# Finding participants based on genotypes

**Emily Perry**

Research Engagement Manager

19th July 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
  - Share your webinar link (we will remove anyone who is here twice)

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

# 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.

Genomics England

# Questions



**Ronnie Rodrigues Pereira**
Bioinformatician -
Research Services



**Alex Stuckey**
Senior Bioinformatician -
Research Services



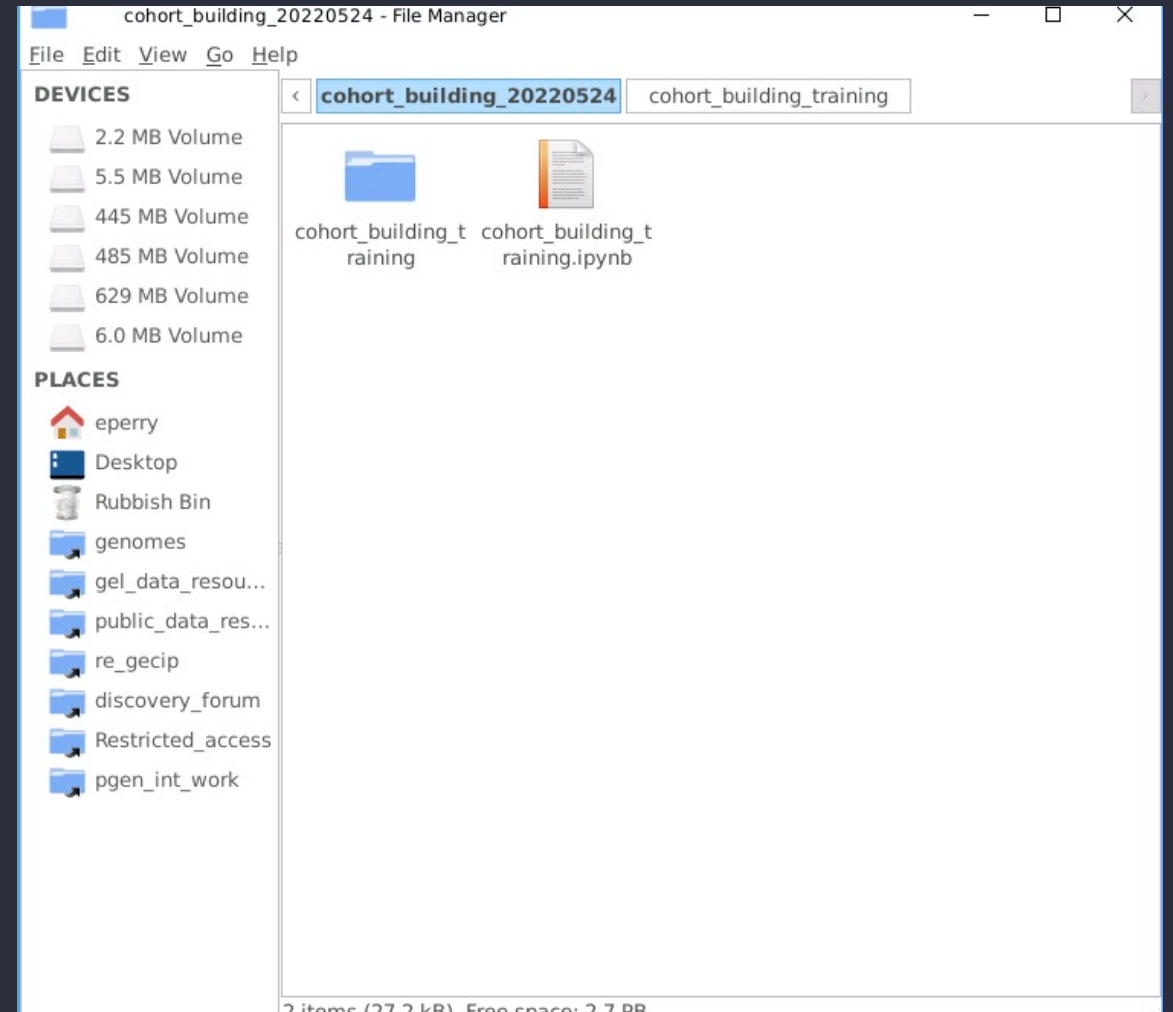**Christian Bouwens**
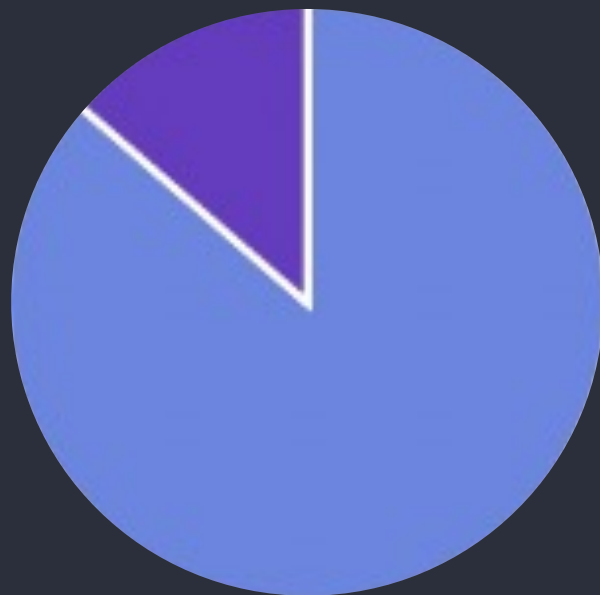Bioinformatician -
Research Services

# Agenda

# Materials

- Slides and video will be sent out to you after the session

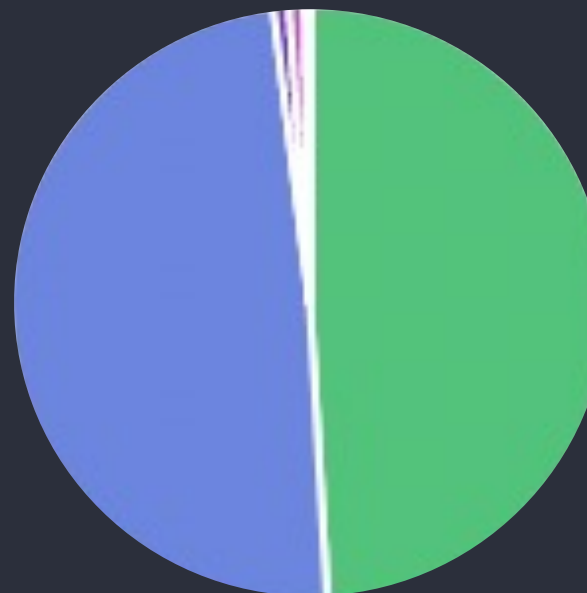- Scripts available in /gel_data_resources/example_scripts/workshop_scripts/cohort_building_20220524

# Genome assembly

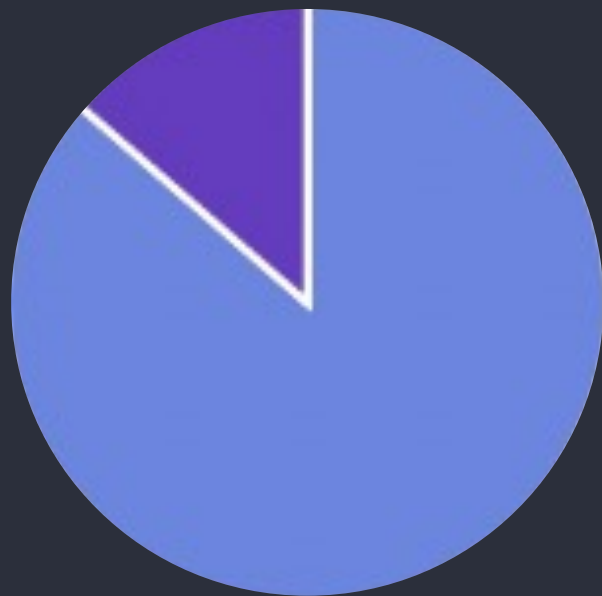## Rare disease



GRCh38 (aka hg38)
GRCh37 (aka hg19)

## Cancer



Somatic GRCh38
Germline GRCh38
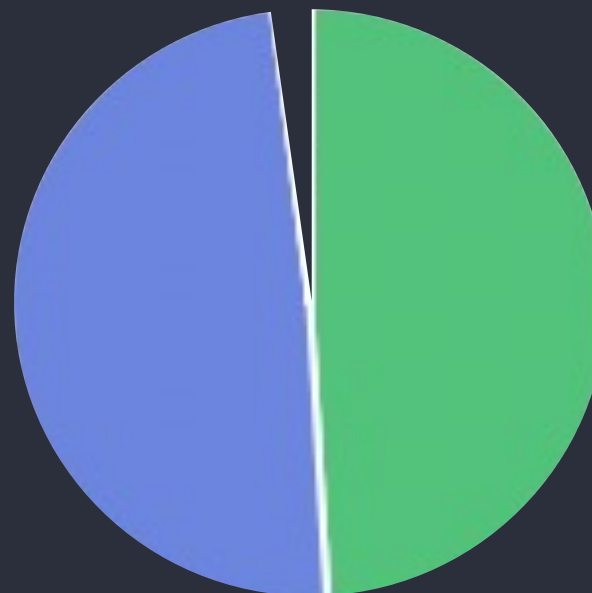Germline GRCh37
Somatic GRCh37

# Genome assembly

## Rare disease



GRCh38 (aka hg38)
GRCh37 (aka hg19)

## Cancer



Somatic GRCh38
Germline GRCh38
~~Germline GRCh37~~
~~Somatic GRCh37~~

Genomics
England

# Genome assembly



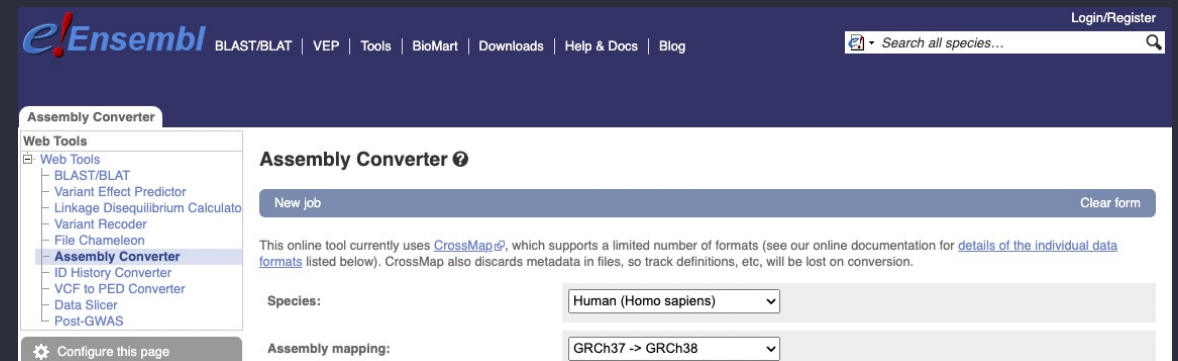chr13 ZAR1L ENST00000345108.6:c.931T>C

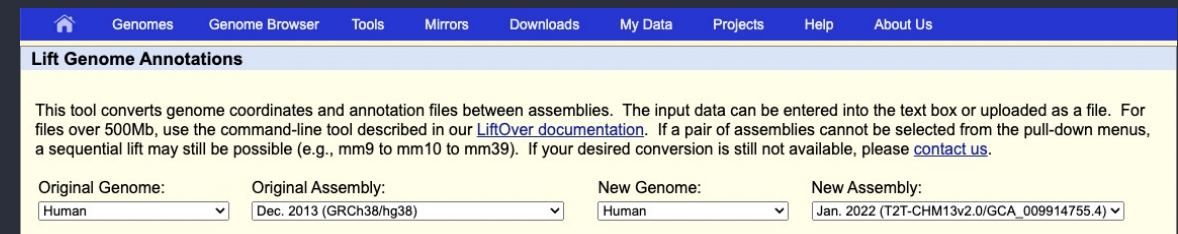| | GRCh37 (hg19) | GRCh38 (hg38) |
|---|---|---|
| ZAR1L | 13:32,877,837-32,889,481 | chr13:32,303,699-32,315,363 |
| ENST00000345108.6:c.931T>C | 13:32,878,051 | chr13:32,303,914 |

# Converting coordinates between assemblies

## Inside the RE:

- Liftover tool on HPC

- Chain files in public_data_resources

## Outside the RE:



https://www.ensembl.org/Homo_sapiens/Tools/AssemblyConverter



https://genome.ucsc.edu/cgi-bin/hgLiftOver

# 2. LabKey tables of variant genotypes

# Labkey

- Participant details and family relationships

- Sample details

- Genomic file locations

- Clinical data

- Bioinformatics analysis results
  - Tiering
  - Exomiser

# Rare disease tiering

# PanelApp

Genomics England

# Exomiser

# Exomiser/tiering assembly

## Rare disease



GRCh38 (aka hg38)
GRCh37 (aka hg19)

## Cancer

~~Somatic GRCh38~~
~~Germline GRCh38~~
~~Germline GRCh37~~
~~Somatic GRCh37~~

Genomics
England

# Exomiser/tiering assembly

chr13      ZAR1L      ENST00000345108.6:c.931T>C

| Search by | |
|---|---|
| gene | Should find all filter-passing variants in the gene on either assembly |
| coordinate(s) | You must also specify the genome assembly |
| HGVS (exomiser only) | Should find all filter-passing variants that match your string on either assembly |

# Demo: Variant data in LabKey

# LabKey API

Labkey API allows you to:

- Combine data and filters from multiple tables
- Work in a variety of programming languages, but most support for Python and R
- Work both locally and on the HPC

# 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 copying and editing a file called .netrc

# Materials

- Slides and video will be sent out to you after the session

- Scripts available in /gel_data_resources/example _scripts/workshop_scripts/gen otypes_20220719

# Accessing the notebooks

Python

```
module load python/3.8.1
jupyter notebook --port
<four digit port number>
```

Open link in browser

R

```
module load R/4.0.2
rstudio
```

Genomics
England

# Demo: Variant data in LabKey API

# 3. Finding genotypes with IVA

# IVA Variant Browser

- Point-and-click interface to explore variants

- Filter by loci, consequences, population frequencies and inheritance.

- Find participant genotypes.

# IVA genome assembly

## Rare disease



OR

GRCh38 (aka hg38)
GRCh37 (aka hg19)

## Cancer



OR

Somatic GRCh38
Germline GRCh38
Germline GRCh37
Somatic GRCh37

Genomics
England

39

# Demo: Finding variants with IVA

# 4. The Gene-Variant and SV/CNV workflows

# Gene-Variant workflow

- Submit a list of genes or regions

- Find all short variants in these genes/regions

- Get participants with these variants

- Choose somatic/germline, cancer/rare disease



① **Make a folder in your own directory**

`mkdir /home/username/my_workflows/`
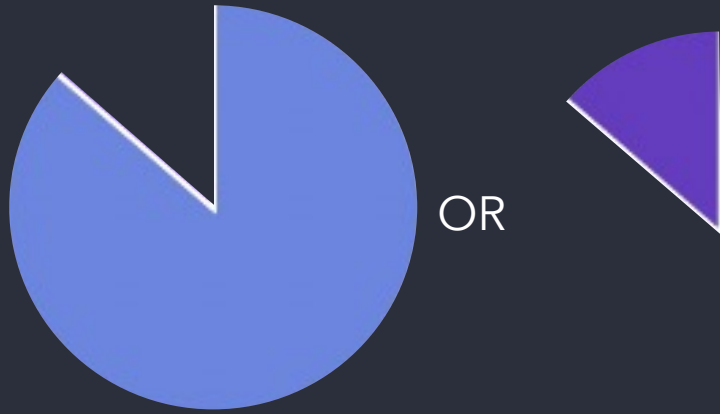
② **Copy workflow to your new folder**

`cp -R /gel_data_resources/workflows/BRS_tools_geneVariantWorkflow/1.6 \`
`/home/username/my_workflows/`

③ **Minimal editing of three files**

gene_list.txt

variant_workflow_inputs.json

submit_workflow.sh

```
nano gene_list.txt               change gene names
nano variant_workflow_inputs.json change project
nano submit_workflow.sh           change project and filepath of workdirectory (-cwd)
```

To save after editing Ctrl+O, type Y to overwrite, and Ctrl+X to exit

④ **Submitting your workflow to the HPC**

`bsub < submit_workflow.sh`

# SV-CNV workflow

- Submit a list of genes or regions

- Find all SVs/CNVs in these genes/regions

- Get participants with these variants

- Choose somatic/germline, cancer/rare disease

# Benefits of using workflows

SV/CNV WORKFLOW IS THE
<span style="color:yellow">ONLY</span> WAY TO GET SV/CNVS
ACROSS ALL PARTICIPANTS
IN RE

BOTH WORKFLOWS QUERY
GENOMES ALIGNED TO
GRCH37 <span style="color:yellow">AND</span> GRCH38

# Workflows and genome assembly

chr13          ZAR1L     ENST00000345108.6:c.931T>C

| Search by | |
|---|---|
| gene | Should find all variants in the gene on either assembly |
| coordinate(s) | You must also specify the genome assembly |

# Demo: Running workflows on the HPC

# 5. Aggregated variant files

# AggV2

combined germline gVCFs → Split into 1371 chunks → Release 8 → GRCh38 → Quality metrics → Functional annotation

⚠️

- No genomes aligned to GRCh37
- Includes participants who have since withdrawn consent – see docs for how to remove

# AggV2 assembly

## Rare disease



## Cancer



GRCh38 (aka hg38)
~~GRCh37 (aka hg19)~~

~~Somatic GRCh38~~
Germline GRCh38
~~Germline GRCh37~~
~~Somatic GRCh37~~

# somAgg

combined somatic VCFs → Split into 1371 chunks → Release 12 → GRCh38 → Quality metrics → 0/0 (not present) or 0/1 (present)

⚠️

- No genomes aligned to GRCh37
- Includes participants who have since withdrawn consent – see docs for how to remove
- Feedback on somAgg

# somAgg assembly

## Rare disease

## Cancer

~~GRCh38 (aka hg38)~~
~~GRCh37 (aka hg19)~~

Somatic GRCh38
~~Germline GRCh38~~
~~Germline GRCh37~~
~~Somatic GRCh37~~

Genomics
England

# Aggregated files code books

- AggV2 and somAgg can be analysed with:
  - Bcftools
  - Split-vep
  - R/Python
  - Bedtools
- Code books available to take you through most common use-cases:
  - AggV2 https://research-help.genomicsengland.co.uk/display/GERE/aggV2+Code+Book
  - somAgg https://research-help.genomicsengland.co.uk/display/GERE/somAgg%3A+Code+Book

# Aggregated files chunks

**Split into 1371 chunks**

- Locus-based queries must query the correct chunk file

- BED file of chunks available

- Create a sorted list of your own regions

- Intersect with BEDtools

- Code books with more information

- Also available in Plink2 format

# 6. Using bcftools on the HPC

# Demo: Using bcftools on the HPC

# Summary and comparison of tools

# Comparison – access type

| | |
|---|---|
| Labkey tables: Tiering and Exomiser | Point-and-click access to tables plus API |
| IVA | Point-and-click |
| Gene-variant and SV/CNV workflows | Command line |
| Aggregated VCFs with bcftools | Command line |

# Comparison – search by

| | |
|---|---|
| Labkey tables: Tiering and Exomiser | Gene, region or HGVS (Exomiser only) |
| IVA | Gene, region or rsID |
| Gene-variant and SV/CNV workflows | Gene or region |
| Aggregated VCFs with bcftools | Region |

# Comparison – variants available

| | |
|---|---|
| Labkey tables: Tiering and Exomiser | Only variants that have passed tiering or exomiser filters |
| IVA | All variants |
| Gene-variant and SV/CNV workflows | All variants |
| Aggregated VCFs with bcftools | All variants present in GRCh38-aligned genomes from release 8 (AggV2) or 12 (somAgg) |

# Comparison – genome assembly

| | |
|---|---|
| Labkey tables: Tiering and Exomiser | Variants on both assemblies in the table, with assembly as a column |
| IVA | Assemblies available as separate datasets |
| Gene-variant and SV/CNV workflows | By-gene searches both assemblies |
| Aggregated VCFs with bcftools | Only GRCh38 |

# Comparison – underlying VCFs

| Labkey tables: Tiering and Exomiser | Platypus |
|---|---|
| IVA | Platypus |
| Gene-variant and SV/CNV workflows | Strelka |
| Aggregated VCFs with bcftools | Strelka |

# Key takeaways

Use IVA for a fast overview

Pre-written workflows for gene-based searches

Aggregated VCFs have code-books for common use-cases

# Key takeaways

Use IVA for a fast overview

! Genome Assembly

Aggregated VCFs have code-books for common use-cases

# 7. Getting help and questions

# 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

# 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

Getting medical history for participants

22 Nov.

20 Sep.

Using the HPC to run jobs

Genomics
England

# Feedback

Thank you