

Building cancer cohorts

Emily Perry

Research Engagement Manager 13th May 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
- If you are joining virtually, 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)

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





Asier Gonzalez-Uriarte Bioinformatician in Data Release Team (DRT)

Matthieu Vizuete-Forster Bioinformatician -Research Services



Agenda

1	Introduction and admin
2	Parameters and considerations for building a cohort
3	Point-and-click cohort building with Participant Explorer
4	Tables for cohort building in cancer
5	Programmatic cohort building in Python and R
6	Getting genomic filepaths for your cohort
7	Using your cohort with aggregate VCFs
8	Help and questions



2. Parameters and considerations for building a cohort





100,000 Genomes cancer



NHS GMS cancer



Cancer characteristics



Cancer treatment



Derive age from dates



Demographics





Smoking Testicular cryptorchidism HRT Endometriosis Menarche/menopause Sunbed use Pregnancy/breastfeeding Dialysis IVF Hepatitis/cirrhosis Height/weight Alcohol Tamoxifen

3. Point-and-click cohort building with Participant Explorer

Participant Explorer

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



Participant Explorer demo



4. Tables for cohort building in cancer

Cancer cohort parameters



Cancer cohort parameters



Derive age from dates





Data dictionary



Tables demo



5. Programmatic cohort building in Python and R

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



Demo notebooks



/gel_data_resources/example_scripts/ workshop_scripts/cancer_cohorts_2025

Programming demo



::: (0)

1

0

1

×

8

•

Amazon WorkSpaces



6. Getting genomic filepaths for your cohort

Genomic files available



Genomes location



File locations

Participant ID	Platekey	Filepath	Filename	File sub- type
12345678	LP12345- DNA_A01	/genomes/by_date/2025-04-08/ 000111222333/LP12345-DNA_A01/ Variations/LP12345-DNA_A01.vcf.gz	LP12345- DNA_A01.vcf.gz	Genomic VCF
12345678	LP12345- DNA_A01	/genomes/by_date/2025-04-08/ 000111222333/LP12345-DNA_A01/ Assembly/LP12345-DNA_A01.cram	LP12345- DNA_A01.cram	CRAM
12345678	LP12345- DNA_A01	/genomes/by_date/2025-04-08/ 000111222333/LP12345-DNA_A01/ Variations/LP12345- DNA_A01.SV.vcf.gz	LP12345- DNA_A01.SV.vcf. gz	Structural VCF

Filepaths demo

Recruited Disease	Year of Birth	Phenotypic Sex	Ethnic Category	Life Status	Genome Build F	amily Group Type/Size	Family ID/Refe	rral ID	
					CDOL00				

7. Using your cohort with aggregate VCFs

Aggregate VCFs

















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

Aggregate demo



000

0

•

4

•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•
•<

C

m _

(.

A

8

100 8

Amazon WorkSpaces

💽 🙆 📥 🮯 🖪 🖻 📴 🛱 🕼 🕑 📾 穼 🔍 😭 Tue 1 Apr 13:26



8. Getting 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>







Introduction to the RE



Materials from past training all online

Training sessions

 8/7 Finding participants based on genotypes 9/9 Getting medical records for participants 14/10 What tools and workflows should Luse to fulfil an overall goal? 11/11 Using GEL data for publications and reports 9/12 Running workflows on the HPC and Cloud 	10/6	Building rare disease cohorts with matching controls	
9/9Getting medical records for participants14/10What tools and workflows should I use to fulfil an overall goal?11/11Using GEL data for publications and reports9/12Running workflows on the HPC and Cloud	8/7	Finding participants based on genotypes	
14/10What tools and workflows should I use to fulfil an overall goal?11/11Using GEL data for publications and reports9/12Running workflows on the HPC and Cloud	9/9	Getting medical records for participants	
11/11Using GEL data for publications and reports9/12Running workflows on the HPC and Cloud	14/10	What tools and workflows should I use to fulfil an overall goal?	
9/12 Running workflows on the HPC and Cloud	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> docs.genomicsengland.co.uk/