



Introduction to the Genomics England Research Environment

Emily Perry

Research Engagement Manager

21st January 2025



Questions



All your
microphones
are 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

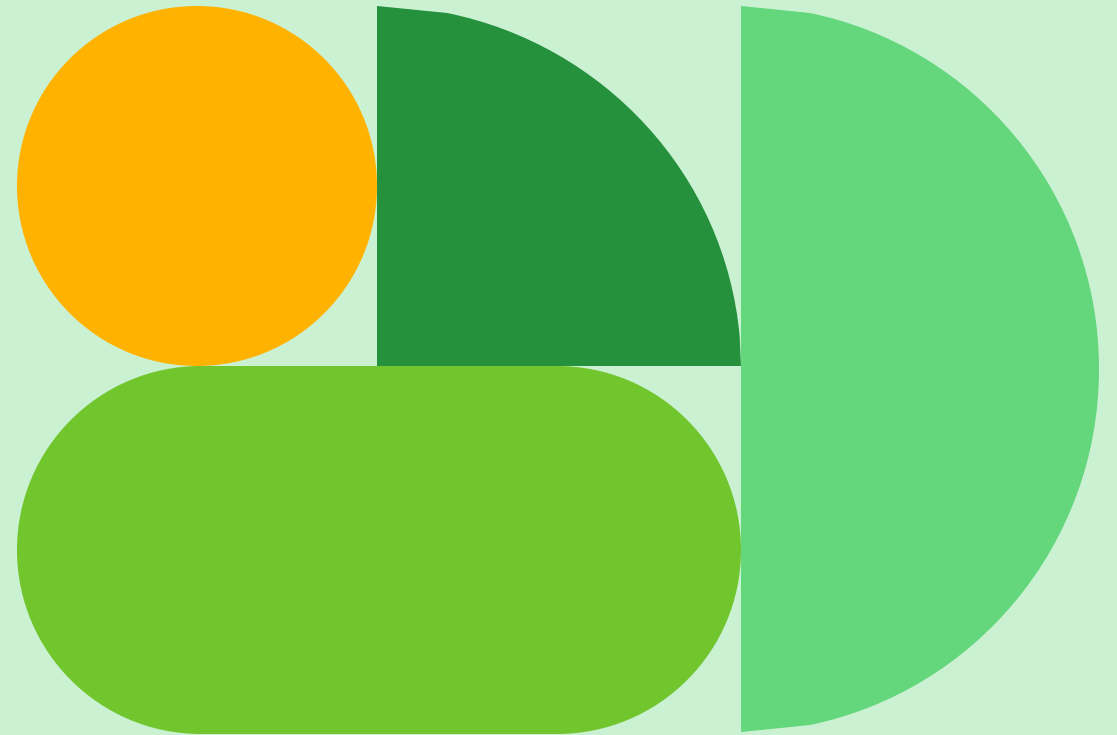
Helpers



**Matthieu
Vizquete-Forster**
Bioinformatician -
Research Services

Agenda

- 1 Introduction and admin
- 2 What is the RE?
- 3 Sources and types of data
- 4 Tools in the Research Environment
- 5 Programmatic access to NGRL data
- 6 Running command line tools and pipelines on the HPC
- 7 Import and export of data and tools
- 8 Help and questions



2. What is the RE?

A virtual machine



Genomic data

Phenotype data

Point and click tools

HPC and command line
tools

A Trusted Research Environment

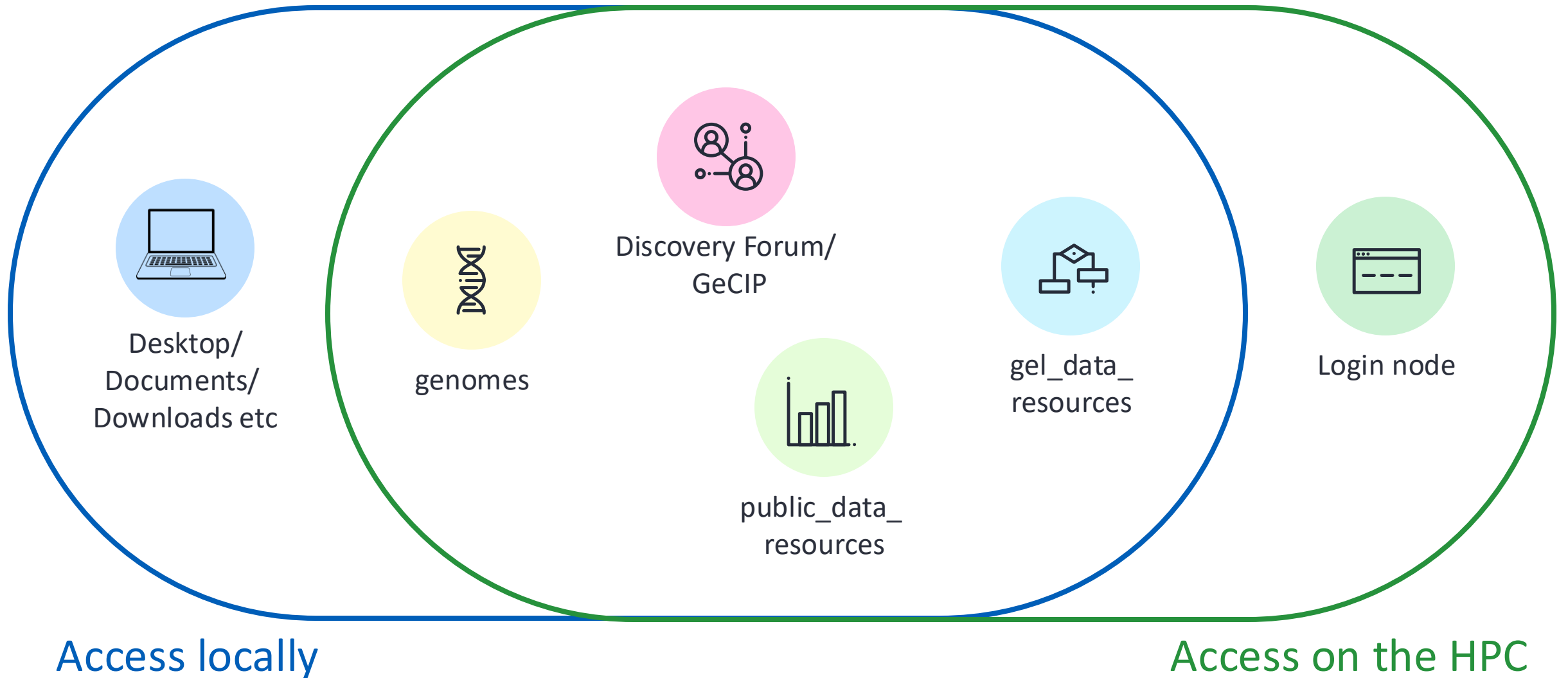
You can...

- Work with the data in the RE
- Copy/paste in
- Bring in Containers
- Access whitelisted websites
- Request to export the results of your analysis

You cannot...

- Share folders between your computer and the VM
- Copy/paste out
- Export files
- Access most of the internet

Files in the RE



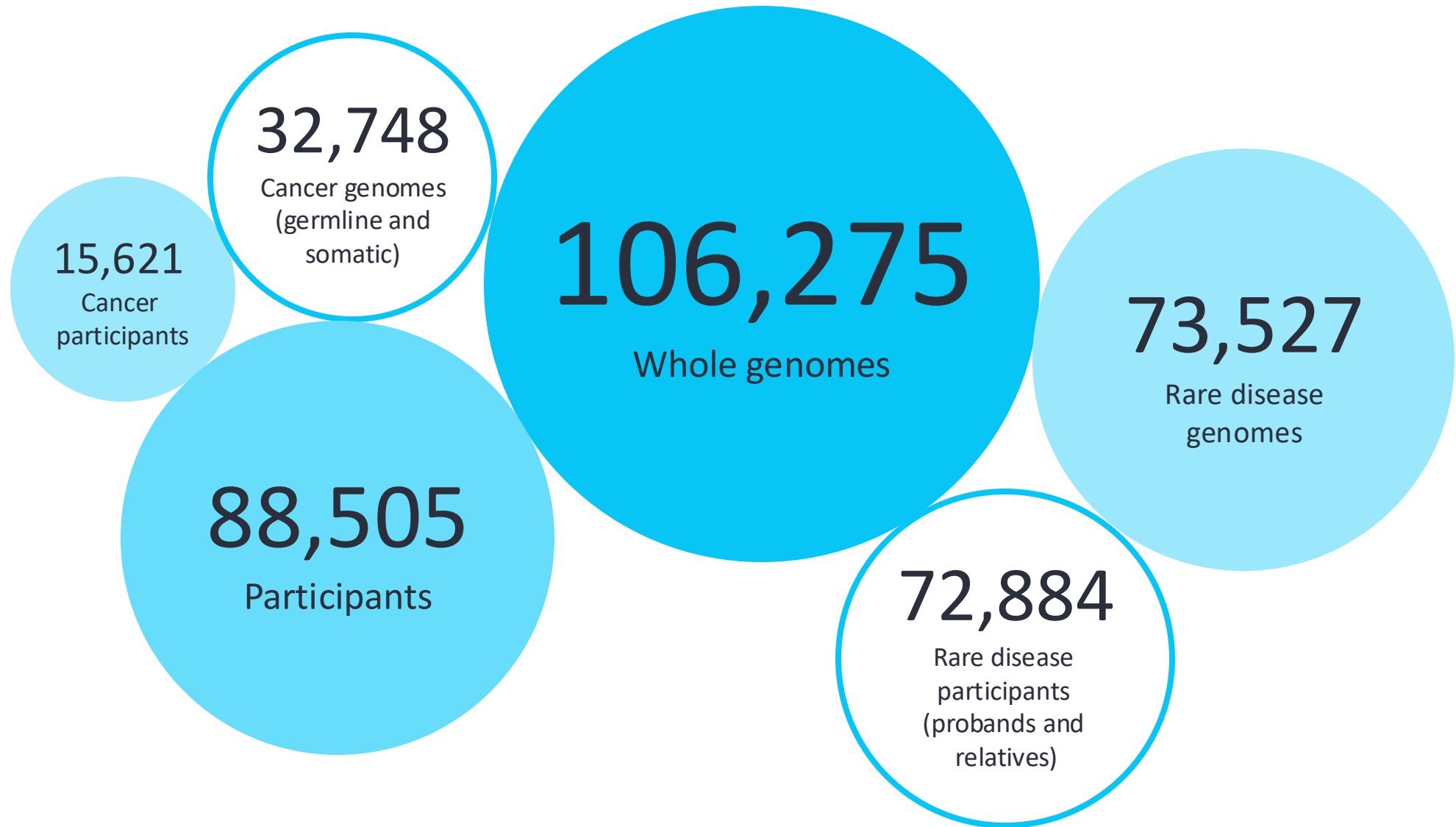
Access locally

Access on the HPC

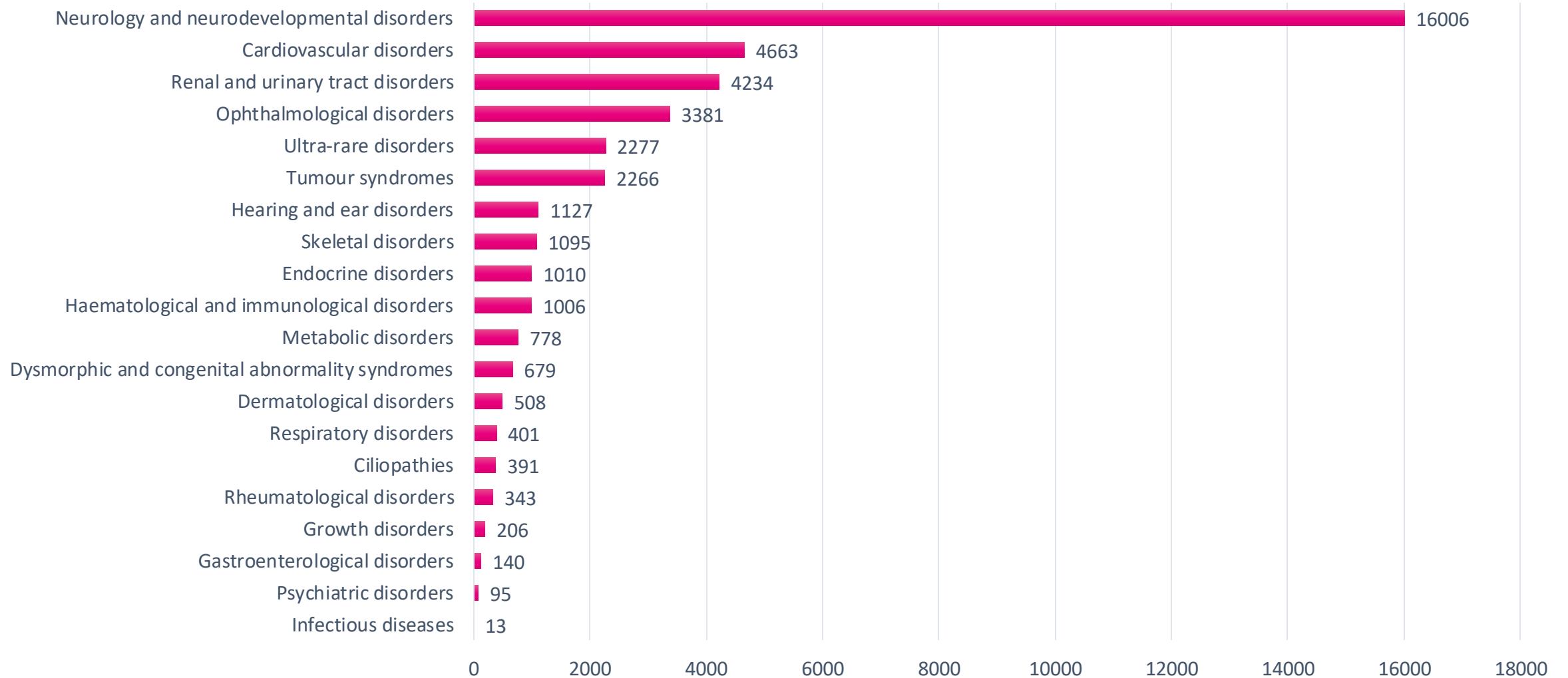
3. Sources and types of data



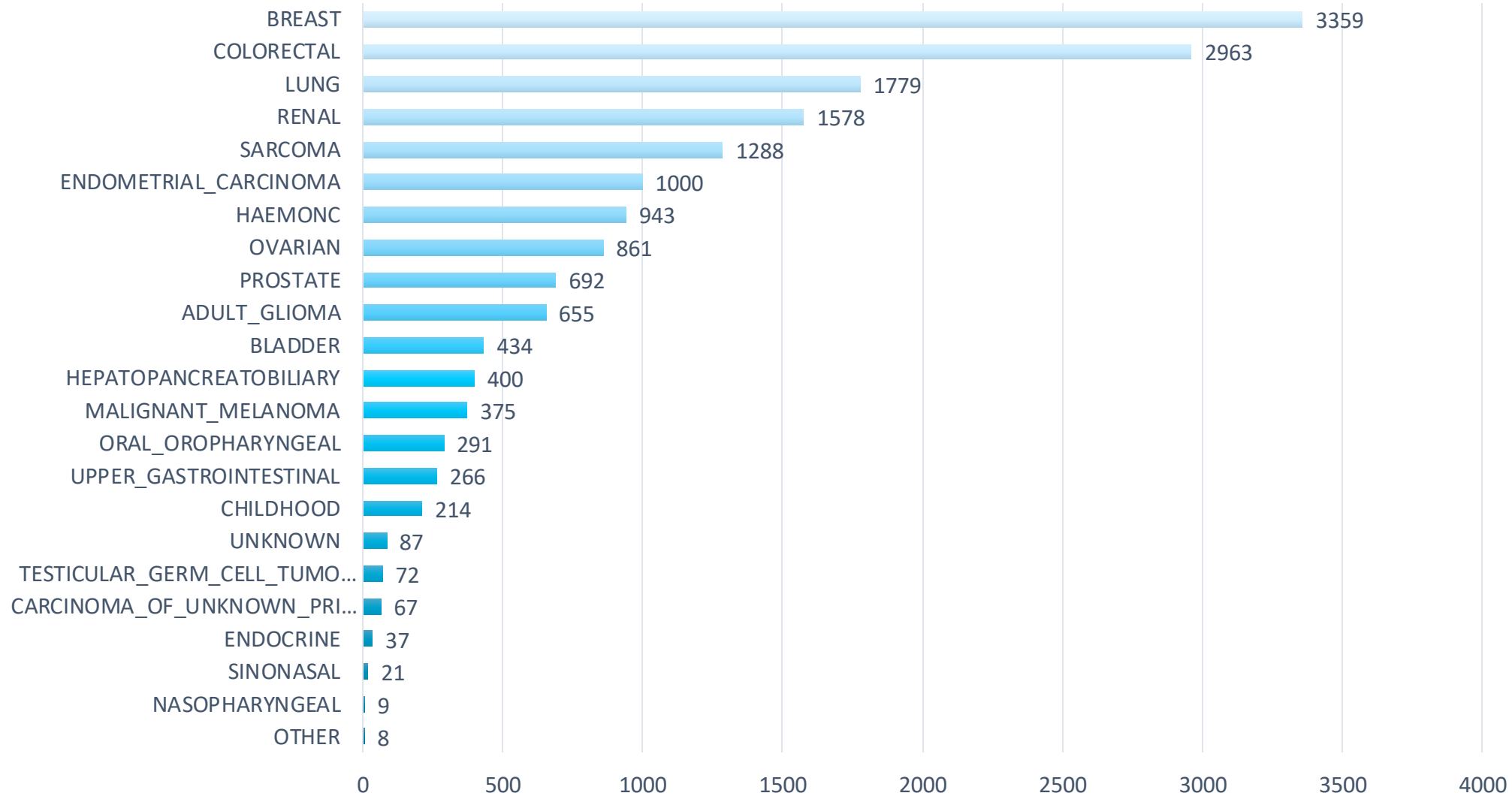
100,000 Genomes Project



100,000 Genomes rare disease



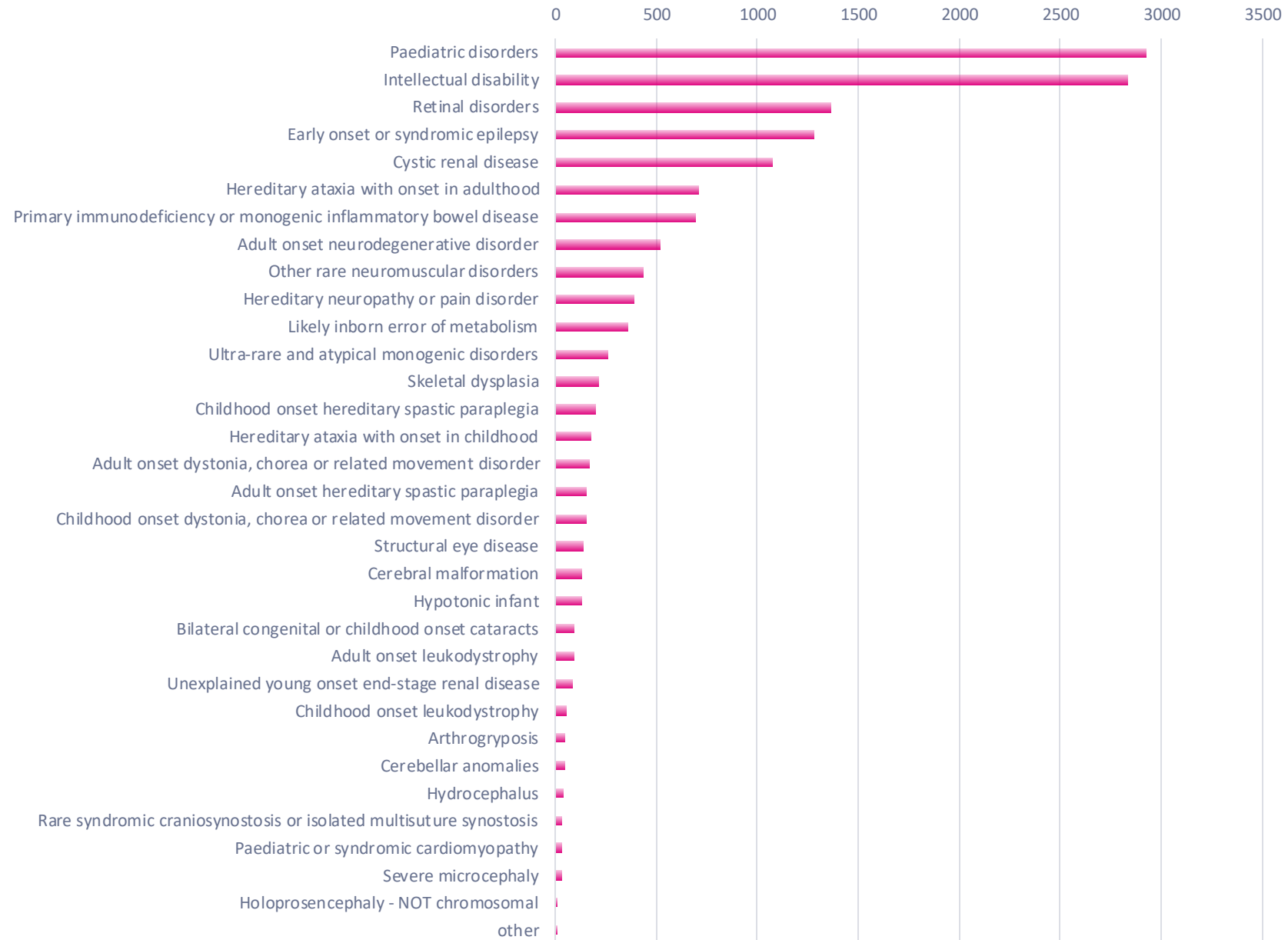
100,000 Genomes cancer



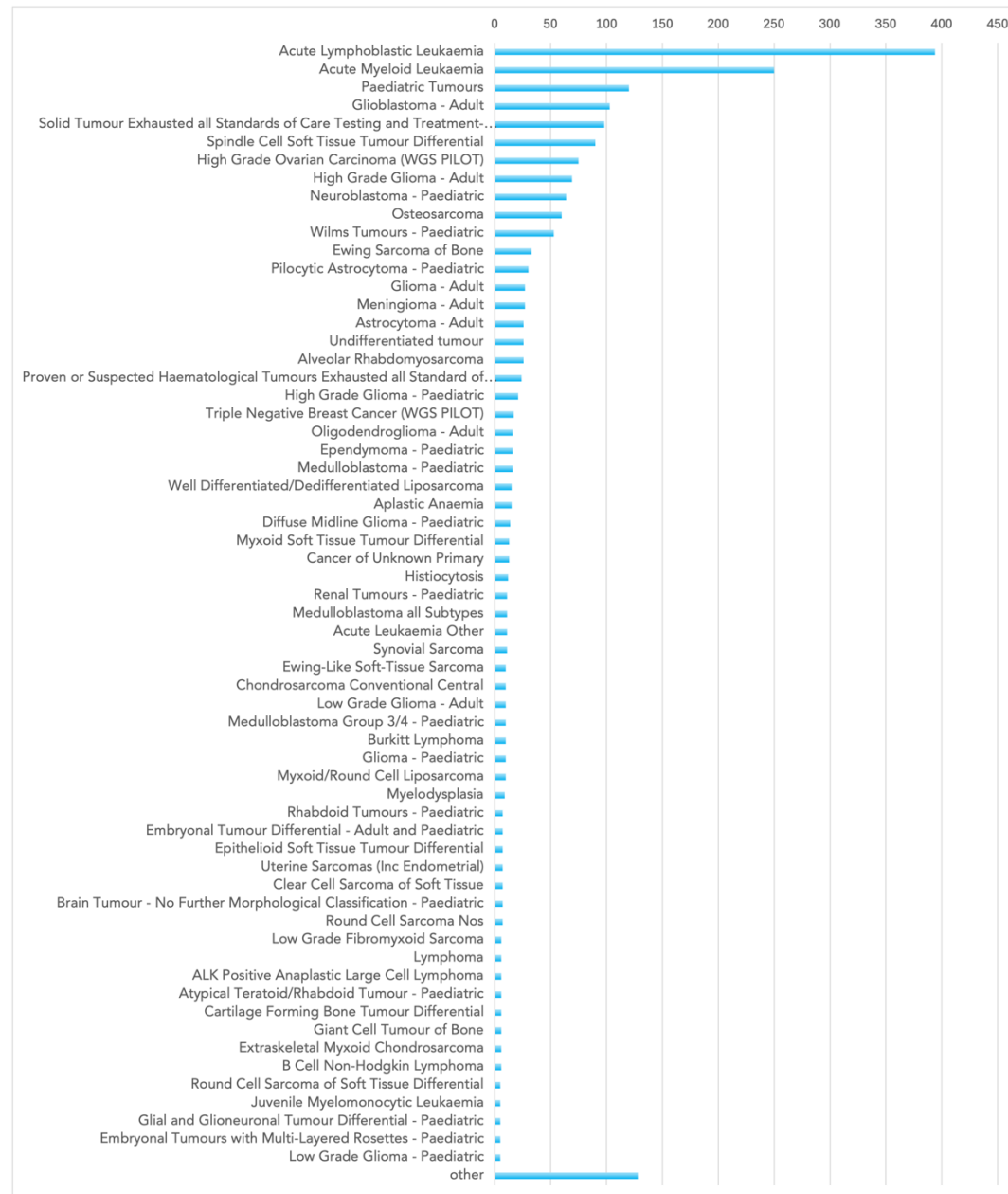
NHS GMS



NHS GMS rare disease



NHS GMS cancer



For **EVERY** genome

Alignment

as BAM or CRAM files

Variant calls

as VCFs, including
gVCF, repeats VCFs,
structural variant VCFs

Analysis

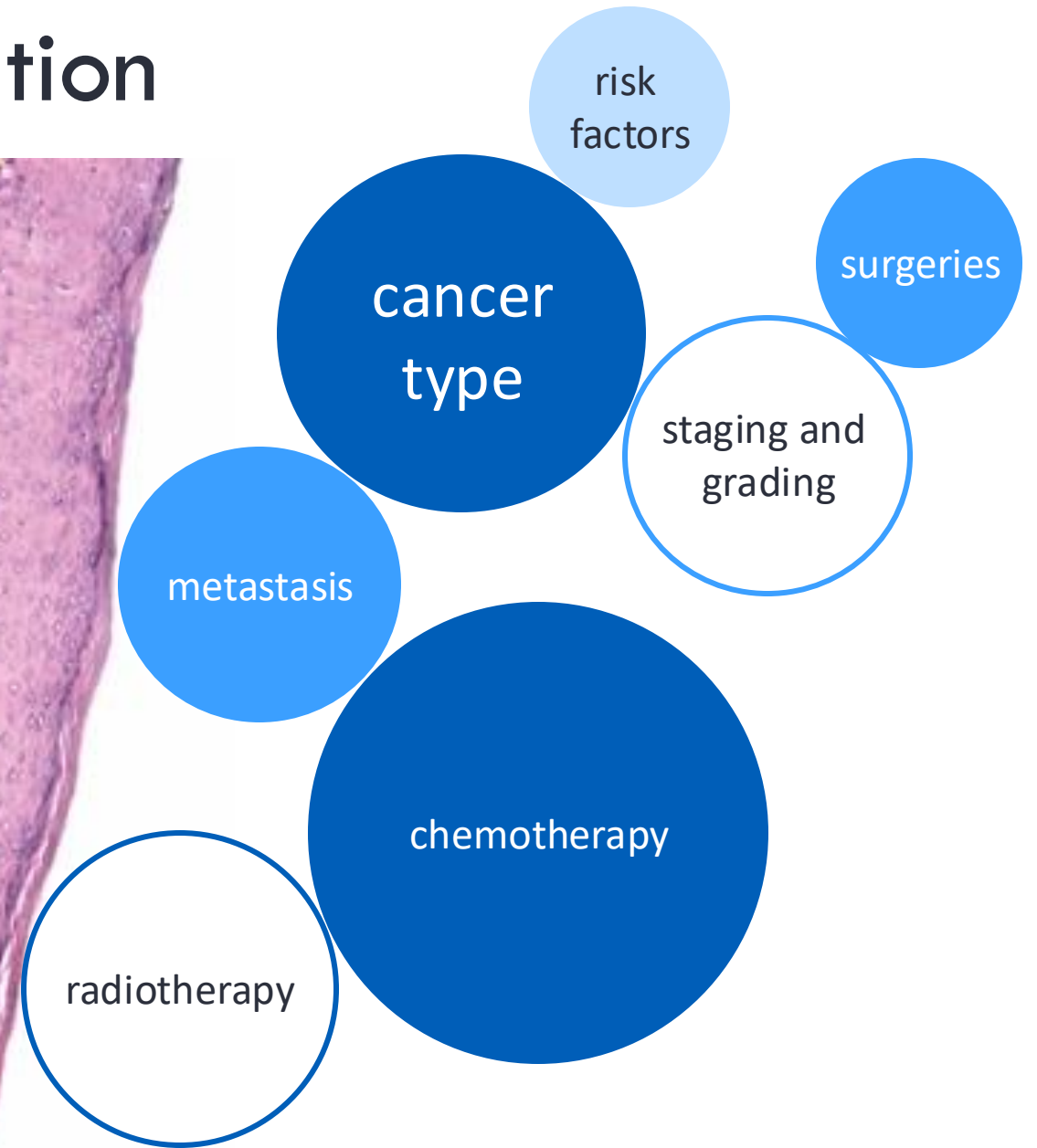
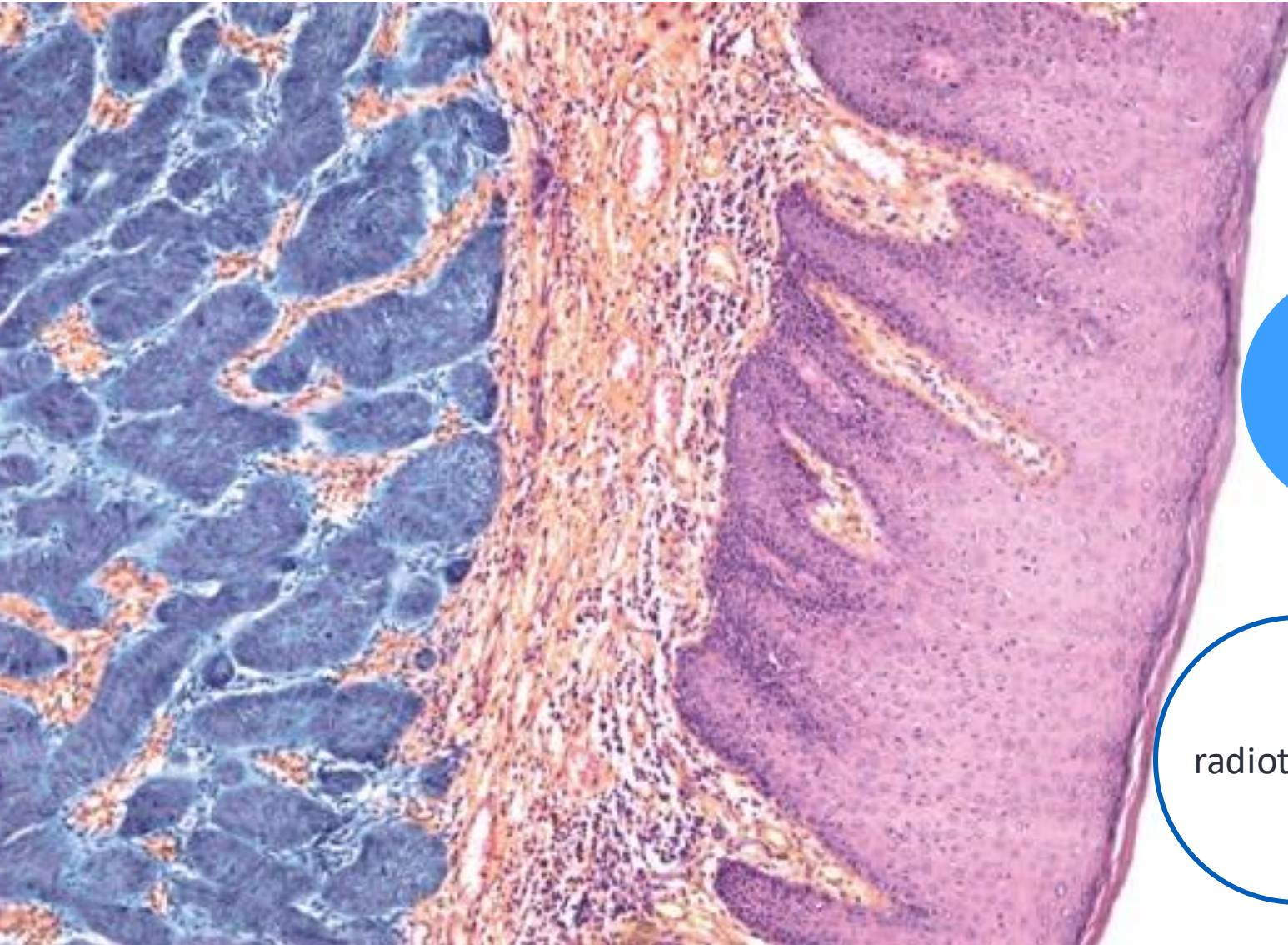
Variant tiering and
tumour mutational
signatures

Rare disease phenotyping

- Disease classification
- HPO terms present/absent
- Measurements and observations (not universal)
 - general measurements
 - early childhood observations
 - details of imaging (but not results)
 - genetic tests
 - lab tests

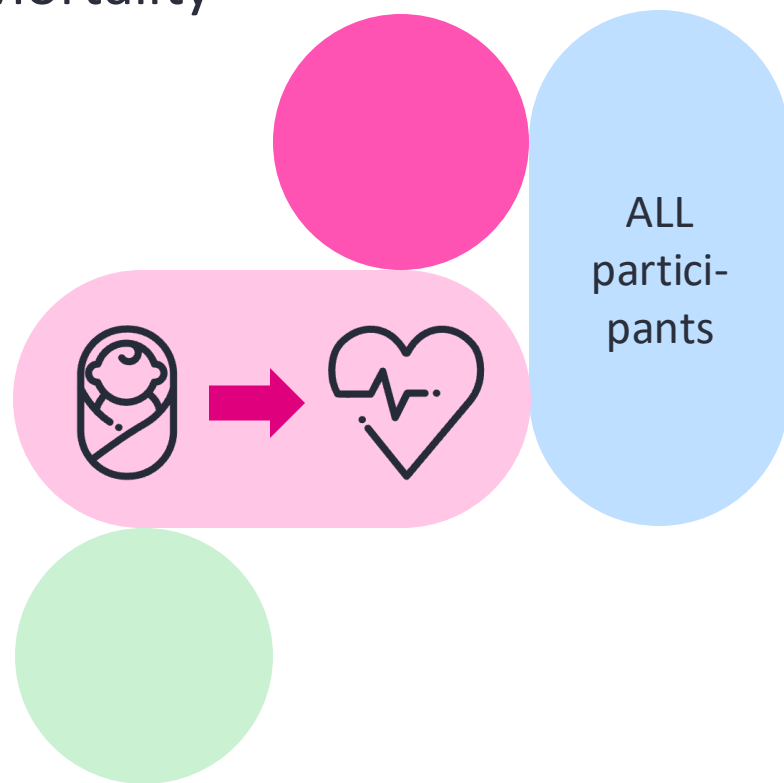


Cancer tumour characterisation



Medical history

- NHSE hospital episode statistics
- Mental health data
- Mortality



Hospital episode statistics

Out-patients

op

Planned day
appointments in
hospital

Admitted
patient care

apc

Overnight hospital
stays

Critical care

cc

Time on life
support

Accident
and
emergency

ae

Unplanned
emergency visits –
walk-in or
ambulance

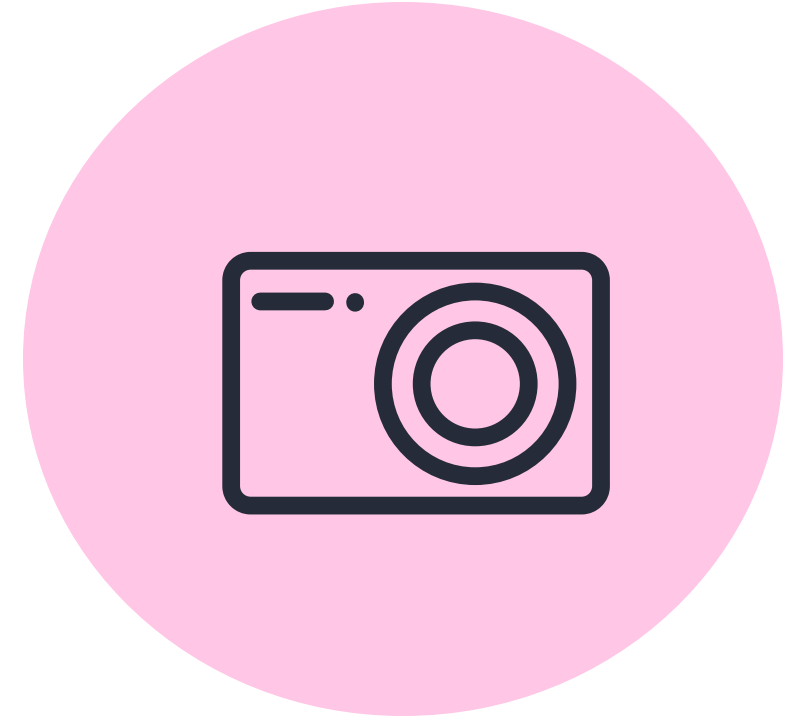
We don't have...



Free text

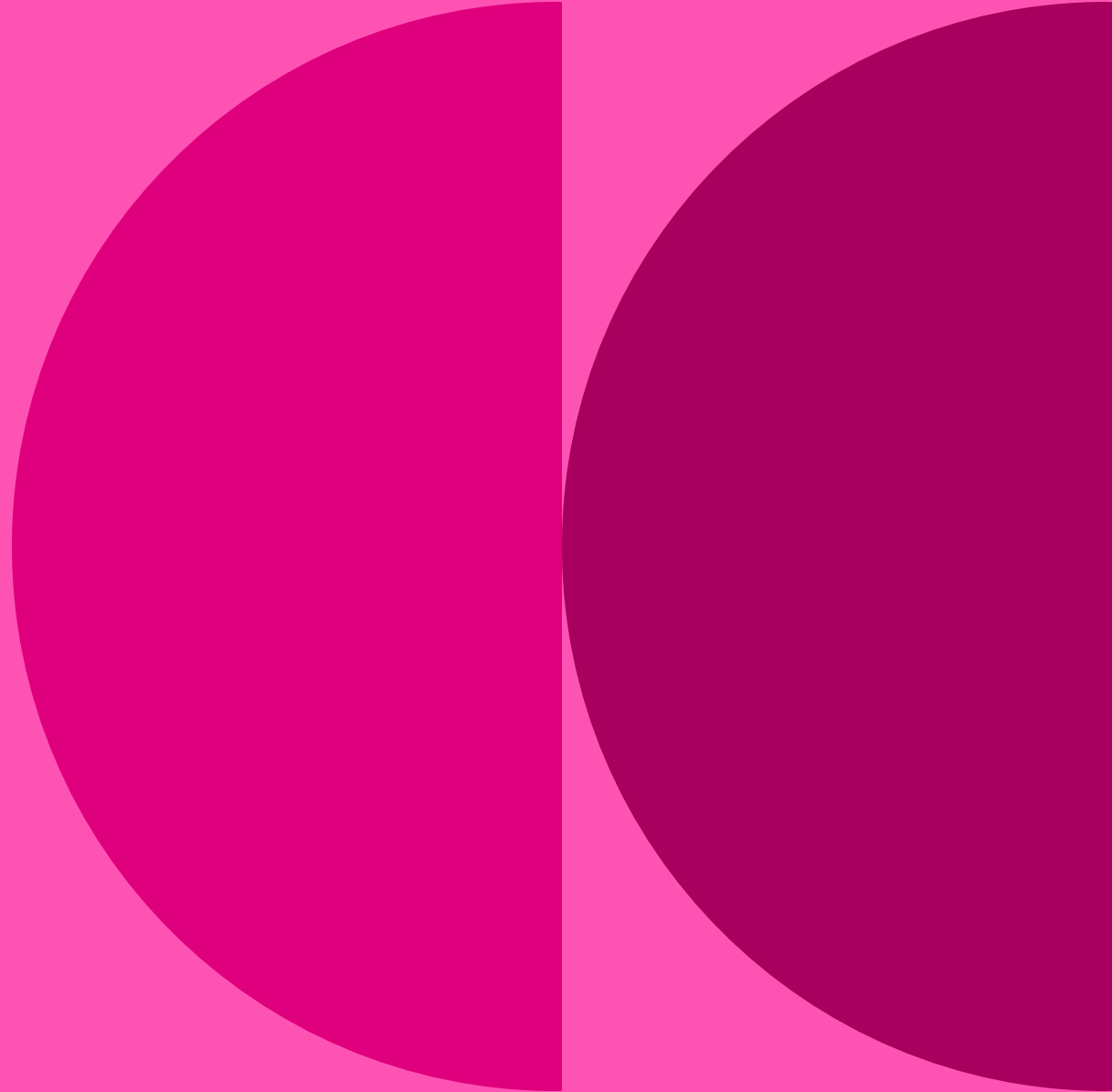


GP data



diagnostic
imaging

4. Tools in the Research Environment



LabKey

- Central database of:
 - Clinical data
 - Results of bioinformatics analysis
 - Locations of genomic files
- Point and click interface
- API



Data dictionary

Table	Field	Short Name	Description	Value
av_imd	participant_id	Participant ID	Participant Identifier (supplied by Genomics England)	participantid, xs:string
av_imd	anon_tumour_id	Pseudonymised tumour ID (IMD)	NCRAS specific ID for the tumour (does not link to GeL tumour_id) Pseudonymised tumour ID. This field replaces tumour_pseudo_id. Note: anon_tumour_id contains a different set of pseudonymised tumour ids to tumour_pseudo_id	xs:string
av_imd	imd	Index of Multiple Deprivation	Measure of deprivation at small area level derived from the IMD domain. Quintiles are weighted equally by the number of LSOAs.	1most deprived 22nd quintile 33rd quintile 44th quintile 5least deprived
	participant_id	Participant ID	Participant Identifier (supplied by Genomics England)	xs:string
	aliasflag	Alias Check Flag	0,1 (Indicates that this patient record has been deduplicated with another patient and the tumour(s) moved to that other patientid)	0,1 (Indicates that the record has been deduplicated with another patient record and the tumour(s) moved to that other patientid)
	birthdateflag	Date Of Birth Check Flag	Date Of Birth Check Flag	0,1,2,3 (Set to 0 if the date of diagnosis is not specified, 1 if the month and year of diagnosis were specified but the day was not specified, 2 if the month and day are not specified but the month and day are not specified but the date was less specific than any of the above, 3 if the date was less specific than any of the above) 0Set to 0 if the date was fully specified 1month and year of diagnosis are known 2 year is fully known, but the month and day are not specified 3date less specific
patient	sex	Person Phenotypic Sex	PERSON_PHENOTYPIC_SEX_CLASSIFICATION, PERSON_GENDER_CODE, which is the most recent	1Male 2Female 9 Indeterminate (unable to be classified as either male or female)

Lists of tables and columns

Value type or meaning of codes

Description of the data

LabKey demo



Participant Explorer

- Search for participants by:
 - IDs
 - Clinical concepts
 - Personal details
- View/compare medical histories



Participant Explorer demo

The desktop environment features a dark theme with a grid of application icons. The icons include: Computer, Text Editor, Airlock, Research Environment Documentation, Welcome Pack, eperry's Home, Data Discovery, Participant Explorer, report.tsv, Trash, Firefox, Visual Studio Code, Ensembl, Document Viewer, IGV Browser, RE Messages, IVA 2.0, R, Terminal Emulator, Rocket Chat, LibreOffice, Panel App, Research Registry, GVim, Labkey, Git GUI, RStudio, Old Firefox Data, Open Targets, Emacs, LibreOffice 7.6, and CloudOS. On the right side, there is a large logo for "Genomics england" with a stylized map of the United Kingdom. Below the logo, there are two overlapping rectangular shapes: a light blue one on top and a light grey one on the bottom.

Interactive variant analysis (IVA)

- Find participants with genetic variants
- Filter variants in a participant by family genotypes
- Filter on genome features



IVA demo



Computer, Text Editor, Airlock, Research Environment Documentation, Welcome Pack

eperry's Home, Data Discovery, Participant Explorer, report.tsv

Trash, Firefox, Visual Studio Code

Ensembl, Document Viewer, IGV Browser, RE Messages

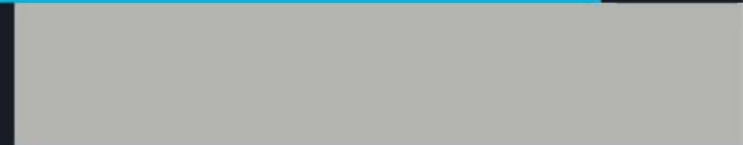
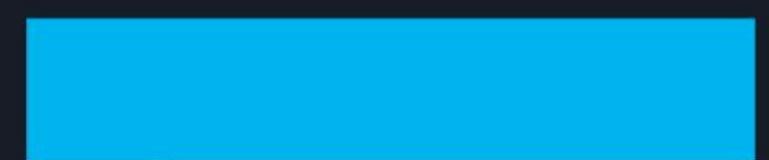
IVA 2.0, R, Terminal Emulator

Rocket Chat, LibreOffice, Panel App, Research Registry

GVim, Labkey, Git GUI, RStudio

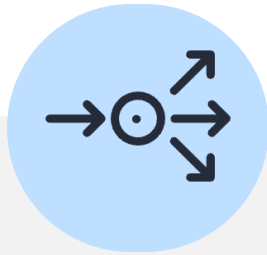
Old Firefox Data, Open Targets, Emacs

LibreOffice 7.6, CloudOS

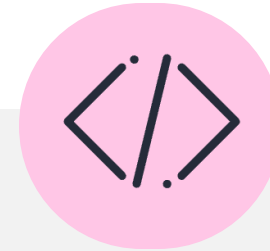


5. Programmatic access to NGRL data

LabKey API



Combine queries between tables



Work in a variety of programming languages
(support for Python and R) using SQL
queries



Replicate queries between releases and
analyses



Work locally and on the HPC

LabKey .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

Programming tools in the RE



LabKey API demo

Using the This notebook will walk you through the steps to set up this notebook for your environment.

Genomics England Research Environment

- [Getting started](#)
- 📄 [How-to guides](#)
- 📄 [Data in the Research Environment](#)
- 🔧 [Tools in the Research Environment](#)
- 🚀 [High Performance Cluster \(HPC\)](#)
- 🔗 [Workflows, scripts and containers](#)
- 🔒 [Data security and Airlock](#)
- 🎓 [Training](#)

🕒 February 12, 2024

ANNOUNCEMENTS

Current data release:

- [/main-programme/main-programme_v18_2023-12-21](#)
- [nhs-gms/nhs-gms-release_v2_2023-02-28](#)

12th February 2024 - [CloudOS access in main RE](#)

- If you are registered for [CloudOS](#) access, you can now access it from the [RE desktop](#), and easily transfer files between CloudOS and the workspace filesystem.

7th February 2024 - [HPC change](#)

- We will be updating to a [new HPC in spring](#).

15th January 2024 - [How to Guides](#)

- Check out the new [How-to Guides](#) section for end-to-end guides, code books and task-based tutorials.

```

[1]: import numpy as np
import functools
import labkey
import pandas as pd

[2]: def labkey_to_df(sql_query, database, maxrows):
    """generate an dataframe from a labkey query"""
    Args:
        sql_query (str): SQL query to execute
        dr (str): GEL database name
    """
    ver = labkey.get_version()

    if ver == '1.2.0':
        server_context = labkey.get_server_context()
        domain = labkey.get_domain()
        container_path = labkey.get_container_path()
        context_path = labkey.get_context_path()
        use_ssl = labkey.get_use_ssl()

    results = labkey.query.execute_sql(
        server_context,

```

6. Running command line tools and pipelines on the HPC

What is an HPC?



What is an HPC?

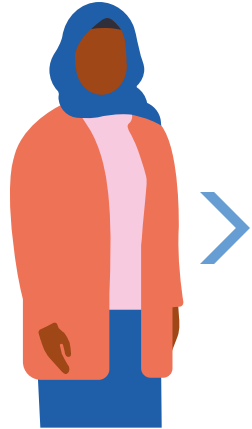


Lots of compute power



Shared with other researchers

What is a "job"?

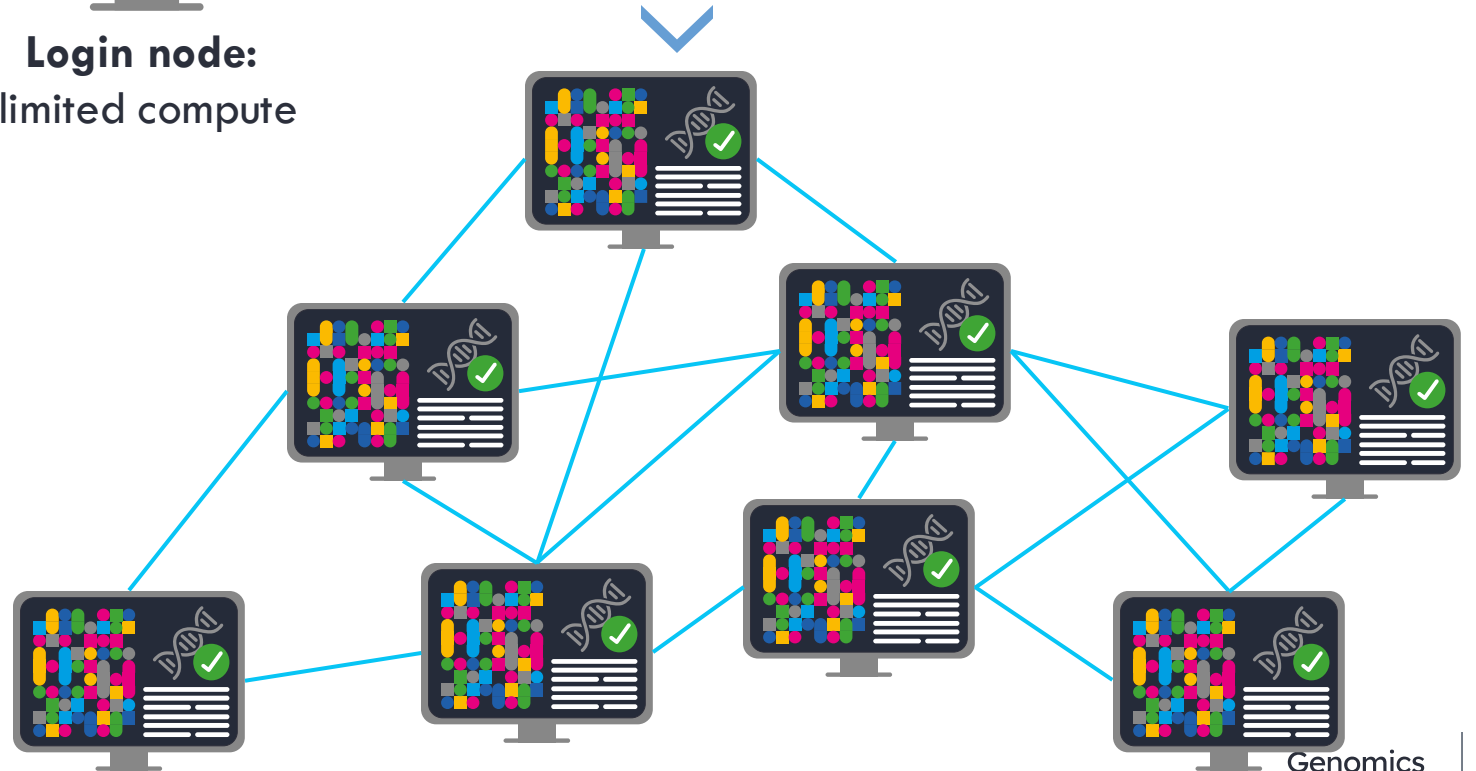


Login to HPC

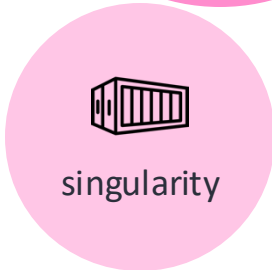
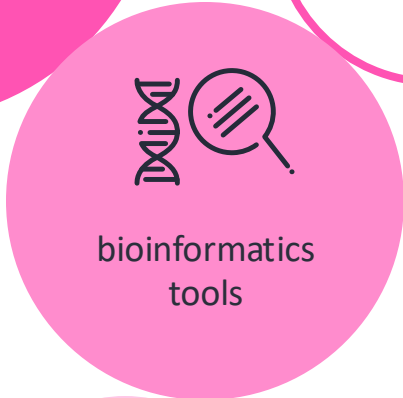


Login node:
limited compute

Create a job



Software on the HPC



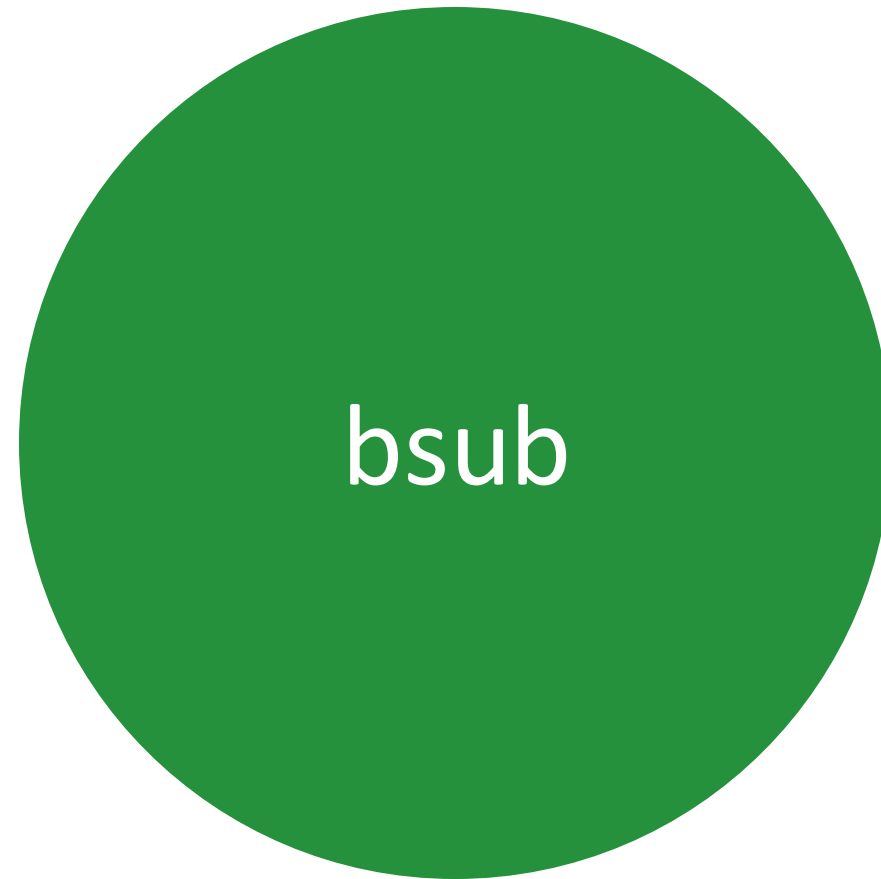
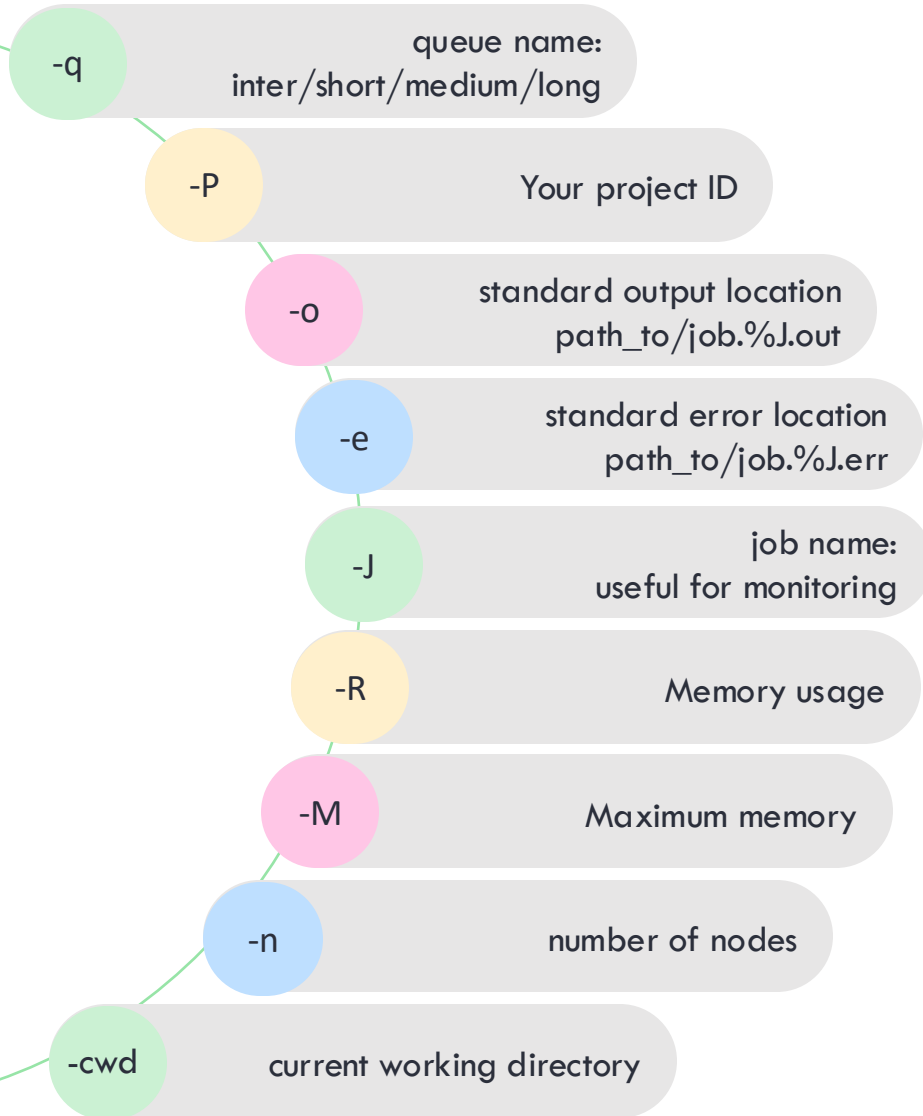
AdapterRemoval/2.3.3
aliview/1.28
ampliconArchitect/1.3.r7
ampliconClassifier/1.1.1
annotSV/3.3.7
annovar/2019Nov
annovar/2024-03-14
ant/1.9.16
apbs/3.4.1
asmc-asmc/2024-02-26
AutoDock_Vina/1.2.5
automake/1.15
aws-cli/2.15
bamtools/2.5.2
bcftools/1.16
beagle/5.4
bedops/2.4.41
bedtools/2.30.0
bedtools/2.31.0
BerkeleyDB/3.01
Bio-DB-HTS/3.01
blast+/2.15
blat/1.0
bolt-imm/2.4.1
boost/1.83
bowtie2/2.5.2
BWA/0.7.17
CADD/1.6
canvas/1.40.0.1613
CaVEMan/1.15.3
circo/0.69-9
clang/16.0.6
cmake/3.24.3
CNView/1.0
CNVnator/0.4.1
cpan/1.7047
cromwell/v65
curl/7.81.0
cython/3.0.8
cytoscape/3.10.1
delly/1.2.6
denovoGear/1.1.1
discover/0.9.5
dotnet/2.0.0
dotnet/8.0.1
drop/1.2.4
eigen/3.3.9
exomiser/13.3.0

exonerate/2.2.0
ExpansionHunter/3.2.2
ExpansionHunter/4.0.2
ExpansionHunterDenovo/0.9.0
fastqc/0.12.1
fetk/1.9.3
ffmpeg/6.0
fribidi/1.0.12
gatk/4.5.0.0
gauchian/1.0.2
gcc/10.4.0
gcta/1.94
gdal/3.7.0
geos/3.12.1
gistic/2.0.23
gmp/6.2.1
gnu-parallel/20190222
gnu/4.4
gradle/8.5
GSL/2.7
guppy/3.4.5
gvcfgenotyper/2019.02.26
haplocheck/1.3.3
hipstr/0.7
hisat2/2.2.1
hla-la/1.0.3
hmftools/2024-02-06
homer/4.11
htslib/1.18
igv/2.17.1
imagemagick/7.1.0
java/1.8
java/11.0.2
java/17.0.2
java/19.0.2
jq/1.7.1
kallisto/0.50.1
king/2.3.2
kraken/1.1.1
kraken2/2.1.3
lapack/3.12.0
ldsc/1.0.1
ldstore/2.0
libdeflate/1.20
libgit2/1.6.2
libgit2/1.6.2
libtiff/3.4
libtiff/4.3.0
libtiff/4.5.0

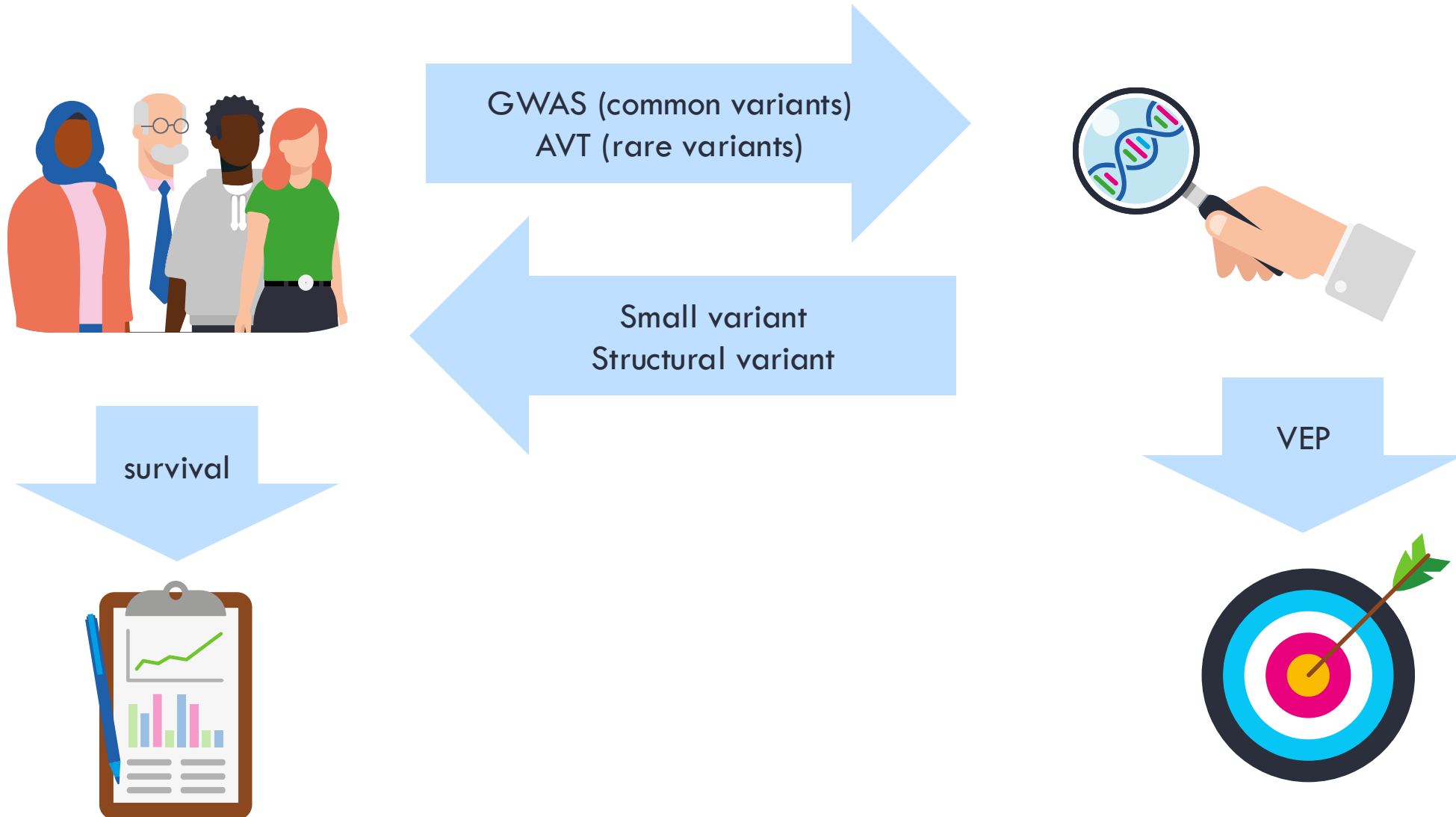
libunwind/1.8.0
liftover/1.0
linasm/1.13
llvm/16.0.6
locuszoom/1.4
lolipop/0.3.0
lumpy/0.3.1
mafft/7.520
magma/1.10
manta/1.6.0
matlab/24.1
matlab/8.1
maven/3.9.6
MEDICC2/1.0.2
meme/5.5.5
metal/1.0
miniconda3/23.11.0
miniforge3/23.11.0-0
minimap2/2.26
mosaicHunter/2024-02-14
MPFR/4.2.0
mplayer/1.5
msisensor-pro/1.2.0
msisensor/0.6
multiqc/1.19
music2/0.2
mutserve/2.0.0-rc15
mutsig2cv/3.11
ncurses/6.4
new_fugue/2010-06-02
nextflow/22.10.5
nextflow/23.04
nextflow/23.10
nextflow/23.10-with-plugins
nextflow/24.04.2-with-plugins
nf-core/0.3.1
nf-test/0.7.3
nf-test/0.8.2
nf-test/0.9.0
nodejs/16.9.0
openrefine/3.7.4
openssl/1.1.1o
pandoc/3.3
perl/5.38.2
picard/3.1.1
pindel/0.2.5b8
platypus/0.8.1
plink_seq/0.10
plink/1.9

plink/2.0
plink/2.00a3.3LM
popdel/1.5.0
proj/8.2.1
prsize-2/2.3.5
pycircos/1.0.2
pysam/0.22.0
python/3.11
python/3.8
python/3.8.1
R/3.6.3
R/4.2.1
R/4.3.3
readline/8.0
regenie/3.4.1
repeatDetector/1.0
REViewer/0.2.7
rtg-tools/3.12.1
rvtests/2.1.0
saige/1.0.9
salmon/1.10.0
samtools/1.16.1
shapeit4/4.2.2
singularity/3.8.3
singularity/4.1.1
sniffles/1.0.11
somalier/0.2.19
sqlite3/3.40.0
squirls/2.0.1
stack/2.15.7
star/2.7.11a
star/2.7.2a
nextflow/23.10
nextflow/23.10-with-plugins
nextflow/24.04.2-with-plugins
nf-core/0.3.1
nf-test/0.7.3
nf-test/0.8.2
nf-test/0.9.0
nodejs/16.9.0
openrefine/3.7.4
openssl/1.1.1o
pandoc/3.3
perl/5.38.2
picard/3.1.1
pindel/0.2.5b8
platypus/0.8.1
plink_seq/0.10
plink/1.9

Creating a job - parameters



Pre-built workflows/scripts to...



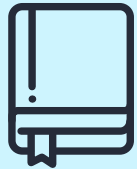
Workflows/scripts provide



Code that runs with only minor tweaks to add your input



Optimised for use on our HPC and with our data



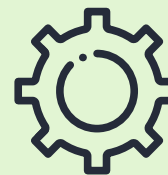
Step-by-step instructions for use



Output in standard or easy-to-interpret formats



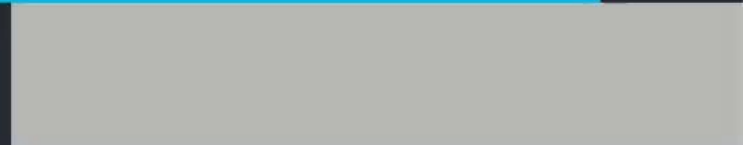
Example input data and submission scripts



Customisable options

HPC demo

- Computer
- Git GUI
- GVim
- Research Registry
- eperry's Home
- IGV Browser
- RStudio
- Link to emily
- IVA
- Terminal Emulator
- Old Firefox Data
- Labkey
- Text Editor
- Airlock
- LibreOffice 7.6
- Visual Studio Code
- CloudOS
- Open Targets
- Welcome Pack
- Desktop.Rproj
- Panel App
- Trash
- Document Viewer
- Participant Explorer
- Emacs
- R
- Ensembl
- RE Messages
- Firefox
- Research Environment Documentation



7. Import and export of data and tools



The Airlock

Data in the RE



Outside world

Our contract with participants

"...although researchers can look at your data and ask questions about it, they can only take away the answers to their questions (their results). They can't copy or take away any of your individual data."



Forms in Airlock



Export findings



Export analysis scripts and software



Contact clinical team and/or report potential diagnosis



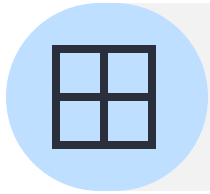
What should go through the Airlock?



Figures



Statistics or numbers for your text



Tables



Notes on the data

Airlock rules

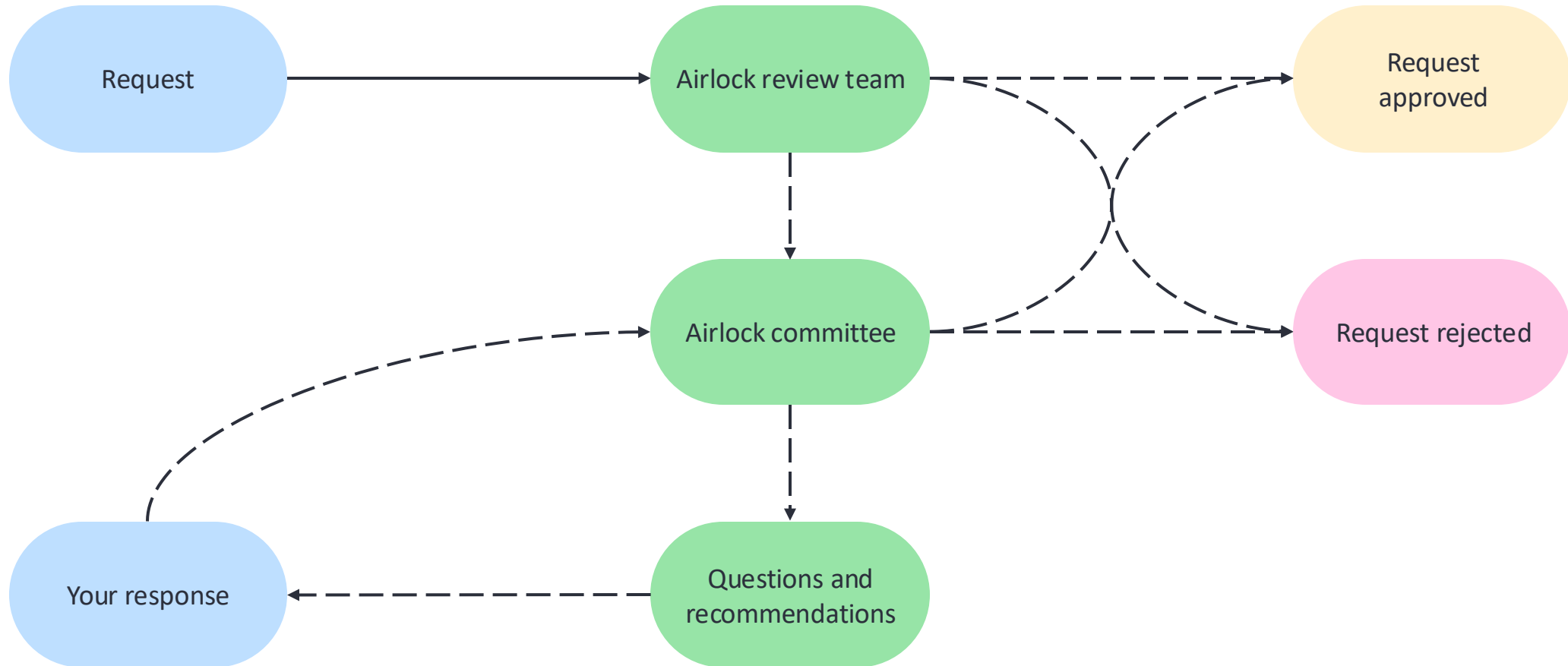


Approved research
project



Participants cannot be
identified

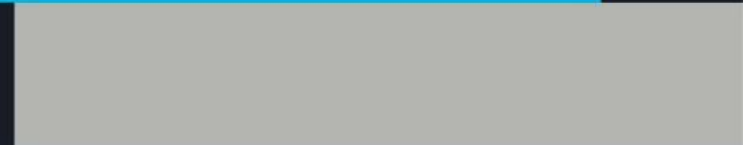
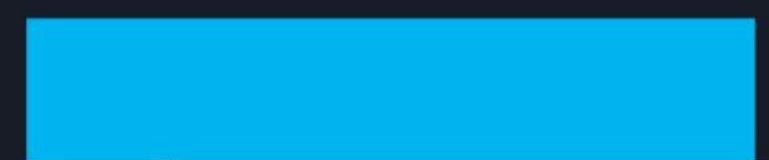
What happens to my request?



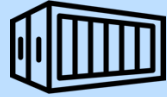
Airlock demo



- Computer
- Text Editor
- Airlock
- Research Environment Documentation
- Welcome Pack
- eperry's Home
- Data Discovery
- Participant Explorer
- report.tsv
- Trash
- Firefox
- Visual Studio Code
- Ensembl
- Document Viewer
- IGV Browser
- RE Messages
- IVA 2.0
- R
- Terminal Emulator
- Rocket Chat
- LibreOffice
- Panel App
- Research Registry
- GVim
- Labkey
- Git GUI
- RStudio
- Old Firefox Data
- Open Targets
- Emacs
- LibreOffice 7.6
- CloudOS



Import



Containers



Copy/paste in

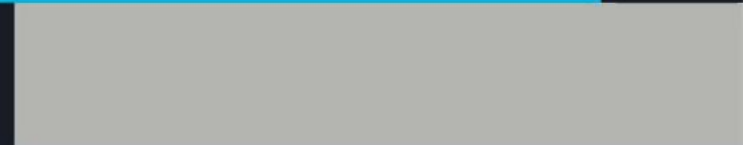
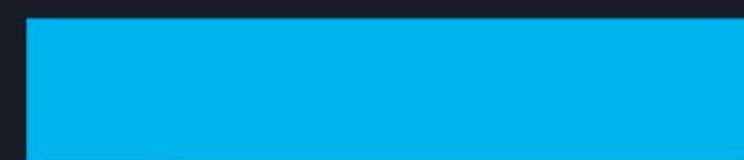


Airlock import

Import demo



- Computer
- Text Editor
- Airlock
- Research Environment Documentation
- Welcome Pack
- eperry's Home
- Data Discovery
- Participant Explorer
- report.tsv
- Trash
- Firefox
- Visual Studio Code
- Ensembl
- Document Viewer
- IGV Browser
- RE Messages
- IVA 2.0
- R
- Terminal Emulator
- Rocket Chat
- LibreOffice
- Panel App
- Research Registry
- GVim
- Labkey
- Git GUI
- RStudio
- Old Firefox Data
- Open Targets
- Emacs
- LibreOffice 7.6
- CloudOS



8. Getting help and questions

Getting help



Check our documentation:
<https://re-docs.genomicsengland.co.uk/>
Click on the documentation icon in the environment



Contact our Service Desk:
<https://jiraservicedesk.extge.co.uk/plugins/servlet/desk>

Thank you

Visit: <https://re-docs.genomicsengland.co.uk/>