# Introduction to the Genomics England Research Environment

**Emily Perry / Roel Bevers**

Research Engagement Manager / Senior Bioinformatician – Research Services

15th March 2022

# Data security

- This training session will include data from the GEL Research Environment

- As part of your IG training you have agreed to not distribute these data in any way

- You are not allowed to:
  - Invite colleagues to watch this training with you
  - Take any screenshots or videos of the training

- We will record this training and distribute the censored video afterwards

Genomics
England

# Questions

Your microphones are all muted

Use the Zoom Q&A to ask questions

Upvote your favourite questions: if we are short on time we will prioritise those with the most votes
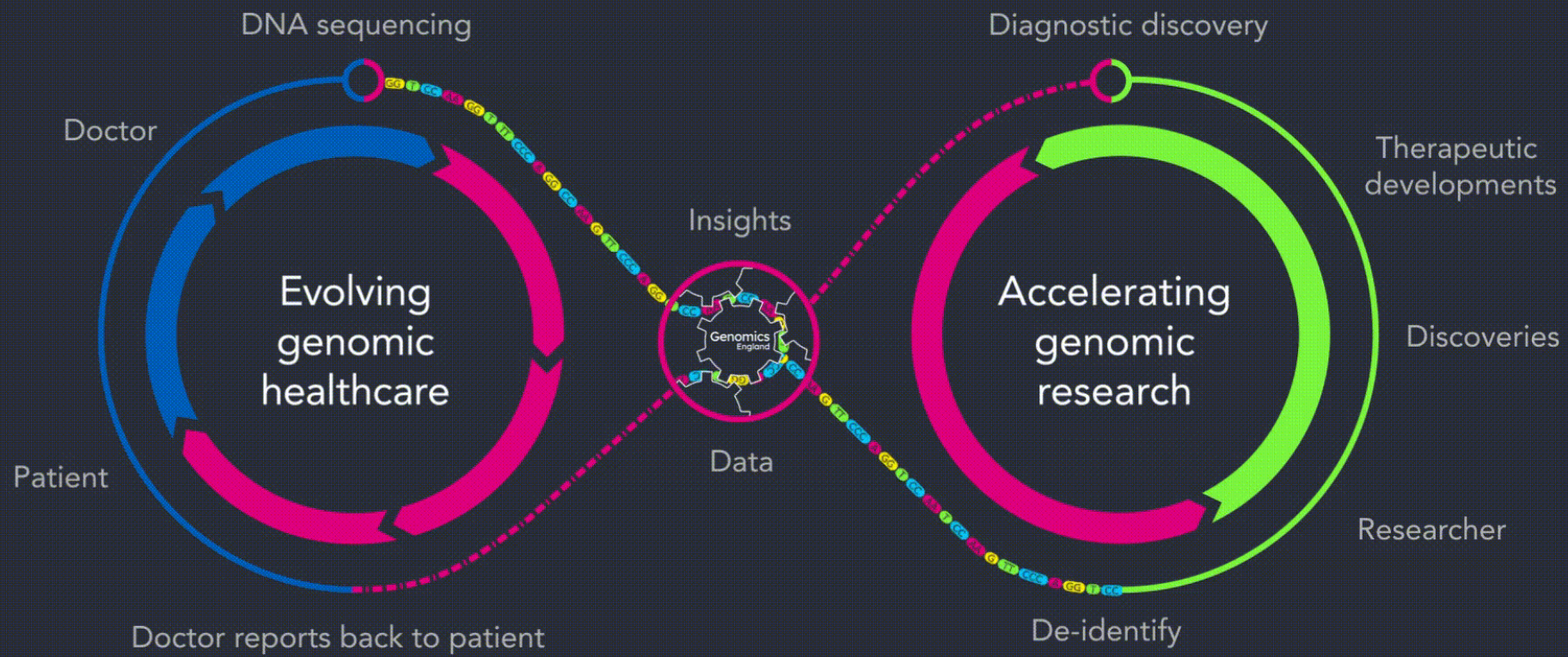
# Agenda

Genomics
England

# Data releases



Central text: V14 2022-01-27

Cycle: Data release → New genomes in → Redacted genomes out → Medical record updates

# 100,000 Genomes Project Data

| Genomics | Cancer | Rare Disease | Total |
|---|---|---|---|
| **Participants** | **17,955** | **72,955** | **90,259**<br>+ 35K COVID in CloudRE only |
| **Genomes** | **42,922**<br>Germline + Tumour<br>30x          100x | **75,526**<br>Germline<br><20% Singleton | **118,488** |

# 100,000 Genomes Project Data

**Genomics**

**Clinical Data**

- HPO terms
  - Rare disease
  - Other conditions
- Tumour staging
- Tumour location
- Histological subtype
- Treatment regimen

- NHS records
  - Hospital Episode Statistics
  - Mental Health Services Data Set

- Mortality data ONS

- Exit questionnaire for rare disease

- COVID-19 status

- Primary Care Data for COVID-19

Genomics England

# 100,000 Genomes Project Data

Genomics

Clinical Data

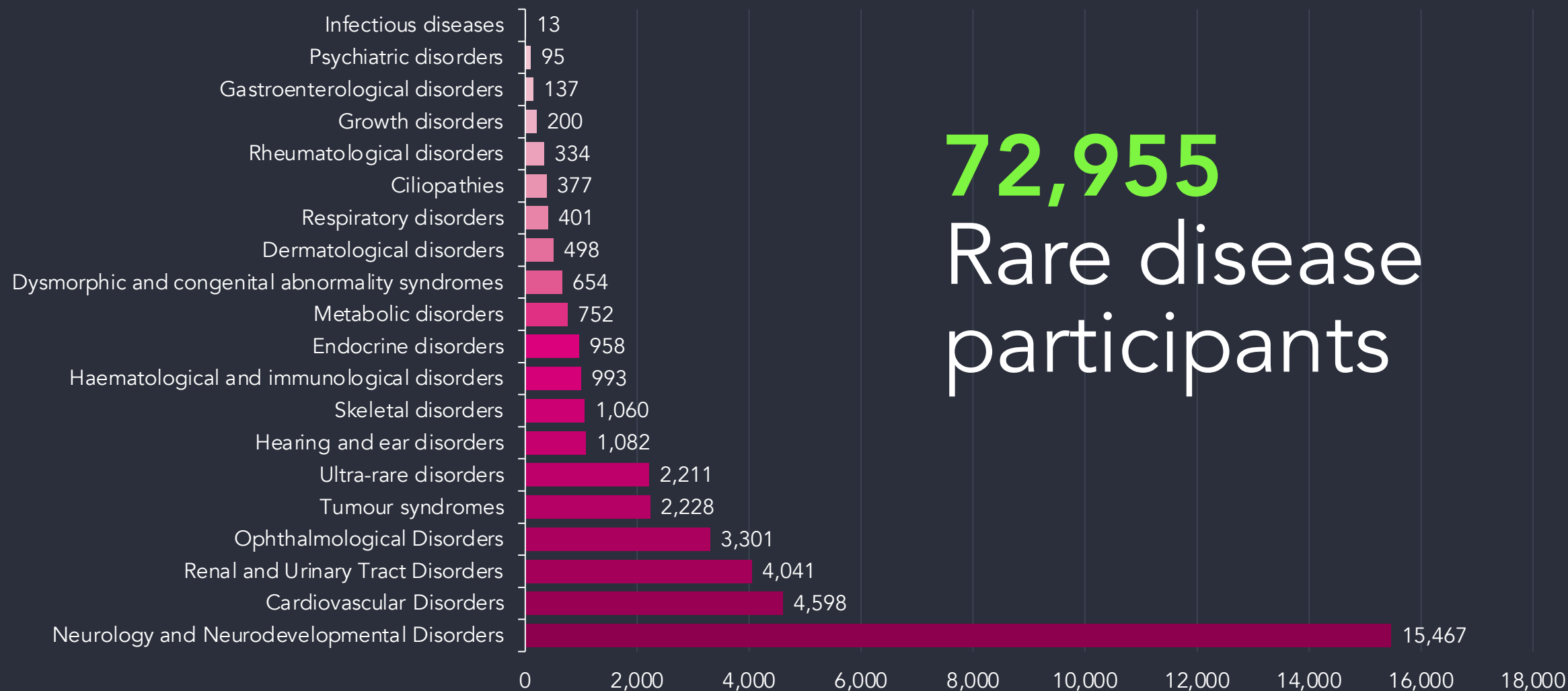Consent

## Clinically accredited pipelines

for diagnostics
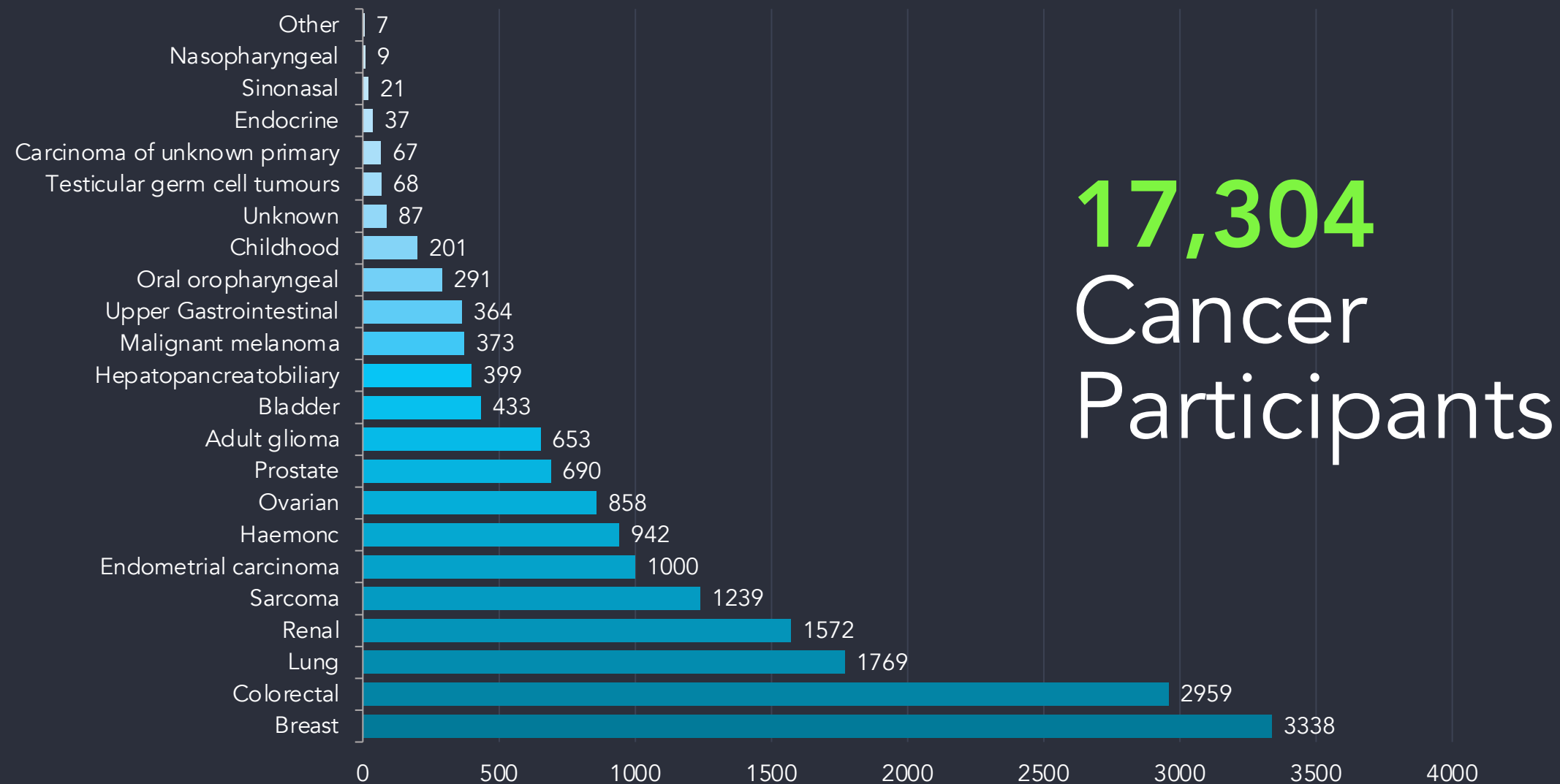
## Lifetime follow-up

+ full retrospective data

## Re-engagement

re-phenotyping

re-sampling

re-cruiting

Genomics England
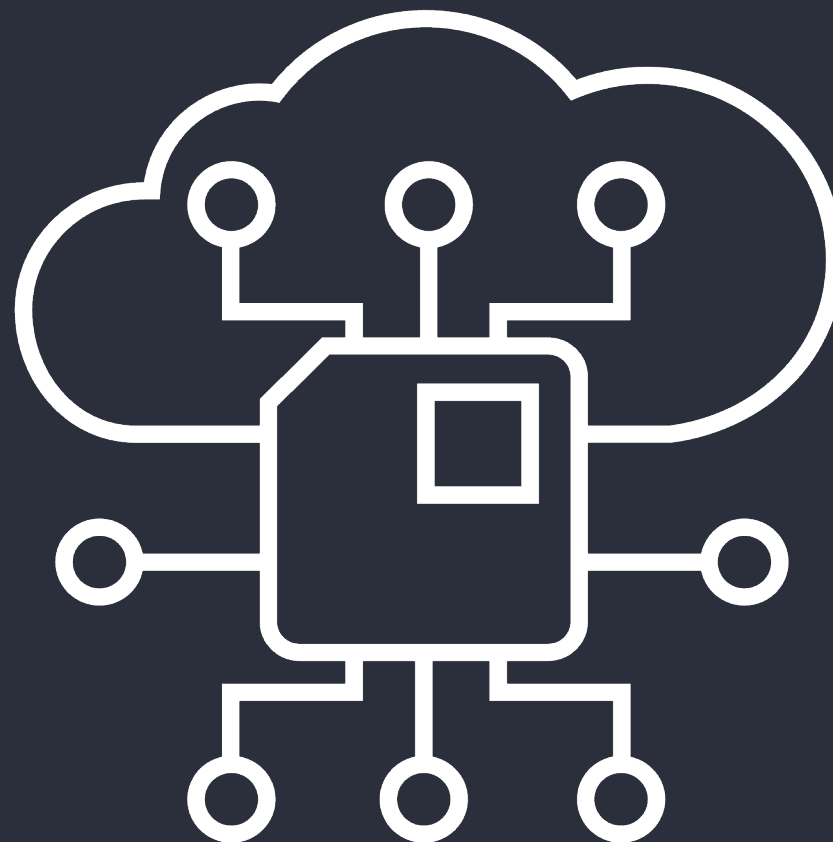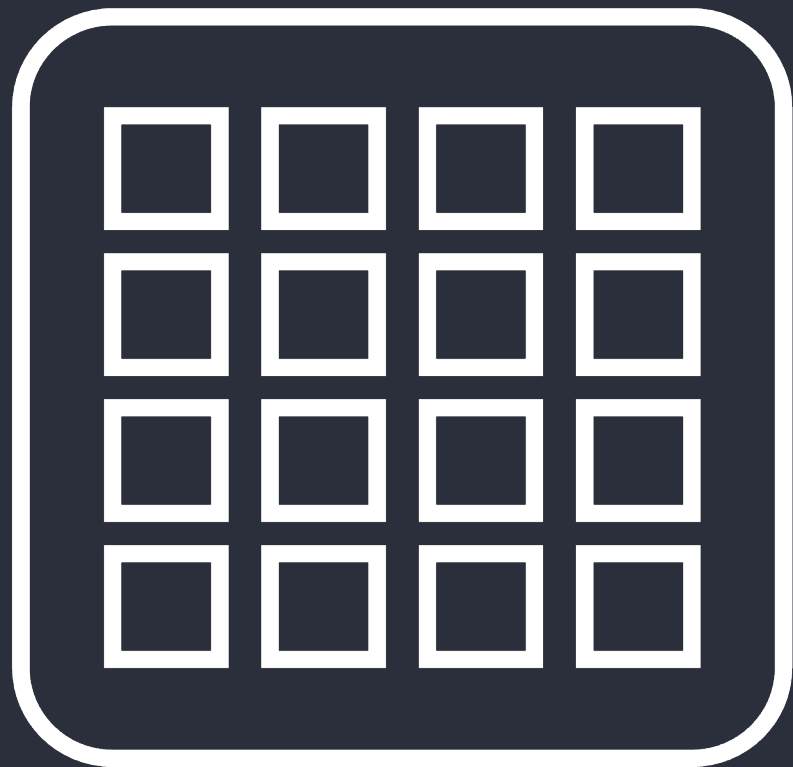
# 100KGP rare disease participants

Infectious diseases — 13
Psychiatric disorders — 95
Gastroenterological disorders — 137
Growth disorders — 200
Rheumatological disorders — 334
Ciliopathies — 377
Respiratory disorders — 401
Dermatological disorders — 498
Dysmorphic and congenital abnormality syndromes — 654
Metabolic disorders — 752
Endocrine disorders — 958
Haematological and immunological disorders — 993
Skeletal disorders — 1,060
Hearing and ear disorders — 1,082
Ultra-rare disorders — 2,211
Tumour syndromes — 2,228
Ophthalmological Disorders — 3,301
Renal and Urinary Tract Disorders — 4,041
Cardiovascular Disorders — 4,598
Neurology and Neurodevelopmental Disorders — 15,467

(x-axis: 0 — 2,000 — 4,000 — 6,000 — 8,000 — 10,000 — 12,000 — 14,000 — 16,000 — 18,000)

**72,955**
Rare disease participants

https://research-help.genomicsengland.co.uk/pages/viewpage.action?pageId=66846734

Genomics
England

10

# 100KGP Cancer participants

| Category | Participants |
|---|---|
| Other | 7 |
| Nasopharyngeal | 9 |
| Sinonasal | 21 |
| Endocrine | 37 |
| Carcinoma of unknown primary | 67 |
| Testicular germ cell tumours | 68 |
| Unknown | 87 |
| Childhood | 201 |
| Oral oropharyngeal | 291 |
| Upper Gastrointestinal | 364 |
| Malignant melanoma | 373 |
| Hepatopancreatobiliary | 399 |
| Bladder | 433 |
| Adult glioma | 653 |
| Prostate | 690 |
| Ovarian | 858 |
| Haemonc | 942 |
| Endometrial carcinoma | 1000 |
| Sarcoma | 1239 |
| Renal | 1572 |
| Lung | 1769 |
| Colorectal | 2959 |
| Breast | 3338 |

**17,304 Cancer Participants**

# Two REs: RE1.0 and CloudRE

# 3. Tools in the Research Environment

**Participant Explorer**
Search for participants by phenotypes or identifiers

IVA
Explore genomic variants and what genotypes GEL participants have for them

IGV
Visualise genomic data

LabKey
Explore the tables of GEL data

# Demo: tools in the RE

Genomics
England

# 4. Programmatic access to Genomics England data

Labkey API allows you to:

- Combine data and filters from multiple tables

- Replicate queries exactly:
  - When new data releases come out
  - Between analyses

- Work in a variety of programming languages, but most support for Python and R

- Work both locally and on the HPC

Genomics
England

# Set up .netrc

- You can access the same data via the LabKey API as you can through other means

- You will need to configure access to the LabKey API with your username and password
  - In your home directory
  - On the HPC

- You do this by editing a file called .netrc

# Labkey API - Python

```python
# Import the needed modules, labkey and pandas (for dataframes)
import labkey
import pandas as pd

# Specify what we are connecting to, and what schema and tables we want
labkey_server = "labkey-embassy.gel.zone"  # The labkey server we are connecting to.
project_name = "main-programme/main-programme_v14_2022-01-27"  # The data we want to access.
context_path = "labkey"
schema_name = "lists"  # The schema we are getting data from.

# Create the SQL query as a string
sql = (
    "SELECT participant.Participant_Id, participant.Programme, sequencing_report.lab_sample_id /
    FROM lists.participant /
    JOIN lists.sequencing_report /
    ON participant.Participant_Id = sequencing_report.Participant_Id /
    WHERE sequencing_report.lab_sample_id between 1018056774 and 1018068634;"
)

# Create an object that will let us connect to the LabKey databases. This does not change.
server_context = labkey.utils.create_server_context(
    labkey_server, project_name, context_path, use_ssl=True
)

# The data are returned and stored in the variable results.
results = labkey.query.execute_sql(server_context, schema_name, sql)

# Data are returned as a dictionary, will all of the table information stored under the key "rows".
# We make a dataframe of all of the table information using pandas.
table_of_data = pd.DataFrame(results["rows"])
```

Change the database version you're accessing

SQL query with standard SQL syntax

Data returned as a dictionary and can be converted to a data-frame

Genomics
England

# Labkey API - R

```r
# Import the labkey library
library(Rlabkey)

# Set the baseURL
labkey.setDefaults(baseUrl= "https://labkey-embassy.gel.zone/labkey/")
Project_name <- "/main-programme/main-programme_v14_2022-01-27"

# Write your SQL query here
query <- "SELECT participant.Participant_Id, participant.Programme, sequencing_report.lab_sample_id
    FROM lists.participant
    JOIN lists.sequencing_report
    ON participant.Participant_Id = sequencing_report.Participant_Id
    WHERE sequencing_report.lab_sample_id between 1018056774 and 1018068634;"

mysql <- labkey.executeSql(
    schemaName="lists",
    colNameOpt = "rname",
    maxRows = 100000000,
    folderPath = project_name,
    sql = query
)
```

Change the database version you're accessing

SQL query with standard SQL syntax

Data returned as a data-frame

Genomics
England

38

# Demo: running LabKey API scripts

# 5. Running command line tools and pipelines using our HPC cluster

# Pre-installed tools on the cluster

| | | | | |
|---|---|---|---|---|
| APBS | CADD | GISTIC | MultiQC | SHAPEIT |
| AdapterRemoval | CNVator | GMAP-GSNAP | NGS | SPAdes |
| AutoDock-Vina | CaVEMAN | HISAT2 | NextGenMap | STAR |
| BCFTools | Canvas | HLA-LA | OMA | SURVIVOR |
| BEDOPS | Centrifuge | HTSlib | OptiType | SVINT |
| BEDTools | Circos | IGV | OrfM | Salmon |
| BLAST | Clustal-Omega | IMPUTE2 | PHYLIP | Sambama |
| BLAT | EIGENSOFT | Jellyfish | PLINK | SeqAn |
| BWA | FASTX-Toolkit | KNIME | Pindel | Trimmomatic |
| BamTools | FASTQC | Kraken | Pysam | UN-CNVc |
| Bio-DB-HTS | FlashPCA2 | LUMPY | Quip | VCFtools |
| BioPerl | GATK | MAFFT | RTG-Tools | VEGAS |
| Bowtie | GD | MetaGeneAnnotator | SAMtools | VEP |

# Pre-installed tools on the cluster

| | | | | |
|---|---|---|---|---|
| Velvet | meRanTK | Tabix | ROOT | Doxygen |
| ViennaRNA | minimap2 | verifyBamID | XML-LibXML | Gradle |
| alleleCount | Mosdept | Vt | XML-Parser | Junit |
| Bam-readcount | Ncbi-vdb | GATE | datamash | LZOM4 |
| Cellbase | New_fugue | GCC | ntCDF | |
| Cromwell | Nextflow | GCCcore | Savvy | |
| Cryptsplice | Picard | LLVM | Shrinkwrap | |
| Ea-utils | Platypus | Ispc | GDB | |
| Fastp | Rvtests | DBD-mysql | Autoconf | |
| Gvcfgenotyper | Seqtk | GDAL | Automake | |
| Kallisto | Singularity | HDF5 | Autotools | |
| liftOver | Snptest | MariaDB | Boost | |
| Locuszoom | Strelka | PyTables | Cmake | |

Genomics
England

# Demo: using standard tools on the HPC

# Creating your own workflows to use on the cluster

- You can incorporate any of the existing tools into your own workflow

- Import scripts via Airlock

- Using containers
  - Singularity
  - Docker
  - Quay.io

# Ready-made scripts/workflows

- Extract variants (small or SV; germline or somatic) by coordinate or gene

- Gene centric SNV reports for cancer

- GWAS

- Survival - cancer

- Aggregate variant testing

- Functional variant annotation

# Demo: running a workflow on the HPC

# SV-CNV workflow

- Submit a list of genes or regions
- Find all SVs/CNVs in these genes/regions
- Choose somatic/germline, cancer/rare disease

# 6. The Airlock, restricted import and export of data



Patient

Evolving genomic healthcare

Researcher

Accelerating genomic research

Unrestricted data sharing

Patients

Healthcare teams

Researchers

Genomics
England

# Airlock: what can you export

- Data matching your approved research project

- Aggregate data for groups ≥5 participants

- Data from a project that has not been approved

- Individual data or data that can be otherwise identified

Genomics
England

# Demo: exporting data using the Airlock

# 7. The future: CloudRE

- GEL data
  - 100,000 Genomes Project: rare disease and cancer
  - COVID-19 – severe and mild cases
- Point-and-click tools
  - Cohort browser
  - Running pipelines
- Compute in the Cloud
  - Flexible options based on budget and speed needed
  - Not limited by load
- Bring in data and tools from outside
  - From Github
  - In Containers
  - In S3 buckets

Genomics
England

Demo: tools in the CloudRE

# If the CloudRE is for you…

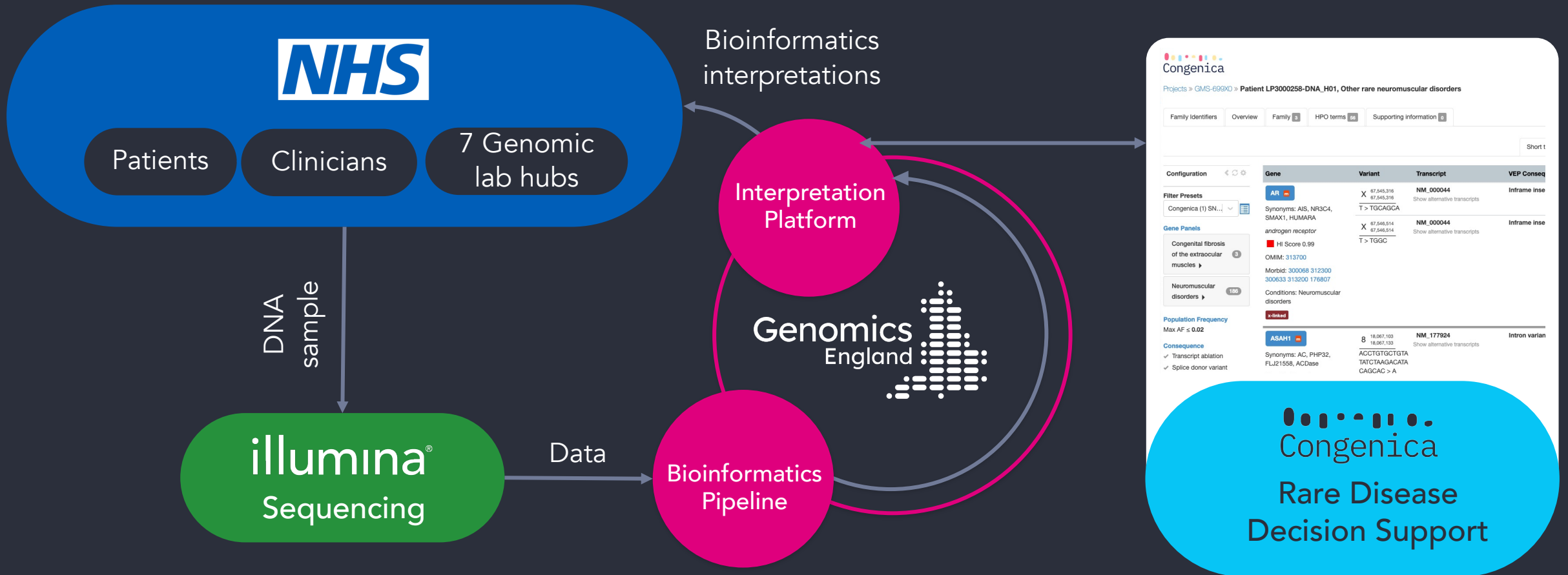Contact gecip-help@genomicsengland.co.uk

Contact your partnerships director (James/Kate)

- Brief description of why you'd like access
  - Use-case
  - Data
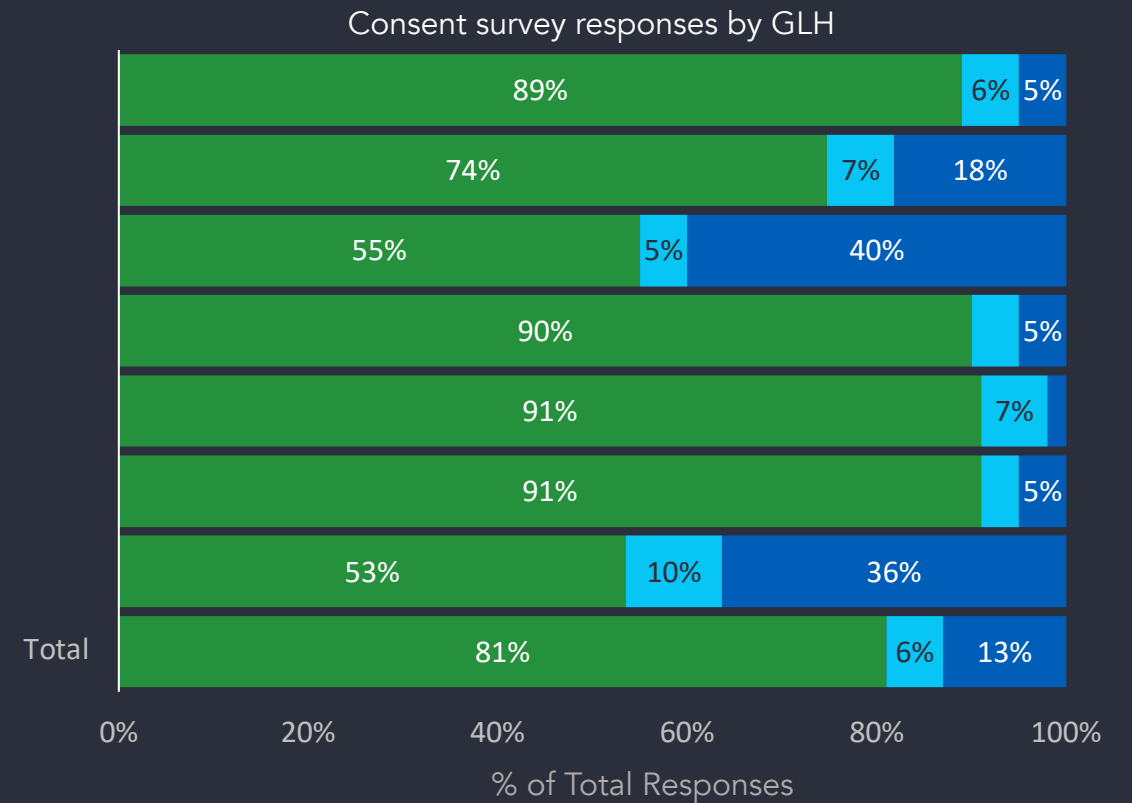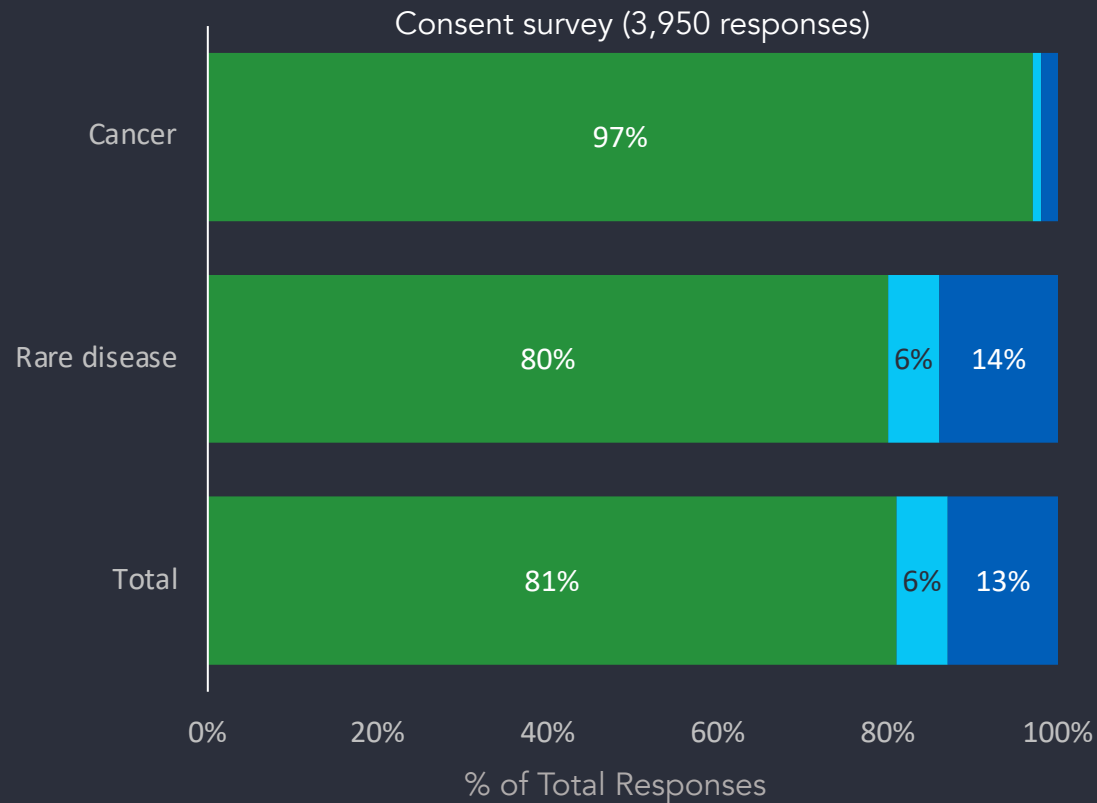- We will get in touch to discuss feasibility

# The future: Genomic Medicine Service

# The future: Genomic Medicine Service

What consent have patients given for research?

Answers    Yes    No    No response

Consent survey (3,950 responses)

Consent survey responses by GLH

Consented that data and samples may be used for research, separate to NHS care

# 8. Getting help

Check our documentation:
- https://research-help.genomicsengland.co.uk/
- Click on the documentation icon in the environment

Contact our Service Desk:

- ge-servicedesk@genomicsengland.co.uk

**Genomics**
England

# 9. Questions



Your microphones are all muted

Use the Zoom Q&A to ask questions

Upvote your favourite questions: if we are short on time we will prioritise those with the most votes

# Future sessions

24th May — Building a cohort based on phenotypes and a matching control cohort

22nd July — Finding participants based on genotypes

20th September — Getting medical history for participants

22nd November — Using the HPC to run jobs

# Thank you