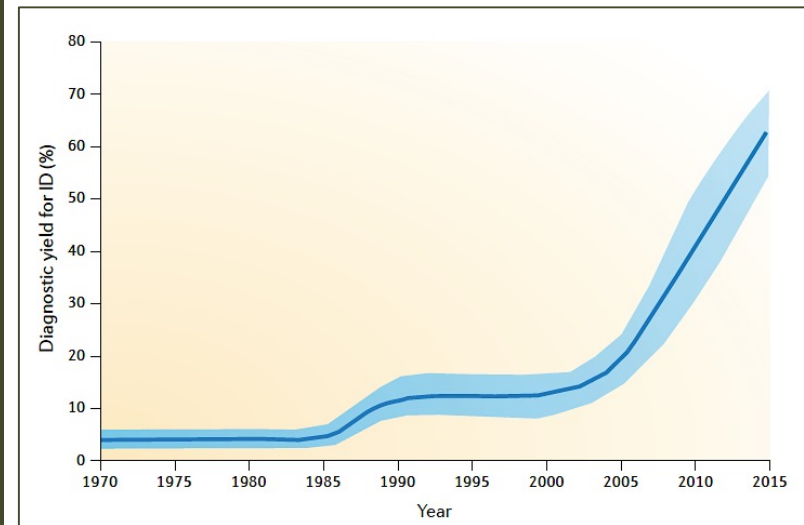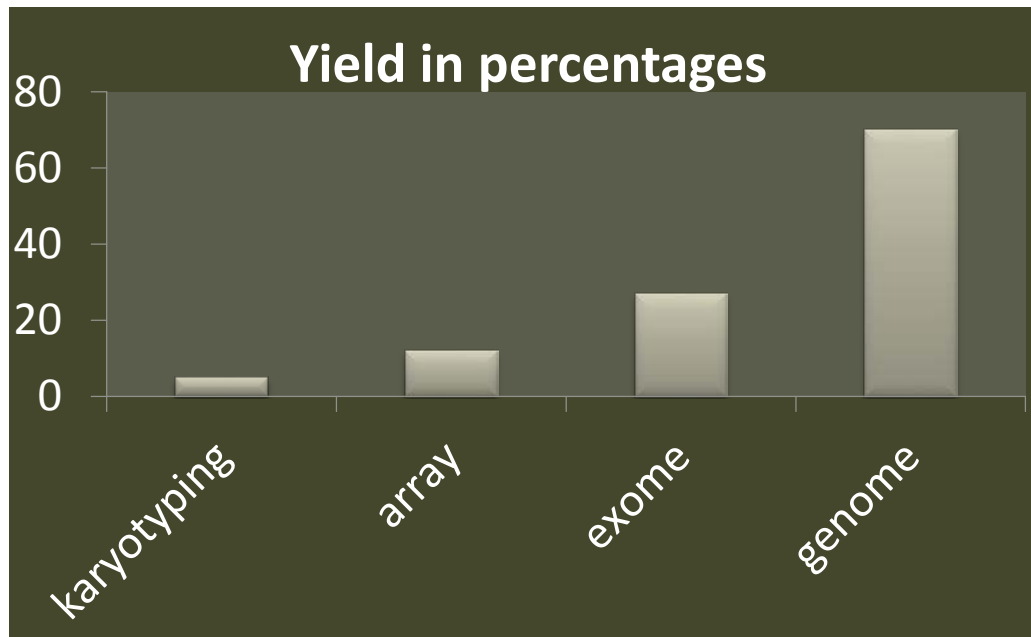
# Data integration in RD-Connect

**Sample data**
(biobank databases)

**Clinical data**
(registries, and phenotypic databases)

**Genomic data**
(WES, WGS)

**Other omics data**
(transcriptomics, metabolomics, proteomics …)

# NGS is becoming affordable

# Number of new genes discovered is increasing

□ Example: intellectual disability

**Yield in percentages**



Vissers et al., Nature Rev Genet 2016

# But: interpretation is still difficult

## Molecular diagnostics in NGS era

### Sample in → Diagnosis out?



# "black box"

# The challenge

## Interpretation of DNA variants: how do I find the pathogenic mutation?
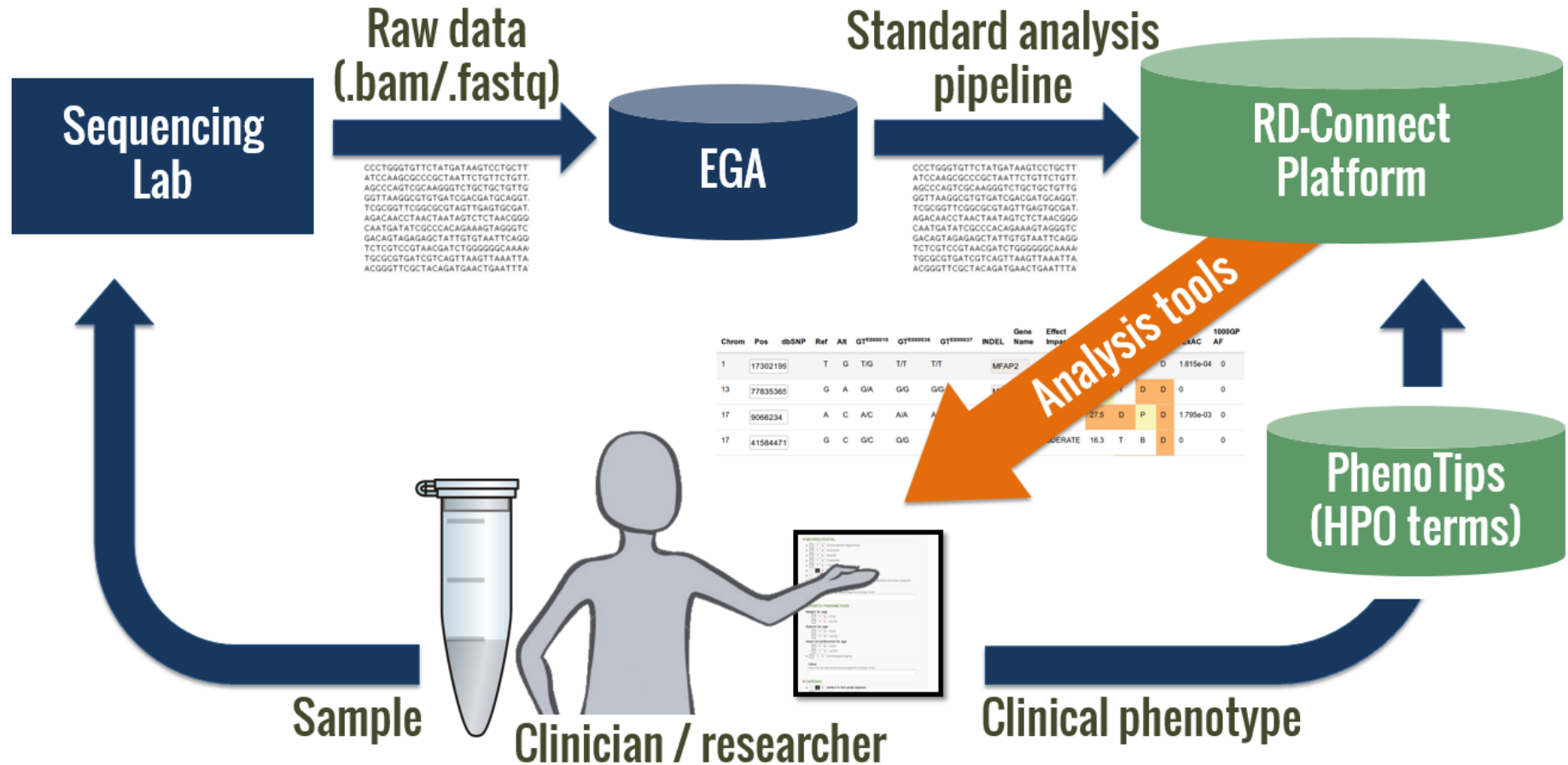
Exome sequencing →

25,000- 50,000 variants ←→ 1 pathogenic mutation

# Genomic data flow in RD-Connect

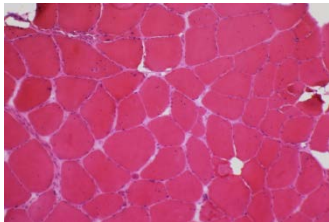# RD-Connect genomic analysis platform

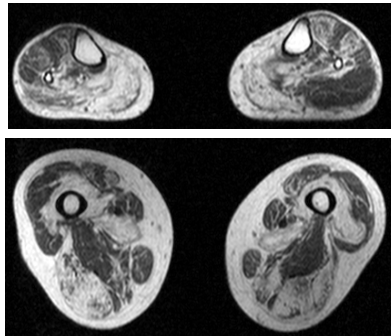# Exome sequencing and data sharing: new congenital myopathy gene

## Newcastle case

- Childhood onset

- Proximal muscle weakness, mainly lower limbs

- Slow progression

- CK: normal or mildly elevated

- Muscle biopsy: internal nuclei, fibre splitting and fibre type 1 predominance

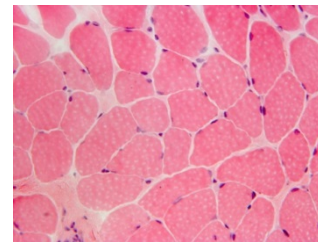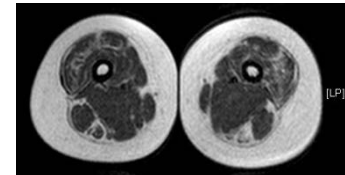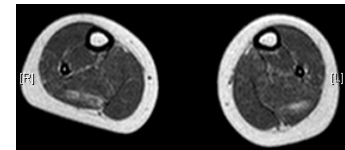- Pattern resembling DNM2 patients



51 years



65 years

## London case

- Antenatal onset with reduced foetal movement

- Proximal muscle weakness, mainly lower limbs

- Axial weakness

- Joint laxity of hands and ankles

- Slow improvement

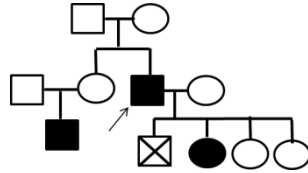- Muscle biopsy: minicores, central cores and some internal nuclei



4 years



4 years

RD Connect

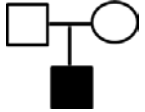# Exome sequencing and data sharing: new congenital myopathy gene

## Newcastle case

- Stop gain
- novel (absent from 62k)
- chrX:153049846 G>A; p.Trp415Ter

## London case

- Essential Splice Site
- novel (absent from 62k)
- chrX:153050629 G>A

## SRPK3

- Serine/arginine protein kinase
- Muscle specific, regulated by myocyte enhancer factor 2 (MEF2)
- Known to regulate mRNA splicing and nuclear lamina proteins
- KO mice develop centronuclear myopathy (Nakagawa et al 2005)
- Preliminary data in zebrafish morpholino knockdown shows slow movement and muscle disorganization (unpublished)
- Four new mutations found (manuscript in preparation)

Ana Töpf
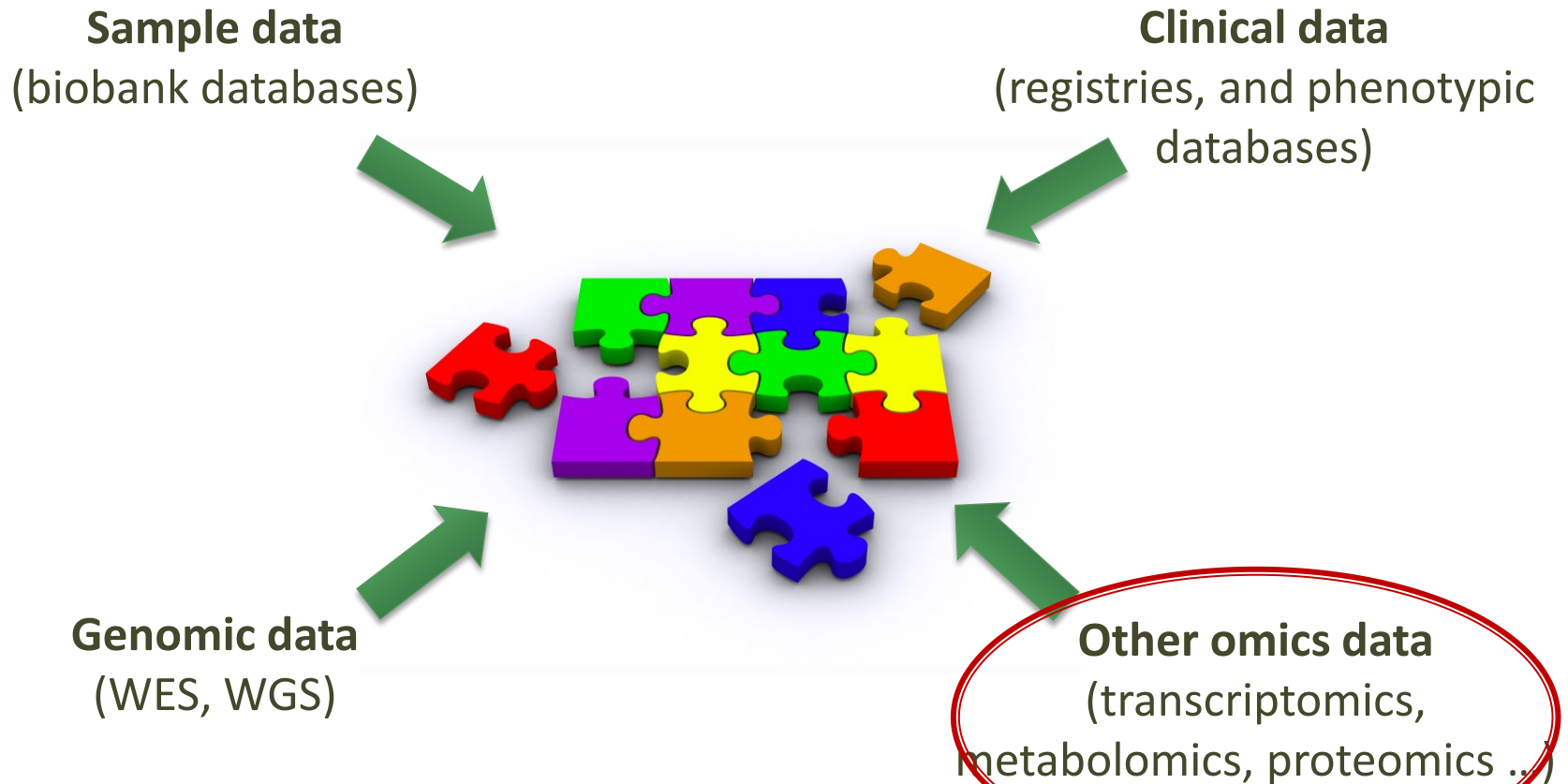
**RD Connect**

# Data integration in RD-Connect

**Sample data**
(biobank databases)

**Clinical data**
(registries, and phenotypic databases)

**Genomic data**
(WES, WGS)

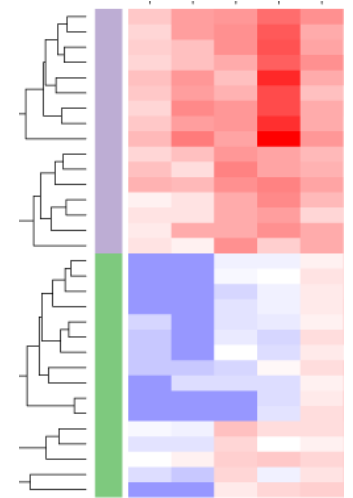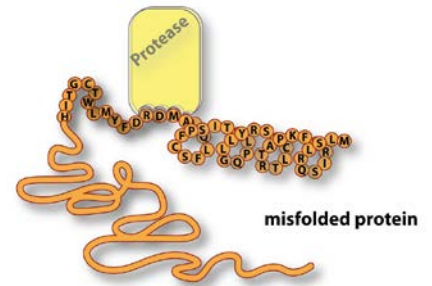**Other omics data**
(transcriptomics, metabolomics, proteomics ...)

# Other omics – work in progress
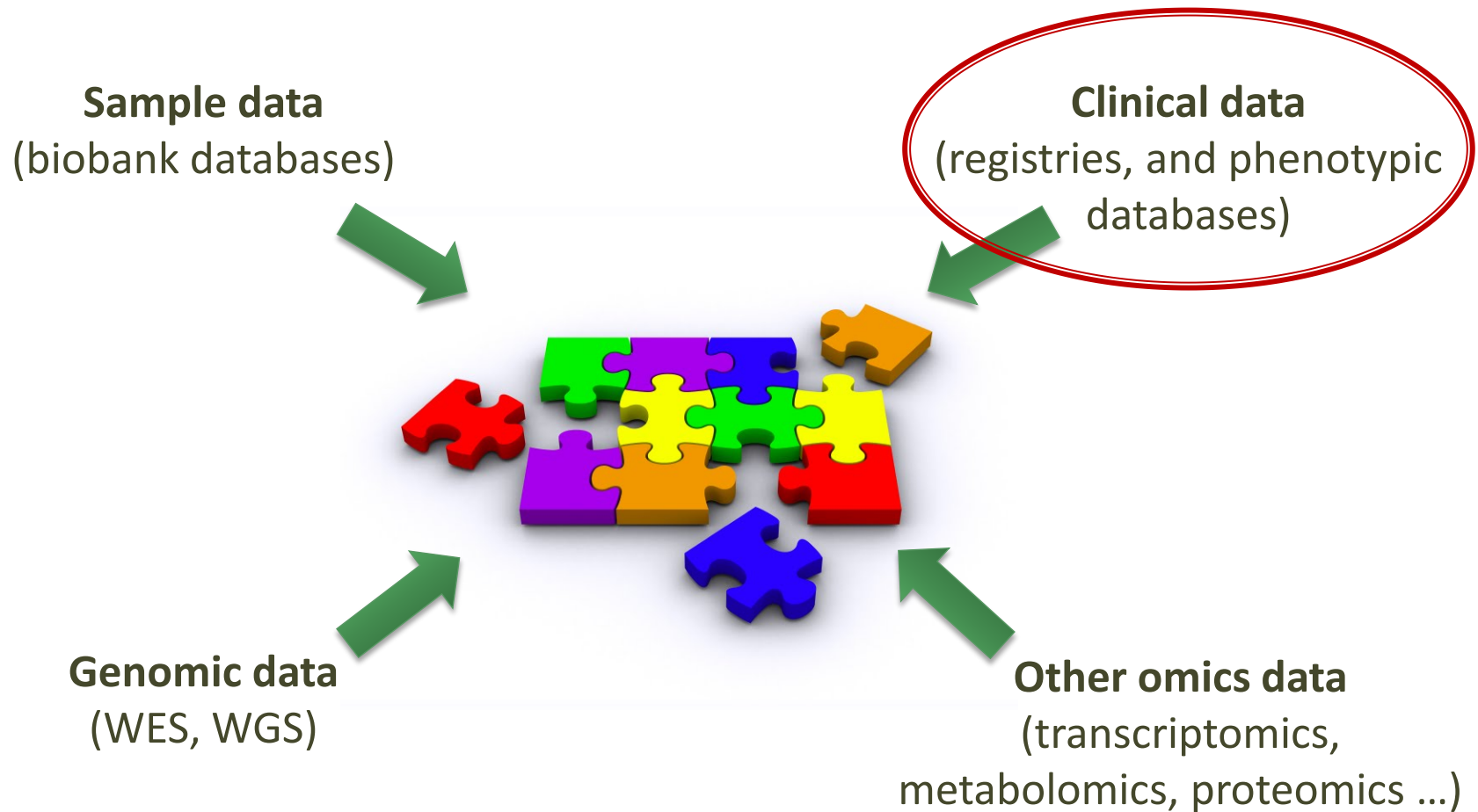
- Integration of other omics data types – transcriptomic, proteomic, lipidomic, metabolomic profiles – is a work in progress

- Challenges with standardization of data done on different machines/from different centres

- Need to work out the multi-omics research questions that people want to answer

- Integration on a per-patient level to allow comparison across data types

misfolded protein

# Data integration in RD-Connect

**Sample data**
(biobank databases)

**Clinical data**
(registries, and phenotypic databases)

**Genomic data**
(WES, WGS)

**Other omics data**
(transcriptomics, metabolomics, proteomics ...)

# Clinical and phenotypic data

- Phenotype is **more important than ever** in the context of clinical outcome measures and next-generation sequencing analysis

- Requires transformation into a "computable" form

- Requires linkage from different sources (multiple registries, phenotypic databases…)

  ➡️ **FAIR DATA**

# What is FAIR data?

**F**indable - (meta)data is uniquely and persistently identifiable. Should have basic machine readable descriptive metadata.

**A**ccessible - data is reachable and accessible by humans and machines using standard formats and protocols.

**I**nteroperable - (meta)data is machine readable and annotated with resolvable vocabularies/ontologies.

**R**eusable - (meta)data is sufficiently well-described to allow (semi)automated integration with other compatible data sources.



13 October 2016

# Common data elements

- Attempts to standardize elements collected in patient registries – an ongoing challenge!



NINDS **Common Data Elements**
Harmonizing Information. Streamlining Research.

EUCERD
Joint Action

**MINIMUM DATA SET FOR RARE DISEASE REGISTRIES**

**PARENT**
cross-border
PAtient REgistries iNiTiative

The **EPIRARE** proposal of a set of indicators and common data elements for the European platform for rare disease registration

Domenica Taruscio, Emanuela Mollo, Sabina Gainotti, Manuel Posada de la Paz, Fabrizio Bianchi, and Luciano Vittozzi

# Ontologies

Consensus on most useful ontologies in rare disease:

- Human Phenotype Ontology (HPO)
  - For phenotypic descriptions (observations)
- Orphanet Rare Disease Ontology (ORDO)
  - For "naming" a disease

Advantages of ontology use:

- Computers understand them
- Tree structure (if x is true then everything above x is also true)
- Allows computational analysis and matchmaking approaches

Figure 4 | **Example of a limb-girdle phenotype hierarchy from the Human Phenotype Ontology (HPO).**

# When data is not prepared for cross-resource analysis



| | C (USA) |
|---|---|
| Education level | C_EDUC: 7 levels |
| Marital status | C_MARSTAT: never, now, separated, divorced, divorced |
| Age/date of birth | Age at baseline in years |

| | R2 (EU) |
|---|---|
| Education level | Edlevel: 9 levels |
| Marital status | Maristat: single, married, partnership, divorced, widowed |
| Age/date of birth | Exact age at visit |

| | R3 (EU) |
|---|---|
| Education level | Isced: 7 levels |
| Marital status | Maristat: single, married, partnership, divorced, widowed |
| Age/date of birth | Exact age at visit |

When data is **not** linkable at the source

Back of the envelope calculation

□ 6 months per data set

□ Reuse: 5x on average, 6x5=30 Months

□ For every RD: 6000x(6x5) = 180000 M

When data is linkable at the source

- 6 months once

- Reuse: 5x on average, +1M, 1x5=5 M (30)

- For every RD: 6000x(6+1x5) = 66000 M

(180000)

**RD Connect**

# How much time do researchers spend on preparing data for integration

- Benefits for cross-resource analysis

- 66% efficiency gain (more time for research)

- Researchers can start analysing 6x faster

# Rare Disease Registry Framework

**29**

**Genetic Test Details**

Are details of genetic testing available `---`

Genetic Test Date

Has the patient received genetic counselling `---`

Has anyone in the patient's family received genetic counselling `---`

**Save**
**Cancel**
**Print**
**Next**
**Previous**

**Motor Function**

Currently able to walk ○ Yes ○ No

Current use of devices to assist with walking `---`

At what age did the patient commence using devices to assist with walking

Current best motor function `---`
*Walking: walking with or without help (orthoses or assistive device or human assistance), inside or outdoors*

Dysarthria `---`

# Patient Archive (HPO)

## Clinical Records

+ Add a Clinical Record

---

**CLINICAL_CONSULTATION**  **In Phenotype Profile**

Annotation Sufficiency ★ ★ ★ ★ ☆

Jun 2, 2016 12 hours ago

🏷 Start annotating  ✏ Edit  🗑 Delete

Unexplained left ventricular hypertrophy (LVH).
Occurs in non-dilated ventricle in the absence of other noticeable cardiac or systemic disease.
Shortness of breath (particularly with exertion), chest pain, palpitations, orthostasis, presyncope.
No syncope.
Recently developed symptoms.

**Chest pain (chest pain)**  **No Syncope (syncope)**  **Palpitations (with pheochromocytoma) (palpitations)**  **Ventriculomegaly (dilated ventricle)**
**Dyspnea (Shortness of breath)**  **Left ventricular hypertrophy (left ventricular hypertrophy)**  **+**

---

**CLINICAL_CONSULTATION**  **In Phenotype Profile**

Annotation Sufficiency ☆ ☆ ☆ ☆ ☆

Jun 2, 2016 12 hours ago

🏷 Start annotating  ✏ Edit  🗑 Delete

Noninvasive cardiac imaging using echocardiography but results are unclear.

**+**

**RD Connect**

# PhenoTips (HPO)

# Data integration

**Sample data**
(biobank databases)

**Clinical data**
(registries, and phenotypic databases)

**Genomic data**
(WES, WGS)

**Other omics data**
(transcriptomics, metabolomics, proteomics …)

RD Connect

# Biosample data

- **(1) Cataloguing** and registration of rare disease biobanks
  - Biobanks can sign up and give details of their biobank in an "ID card"
  - Allows biobanks to participate in RD-Connect infrastructure and research
  - Standardised assessment procedure, MTAs etc.

- **(2)** Sharing **sample-level data** in a common database
  - Not just sample numbers but drill-down right to individual samples
    - Researchers can find the samples they need for their research
    - Allows data from omics experiments to be traced back to the sample it came from for further research

RD Connect

# Biosample database

☐ RD-Connect biosample database contains sample-level data from all participating biobanks



RD Connect

# Patient (research participant) identifier

- Platform cannot store personally identifiable information (for obvious privacy reasons)

- How do we track that different data items (biosample, natural history data, exome sequence) all come from the same patient?

  - Assign an identifier (e.g. EURenOmics case: HEIDELBERG1234)

    - Advantage: Simple

    - Limitation: requires a central point (e.g. clinician) who knows the link between the patient and the identifier for all datasets

  - Generate an identifier from personally identifiable information

    - Advantage: the same patient will always have the same ID even if clinician A (who stored the biosample) doesn't know that clinician B uploaded an exome for the same patient

    - Limitation: requires consensus on a set of PII sufficient for generating a unique identifier – may be hard to do retrospectively if this info was not available

RD⬤Connect

# Existing systems for identifier

## US NIH GUID

☐ Originally used by National Database for Autism Research; concept now extended to several other NIH projects, with plans for a RD GUID

☐ Based on a standardised set of PII (including first, middle and surname as on birth certificate,  date of birth, city of birth as on birth certificate)

☐ Participant PII is entered into a Java webservice client application, which generates a one-way hash

☐ Hash is sent to central NIH server, which returns a GUID for that participant

RD Connect

# Plan moving forward for identifier

- At least in the interim, RD-Connect will establish an ID system for European RD projects contributing data to RD-Connect (partner projects can assign all patients an RD-ID)

- BUT use the SAME set of PII used in NIH (and Huntington) systems (interoperable)

- Continue to enable linkage of data in the platform by other mechanisms (e.g. manually generated ID) where it is not possible to generate an ID due to lack of PII (legacy data)

- Contribute to the task force jointly set up by IRDiRC and GA4GH and implement its output when ready

RD Connect

# Questions/feedback

Coordinator:

Hanns Lochmüller – hanns.lochmuller@ncl.ac.uk

Data helpdesk – personalised support!

John Dawson – john.dawson@ncl.ac.uk

Technical questions:

Sergi Beltran – sergi.beltran@cnag.crg.eu

Research questions:

Rachel Thompson – rachel.thompson@ncl.ac.uk

All other questions:

Emma Heslop – emma.heslop@ncl.ac.uk

@ConnectRD
@treat_nmd
@bushbykate

rachel.thompson@ncl.ac.uk
hanns.lochmuller@ncl.ac.uk

# #Brexit – thanks for all the support!