

Finding participants based on genotypes

Emily Perry

8th July 2025



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

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

Questions



Magdalena Drożdż Bioinformatician -Research Services

Agenda

1	Introduction and admin
2	Genome assembly
3	LabKey tables of variant genotypes
4	Finding genotypes with IVA and Cohort Browser
5	The Small Variant and Structural Variant workflows
6	Aggregated variant files
7	When/why you would use each method
8	Help and questions

Materials

- Slides and video will be sent out to you after the session
- Scripts available in /gel_data_resources/examp le_scripts/workshop_scrip ts/genotypes_2025



2. Genome Assembly

100,000 Genomes Project



100,000 Genomes Project



NHS GMS



Germline GRCh38 (hg38) Somatic GRCh38 (hg38)

Cancer

Genome assembly coordinates



	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 between assemblies

Inside the RE

- Liftover tool on HPC
- Chain files in public_data_resources

Outside the RE

- Ensembl Assembly converter
- UCSC Liftover

3. LabKey tables of variant genotypes

LabKey

- Participant details and family relationships
- Sample details
- Genomic file locations
- Clinical data
- Bioinformatics analysis results
 - Rare disease tiering
 - Cancer tiering
 - Exomiser



Rare disease tiering



Rare disease tiering based on PanelApp genes

List 🛧	Entity	Reviews	Mode of inheritance	Details	
Filter Ent	ities				8 Entities
Green	ATP1A3	1 review 1 green	MONOALLELIC, autosomal or pseudoautosomal, NOT imprinted	Sources • Expert Review • Expert Review Green Phenotypes • 601338 • 614820 Tags	
Green	DFNB59	2 reviews 1 green	BIALLELIC, autosomal or pseudoautosomal	Sources Expert Review Expert Review Green Phenotypes 610219 Tags new-gene-name	
Green	OPA1	2 reviews 1 green	MONOALLELIC, autosomal or pseudoautosomal, NOT imprinted	Sources Eligibility statement prior genetic testing Expert Review Green Phenotypes Optic atrophy 1, OMIM:165500 Optic atrophy plus syndrome, OMIM:125250 Tags	
Green	OTOF	1 review 1 green	BIALLELIC, autosomal or pseudoautosomal	Sources Expert Review Green Radboud University Medical Center, Nijmegen Phenotypes 601071 Tags	
Amber	DIAPH3	3 reviews 1 red	BOTH monoallelic and biallelic, autosomal or pseudoautosomal	Sources Expert Review Amber Radboud University Medical Center, Nijmegen Phenotypes Auditory neuropathy, autosomal dominant, 1, 609129 Tags	

Rare disease Exomiser

Exomiser





100k tiering



Cancer Germline GRCh38 (hg38) Somatic GRCh38 (hg38)

GRCh38 (hg38)

Tiering tables genome assembly – 100k



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

NHS GMS tiering

Rare disease



Variant data in LabKey demo



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

Materials

- Slides and video will be sent out to you after the session
- Scripts available in /gel_data_resources/examp le_scripts/workshop_scrip ts/genotypes_2025



Variant data in LabKey API demo

Amazon WorkSpaces View Settings Connections Support

.

• .

0

•

1

-

-

Ð

1.2

C

6

•

m

-

-

.

12.0

1

1.00

5

1

📥 🔘 🌔 🍓 🗾 📕 🔯 🖇 📫 🗢 Q 🚍 Tue 20 Jun 09:57



4. Finding genotypes with IVA and Cohort Browser

Interactive variant analysis (IVA)

- Point-and-click interface to explore variants
- Filter by loci, consequences, population frequencies and inheritance
- Find participant genotypes



100k in IVA



GRCh37 (hg19) GRCh38 (hg38) Germline GRCh38 (hg38) Somatic GRCh38 (hg38)

IVA demo



5. The Small Variant and Structural Variant workflows

Small variant workflow

Submit a list of genes

Find all short variants in these genes

Get 100k participants with these variants

Structural variant workflow

Submit a list of genes or regions

Find all structural variants overlapping these genes

Get 100k participants with these variants

Workflows genome assembly – 100k

Search by	
gene	Should find all variants in the gene(s) on either assembly
coordinates (structural only)	You must also specify the genome assembly

Running workflows on the HPC demo

. . .

-

A.

0

•

.

1

-

0

C

m

-

.

-

-

IF .

- 6

1

Amazon WorkSpaces

6. Aggregate variant files

https://re-docs.genomicsengland.co.uk/aggv2/ https://re-docs.genomicsengland.co.uk/somAgg/

AggV2 – germline samples

Rare disease

GRCh37 (hg19) GRCh38 (hg38) Germline GRCh38 (hg38) Somatic GRCh38 (hg38)

43

https://re-docs.genomicsengland.co.uk/somAgg/

Aggregate VCFs

https://re-docs.genomicsengland.co.uk/aggv2_code_book/ https://re-docs.genomicsengland.co.uk/somAgg_code_book/

Aggregate VCF chunks

- Locus-based queries must query the correct chunk file
- BED file of chunks available
- Create a sorted BED file of your own regions
- Intersect with BEDtools
- Code books with more information
- Also available in Plink2 format

https://re-docs.genomicsengland.co.uk/aggv2_code_book/ https://re-docs.genomicsengland.co.uk/somAgg_code_book/

Using bcftools on the HPC demo

7. When/why you would use each method

Genomics England

Genome assembly

Tiering and exomiser tables

GRCh37 and GRCh38 Assembly as a separate column IVA

GRCh37 and GRCh38 in separate databases

Small/ Structural variant workflows

GRCh37 and GRCh38 queries simultaneously

Querying the aggregates

GRCh38 only

Tiering and exomiser tables

Rare disease: Platypus Cancer: Strelka

> Small/ Structural variant workflows

> > Strelka

Rare disease: Platypus Cancer: Strelka

IVA

Querying the aggregates

Strelka

Key takeways

8. Help and questions

Getting help

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

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

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

Introduction to the RE

Materials from past training all online

Training sessions

9/9	Getting medical records for participants
14/10	What tools and workflows should I use to fulfil an overall goal?
11/11	Using GEL data for publications and reports
9/12	Running workflows on the HPC and Cloud

Materials from past training all online

Feedback

Thank you

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