

Using genomic data to build cohorts

Emily Perry

9th June 2026



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



**Matthieu
Vizquete-Forster**
Learning designer

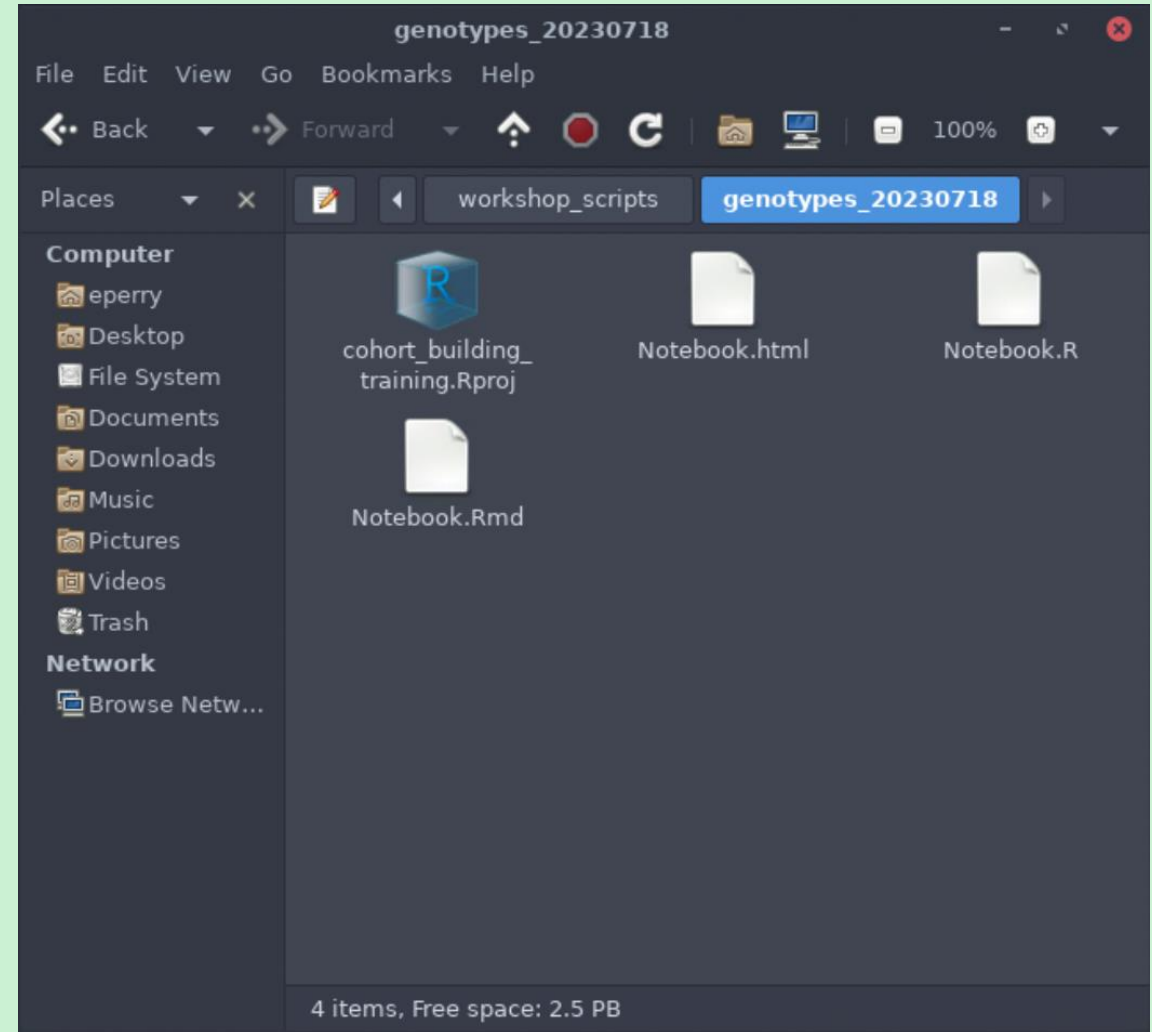
Agenda

- 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/example_scripts/workshop_scripts/genotypes_2026



2. Genome Assembly



100,000 Genomes Project

Rare disease



GRCh37 (hg19)
GRCh38 (hg38)

Cancer



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

100,000 Genomes Project

Rare disease



GRCh37 (hg19)
GRCh38 (hg38)

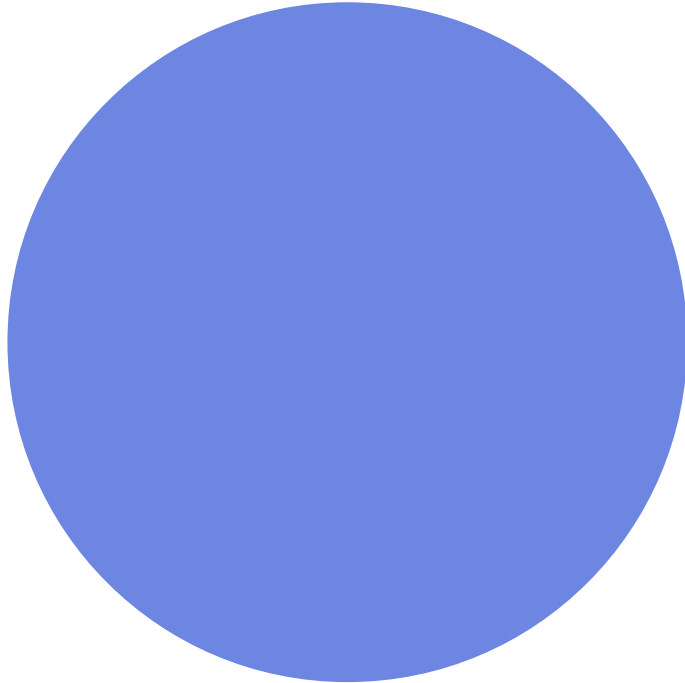
Cancer



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

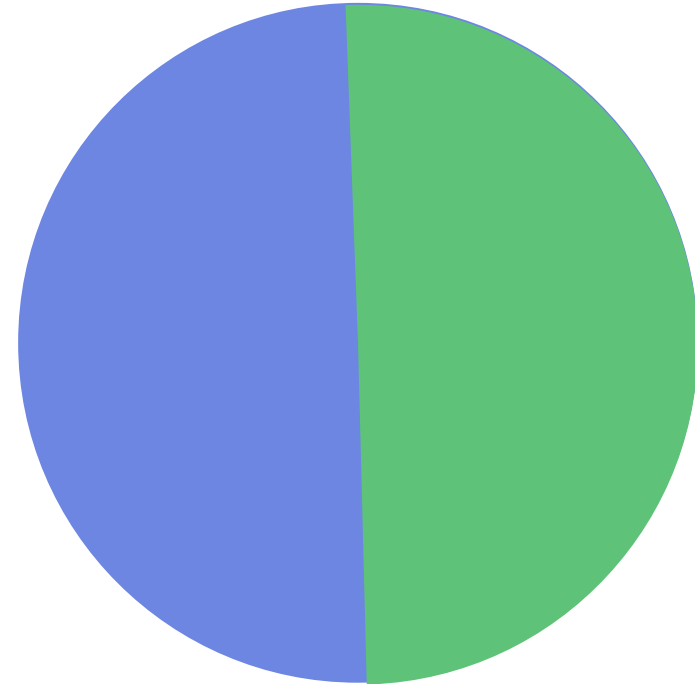
NHS GMS

Rare disease



GRCh38 (hg38)

Cancer



Germline GRCh38 (hg38)

Somatic GRCh38 (hg38)

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

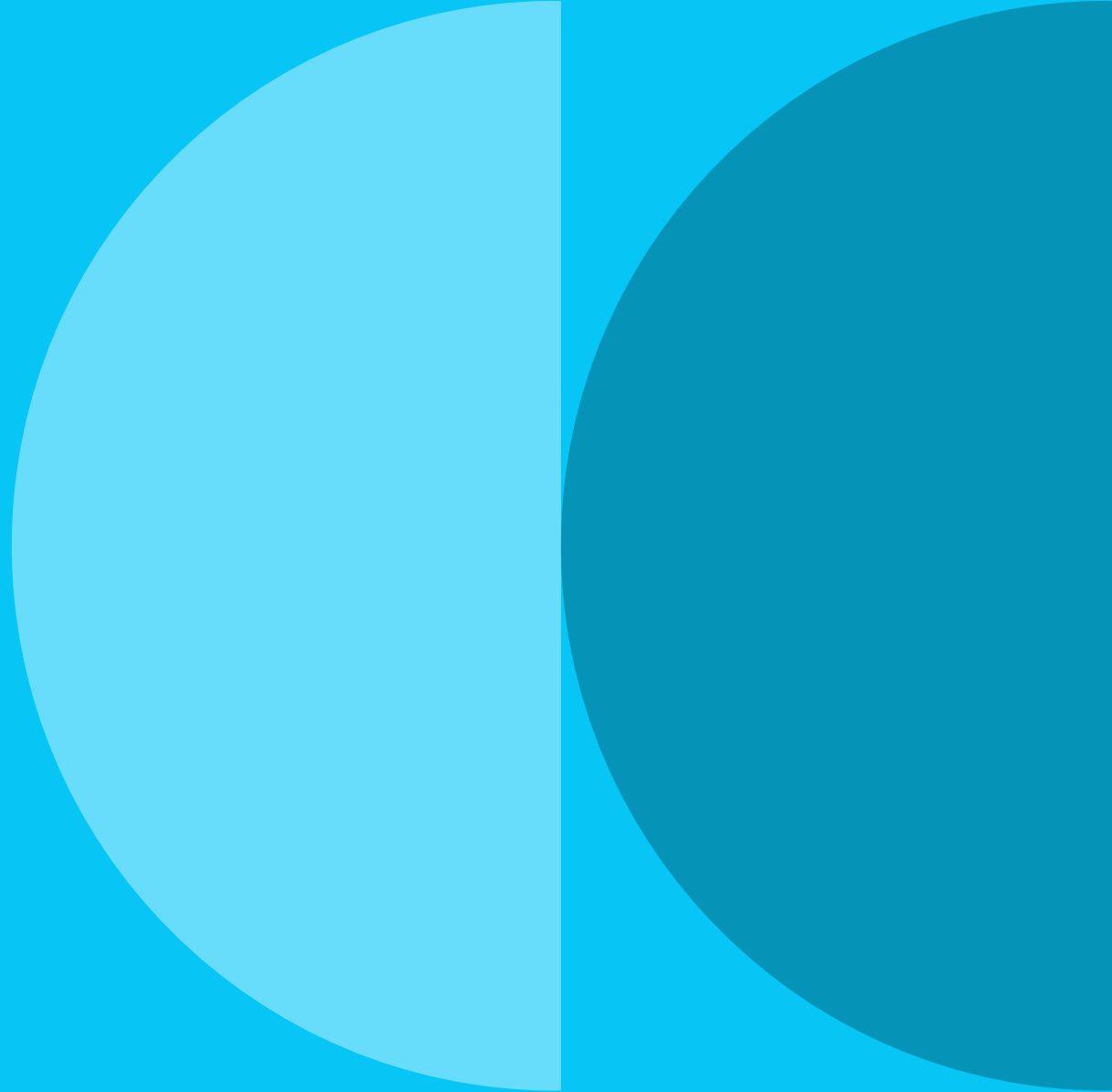
Inside the RE

- Liftover tool on HPC
- Chain files in `public_data_resources`

Outside the RE

- Ensembl Assembly converter
- UCSC Liftover

3. LabKey tables of variant genotypes

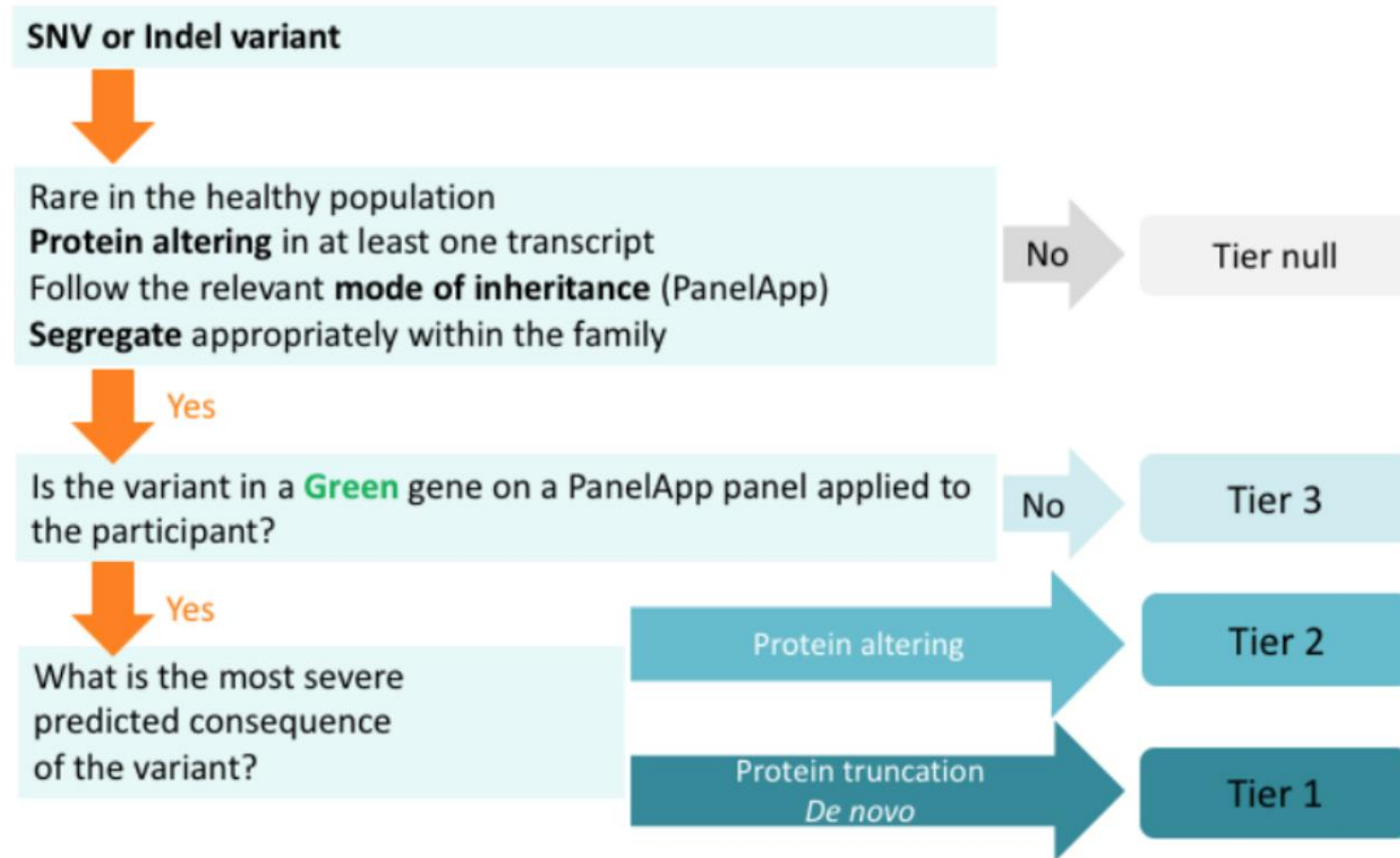


LabKey

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



Rare disease tiering

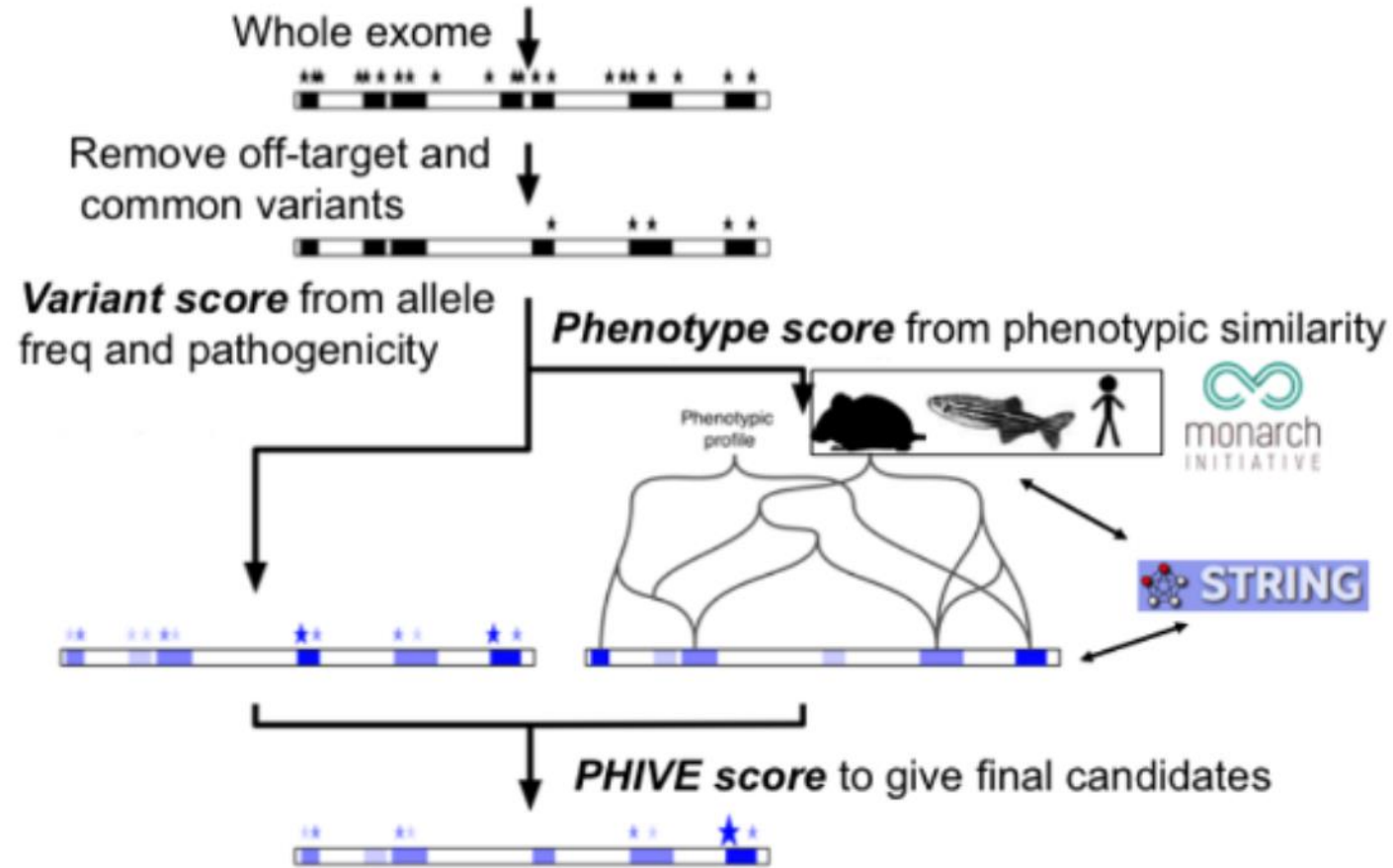


Rare disease tiering based on PanelApp genes

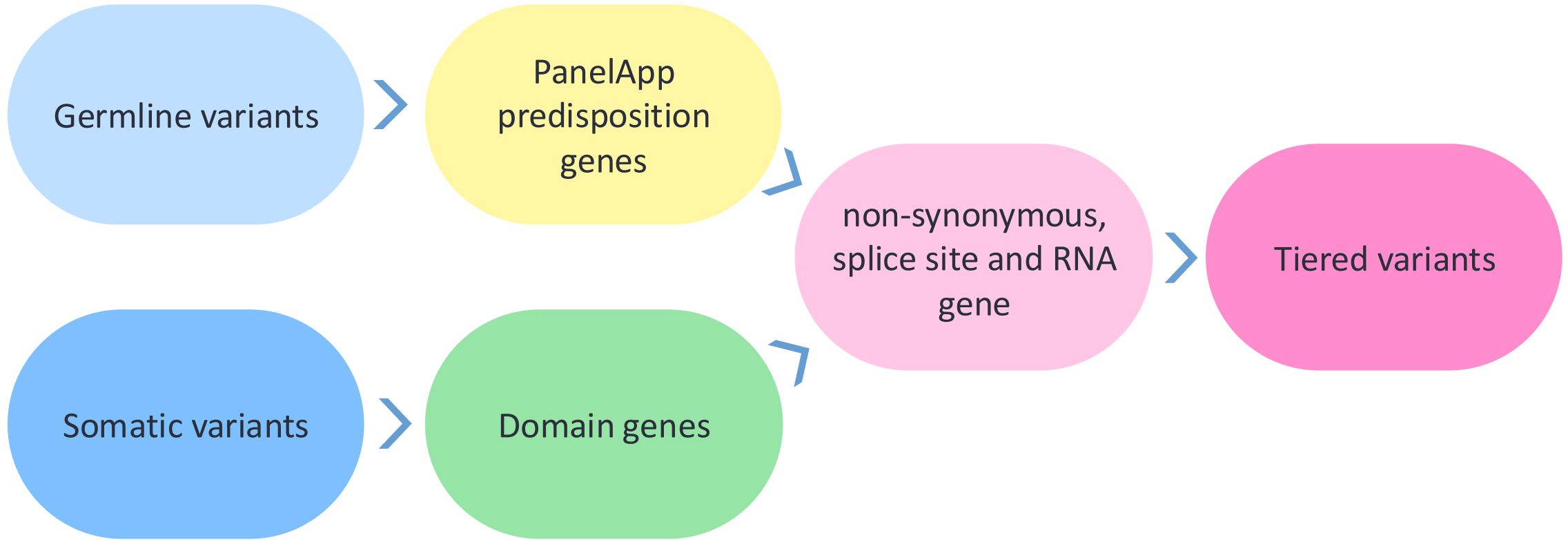
List ↑	Entity	Reviews	Mode of inheritance	Details
	Filter Entities			8 Entities
Green	ATP1A3	1 review 1 green	MONOALLELIC, autosomal or pseudoautosomal, NOT imprinted	Sources <ul style="list-style-type: none"> Expert Review Expert Review Green Phenotypes <ul style="list-style-type: none"> 601338 614820 Tags
Green	DFNB59	2 reviews 1 green	BIALLELIC, autosomal or pseudoautosomal	Sources <ul style="list-style-type: none"> Expert Review Expert Review Green Phenotypes <ul style="list-style-type: none"> 610219 Tags <input type="text" value="new-gene-name"/>
Green	OPA1	2 reviews 1 green	MONOALLELIC, autosomal or pseudoautosomal, NOT imprinted	Sources <ul style="list-style-type: none"> Eligibility statement prior genetic testing Expert Review Green Phenotypes <ul style="list-style-type: none"> Optic atrophy 1, OMIM:165500 Optic atrophy plus syndrome, OMIM:125250 Tags
Green	OTOF	1 review 1 green	BIALLELIC, autosomal or pseudoautosomal	Sources <ul style="list-style-type: none"> Expert Review Green Radboud University Medical Center, Nijmegen Phenotypes <ul style="list-style-type: none"> 601071 Tags
Amber	DIAPH3	3 reviews 1 red	BOTH monoallelic and biallelic, autosomal or pseudoautosomal	Sources <ul style="list-style-type: none"> Expert Review Amber Radboud University Medical Center, Nijmegen Phenotypes <ul style="list-style-type: none"> Auditory neuropathy, autosomal dominant, 1, 609129 Tags

Rare disease Exomiser

Exomiser



Cancer tiering



100k tiering

Rare disease



GRCh37 (hg19)
GRCh38 (hg38)

Cancer



Germline GRCh38 (hg38)
Somatic 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

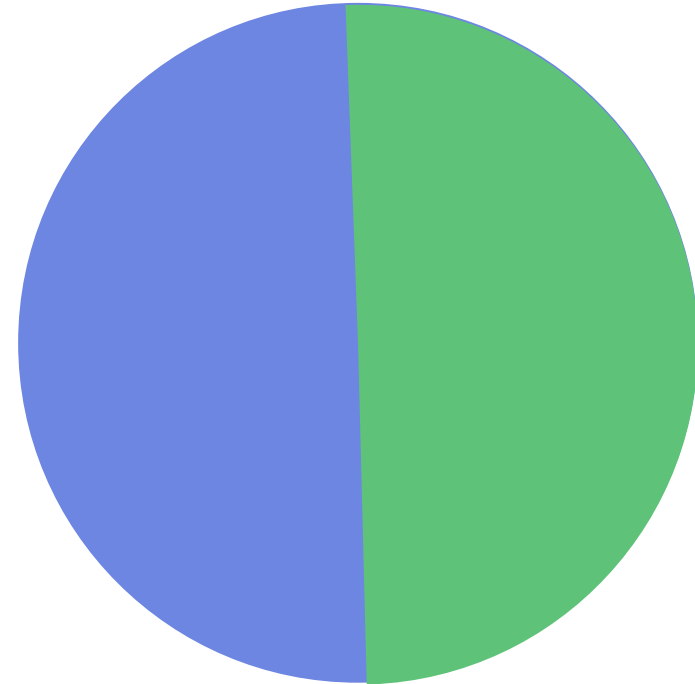
NHS GMS tiering

Rare disease



GRCh38 (hg38)

Cancer



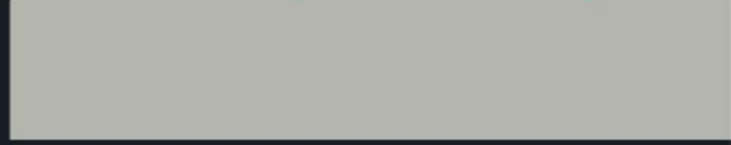
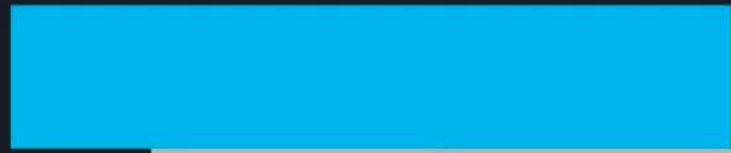
Germline GRCh38 (hg38)

Somatic GRCh38 (hg38)

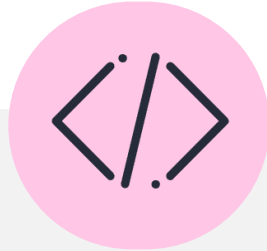
Variant data in LabKey demo



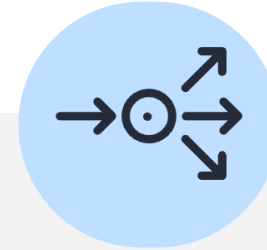
- Computer
- eperry's Home
- Trash
- Text Editor
- Rocket Chat
- Gvim
- Document Viewer
- IVA 2.0
- Airlock
- Data Discovery
- Open Targets
- R
- Participant Explorer
- LibreOffice
- Labkey
- IGV Browser
- Research Environment Documentation
- Research Registry
- Terminal Emulator
- Visual Studio Code
- Welcome Pack
- Panel App
- Git GUI
- RStudio
- Emacs



LabKey API



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



Combine queries between tables



Replicate queries between releases and
analyses



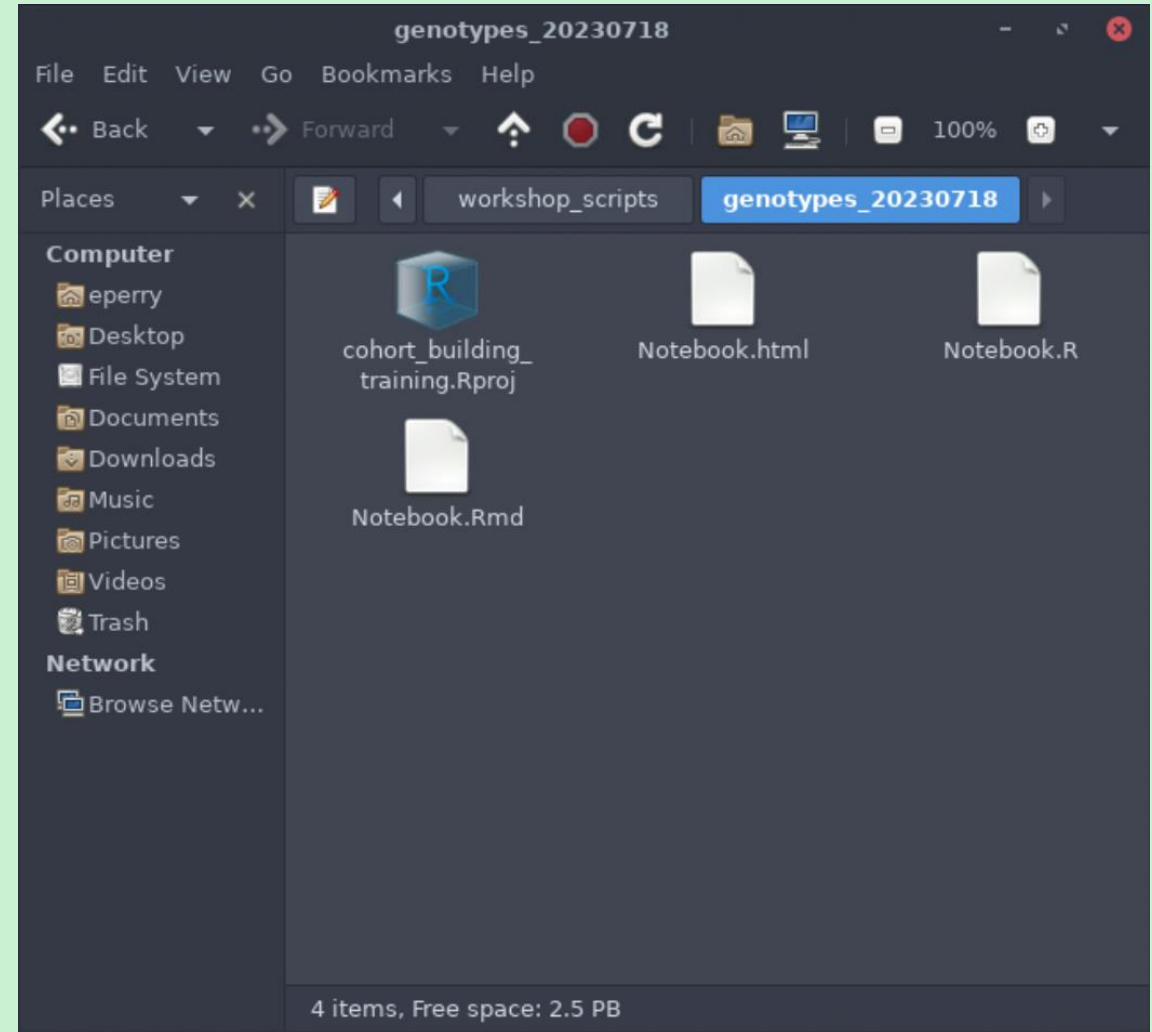
Work locally and on the HPC

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/example_scripts/workshop_scripts/genotypes_2026`



Variant data in LabKey API demo

genotypes_tr... (3) - JupyterLab — Mozilla Firefox

Jupyter Lab on Helix - Ge... R Notebook x genotypes_tr... (3) - Jupy x +

127.0.0.1:9537/lab/tree/genotypes_training.ipynb

File Edit View Run Kernel Tabs Settings Help

GENOTYPES_TRAINING.IPYNB

rd_cohort_building_trainin x genotypes_training.ipynb x cohort_building_training.ip...

Python 3 (ipykernel)

Finding participants by genotypes in Python

Contents:

- Use the LabKey API
 - Import Python modules you need
 - Helper function to access the LabKey API with Python
 - Querying the rare disease tiering_data table
 - Querying the exomiser table
 - Querying the cancer_tier_and_domain_varian table
 - Querying the NHS GMS tiering_data table
- Running workflows on the HPC
 - Small variant workflow
 - SV/CNV workflow
- Using bcftools on the HPC
- Optional exercise
 - A possible solution

Finding participants by genotypes in Python

This notebook will walk you through finding participants by genotypes. You are welcome to copy/paste any code from this notebook for your own scripts/notebooks.

Contents:

- Use the LabKey API
 - Import Python modules you need
 - Helper function to access the LabKey API with Python
 - Querying the tiering_data table
 - Querying the exomiser table
 - Querying the cancer_tier_and_domain_variants table
 - Querying the NHS GMS tiering_data table
- Running workflows on the HPC
- Using bcftools on the HPC
- Optional exercise

Use the LabKey API

Import Python modules you need

```
[32]: import numpy as np
import functools
import labkey
import pandas as pd
```

Helper function to access the LabKey API with

Simple 0 3 Python 3 (ipykernel) | Idle Mode: Command Ln 1, Col 1 genotypes_training.ipynb

R Notebook — Mozilla Firefox

R Notebook x +

file:///nas/weka.gel.zone/pgen_int_work/BRS/emily/genotypes_2023/genoty...

R Notebook

Finding participants by genotypes in R

This notebook will walk you through finding participants by genotypes. You are welcome to copy/paste any code from this notebook for your own scripts/notebooks.

Contents:

- Use the LabKey API
 - Import R libraries you need
 - Helper function to access the LabKey API with R
 - Querying the tiering_data table
 - Querying the exomiser table
 - Querying the cancer_tier_and_domain_variants table
 - Querying the NHS GMS tiering_data table
- Running workflows on the HPC
- Using bcftools on the HPC
- Optional exercise

Use the LabKey API

Import R libraries you need

```
library(tidyverse)
```

```
## — Attaching core tidyverse packages — tidyverse 2.0.0 —
## ✓ dplyr 1.1.2 ✓ readr 2.1.4
## ✓ forcats 1.0.0 ✓ stringr 1.5.0
## ✓ ggplot2 3.4.2 ✓ tibble 3.2.1
## ✓ lubridate 1.9.2 ✓ tidyr 1.3.0
## ✓ purrr 1.0.1
## — Conflicts — tidyverse_conflicts() —
## * dplyr::filter() masks stats::filter()
## * dplyr::lag() masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors
```

```
library(Rlabkey)
```

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

Rare disease



GRCh37 (hg19)
GRCh38 (hg38)

Cancer

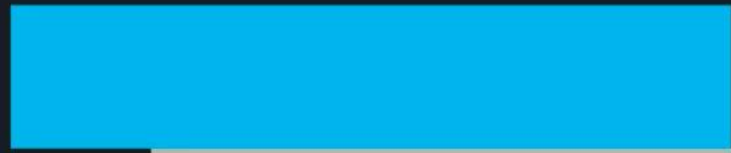


Germline GRCh38 (hg38)
Somatic GRCh38 (hg38)

IVA demo



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



CloudOS Cohort Browser – genotypes

Filter genotypes and phenotypes in one interface

Point-and-click interface to explore variants

Filter by loci, consequences, population frequencies and inheritance.

Find participant genotypes.

The screenshot displays the CloudOS Cohort Browser interface. A 'Global genomic filters' dialog box is open, showing various filter categories: Genes Location, Consequence, Population frequency, Clinical, Phenotype, Deleteriousness, and Conservation. The 'Genes Location' field contains '3,444-55555, 1,1-100000'. The 'Feature loci' field is empty. The 'Gene biotype' dropdown is set to 'Select option(s)'. Below these fields, a 'VARIANT' section lists options: SNV, INDEL, CNV, INSERTION, and DELETION. The main interface shows a cohort named 'Example_cohort_training' with 18 of 1186 samples. A table of variants is displayed under the 'ATF0-filter' filter. The table has columns for Location, Reference, Alternate, Chromosome, DISCOVER_ALL, GNOMAD_EXOMES_ALL, and GNOMAD_EXOMES_FIN. The table shows several variants on chromosome 1.

Location	Reference	Alternate	Chromosome	DISCOVER_ALL	GNOMAD_EXOMES_ALL	GNOMAD_EXOMES_FIN
1:10948349	G	C	1	0.0005	0.000531202	0.000449559
1:10948360	G	A	1	0.00154053	0.00262135	0.00107778
1:10948374	CTC		1			
1:10948392	G	A	1	0.0005	0.0000204369	
1:10948418	C	T	1		0.00000407967	
1:10948430	C	T	1	0.0005	0.0000977868	
1:10948431	C	A	1		0.0000163043	

CloudOS omics demo

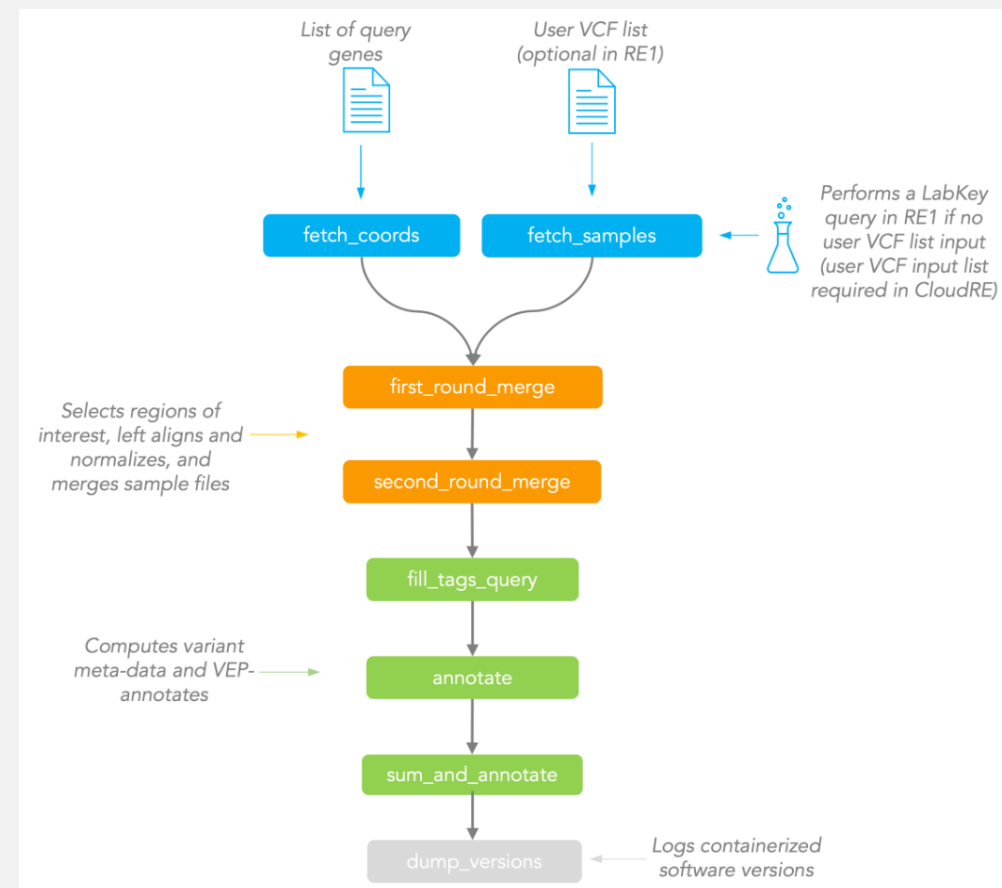
5. The Small Variant and Structural Variant workflows

Small variant workflow

Submit a list of genes

Find all short variants in these genes

Get 100k participants with these variants

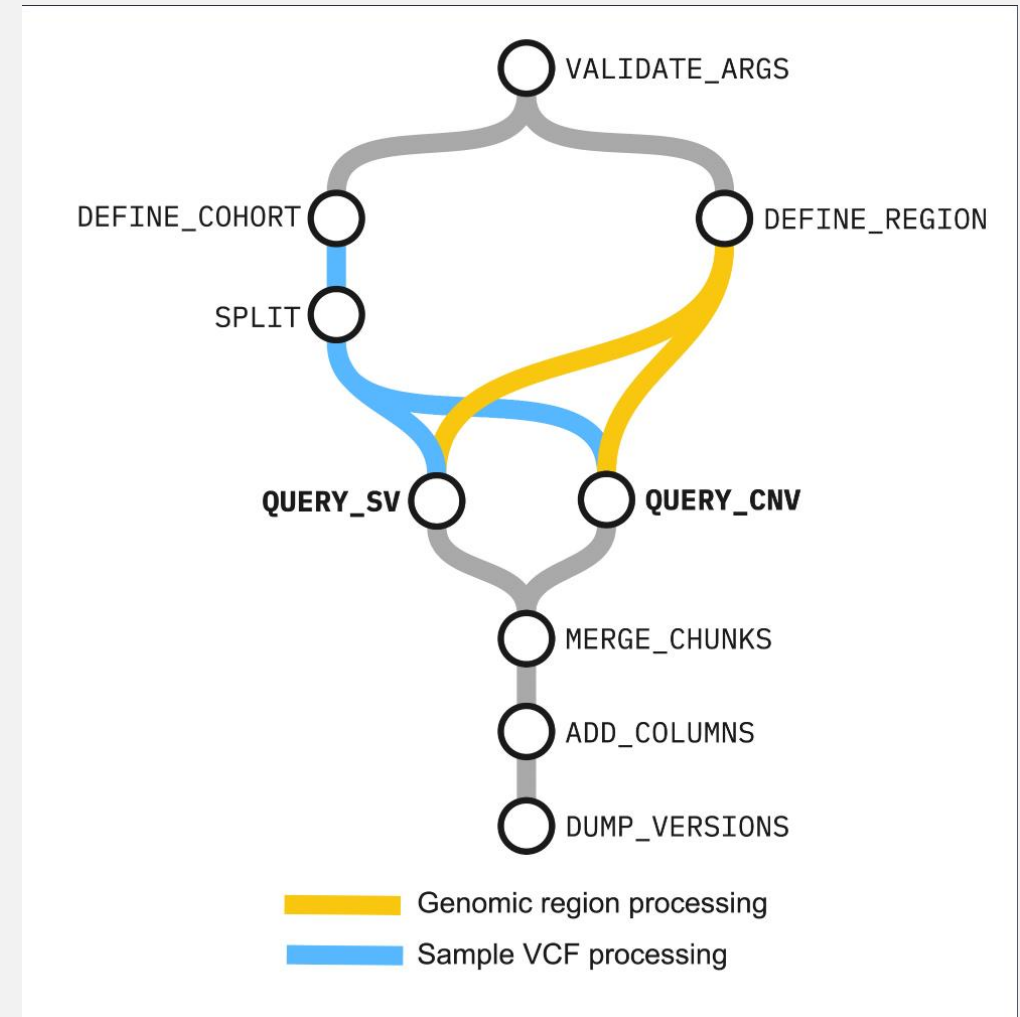


Structural variant workflow

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

genotypes_tr... - JupyterLab — Mozilla Firefox

127.0.0.1:8387/lab/tree/genotypes_training.ipynb

File Edit View Run Kernel Tabs Settings Help

Launcher x genotypes_training.ipynb x

604 rows x 7 columns

Running workflows on the HPC

Small variant workflow

- ssh to the HPC with your usual login credentials
- cd into your working directory
- Make and cd into your working directory `mkdir small_variant_demo cd small_variant_demo`
- Copy the Small Variant workflow submission script into your folder `cp /pge_int_data_resources/workflows/rdp_small_variant/main/submit.sh .`
- Make the file `gene_list.txt` and add your list of genes to it `vi gene_list.txt`
Add `SMC3` to the file `Esc :wq`
- Edit the submission script: `vi submit.sh`
Change the `project_code` to your code Change the `gene_input` to `gene_list.txt` `Esc :wq`
- Run the workflow `bsub < submit.sh`
- Find your results when the job is finished

```
[18]: small_variant_results = pd.read_csv('small_variant/results/GRCh38
small_variant_results

/resources/conda/miniconda3/envs/2022_base/lib/python3.7/site-pac
kages/IPython/core/interactiveshell.py:3186: DtypeWarning: Column
s (101) have mixed types.Specify dtype option on import or set lo
w_memory=False.
interactivity=interactivity, compiler=compiler, result=result)
```

Simple 0 s. 1 Python 3 (ipykernel) | Idle Mode: Command Ln 1 eperry@a-avtl60awn827:~

/nas/weka.gel.zone/pge_int_work/BRS/emily/genotypes_2024/Notebook.html

Notebook.html Open in Browser Find Publish

Running workflows on the HPC

Small variant workflow

- ssh to the HPC with your usual login credentials
- cd into your working directory
- Make and cd into your working directory `mkdir small_variant_demo cd small_variant_demo`
- Copy the Small Variant workflow submission script into your folder `cp /pge_int_data_resources/workflows/rdp_small_variant/main/submit.sh .`
- Make the file `gene_list.txt` and add your list of genes to it `vi gene_list.txt` Add `SMC3` to the file `Esc :wq`
- Edit the submission script: `vi submit.sh` Change the `project_code` to your code Change the `gene_input` to `gene_list.txt` `Esc :wq`
- Run the workflow `bsub < submit.sh`
- Find your results when the job is finished

```
small_variant_results <- readr::read_tsv('small_variant/results/GRCh38_SMC3_ENSG00000108055_annotated_va
riants.tsv')
```

```
## Rows: 58169 Columns: 102
## — Column specification —————
## Delimiter: "\t"
## chr (91): CHROM_variant, ID_variant, REF_variant, ALT_variant, Location_anno...
## dbf (11): POS_variant, AN_variant, AC_variant, AC_Hom_variant, AC_Het_varian...
##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

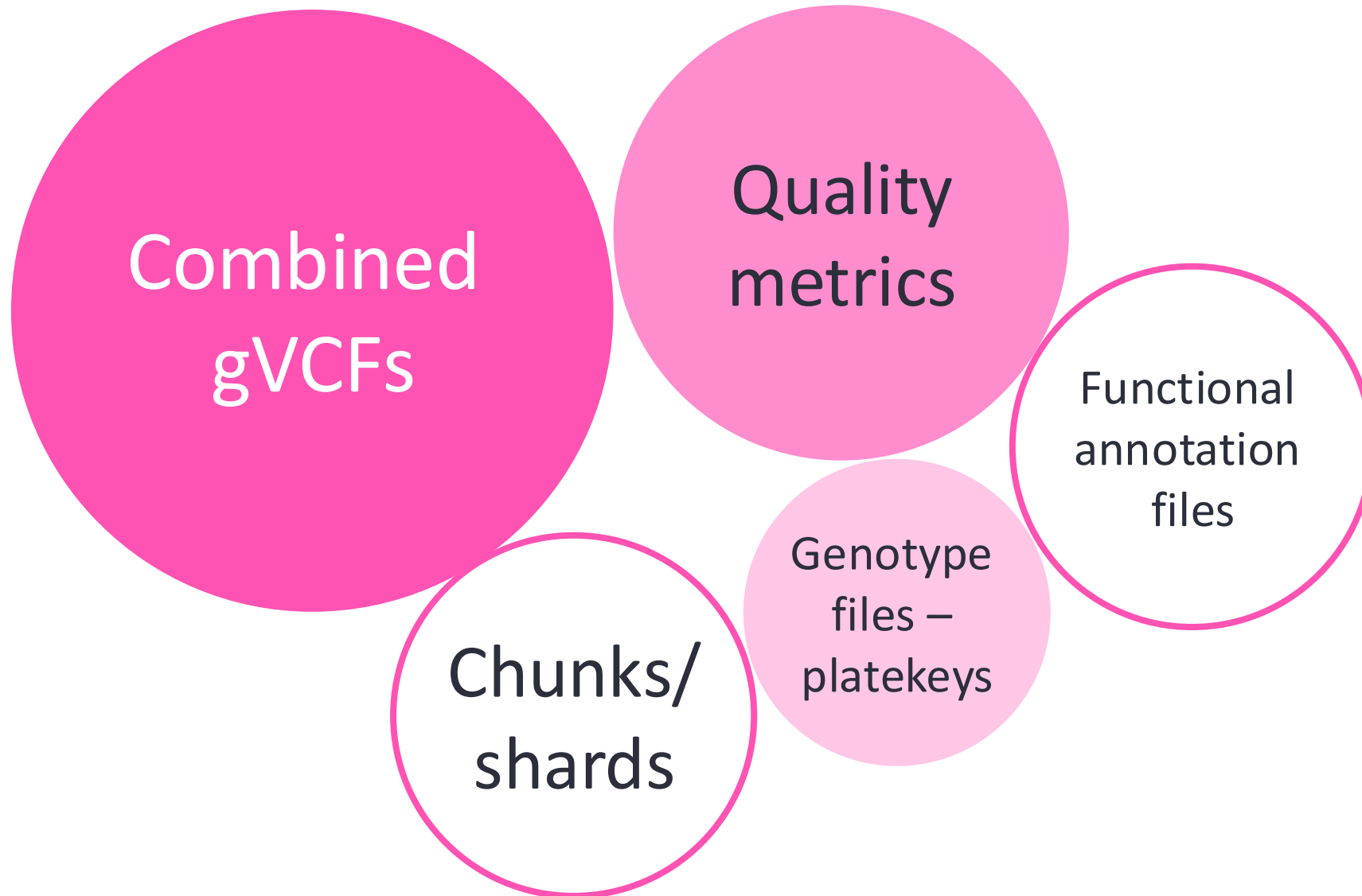
```
head(small_variant_results)
```

CHROM_variant	POS_variant	ID_variant	REF_variant	ALT_variant	AN_variant	AC_variant
<chr>	<dbl>	<chr>	<chr>	<chr>	<dbl>	<dbl>
chr10	110567687	chr10_110567687_1	G	A	242	121
chr10	110567687	chr10_110567687_1	G	A	242	121
chr10	110567687	chr10_110567687_1	G	A	242	121
chr10	110567687	chr10_110567687_1	G	A	242	121
chr10	110567687	chr10_110567687_1	G	A	242	121
chr10	110567687	chr10_110567687_1	G	A	242	121

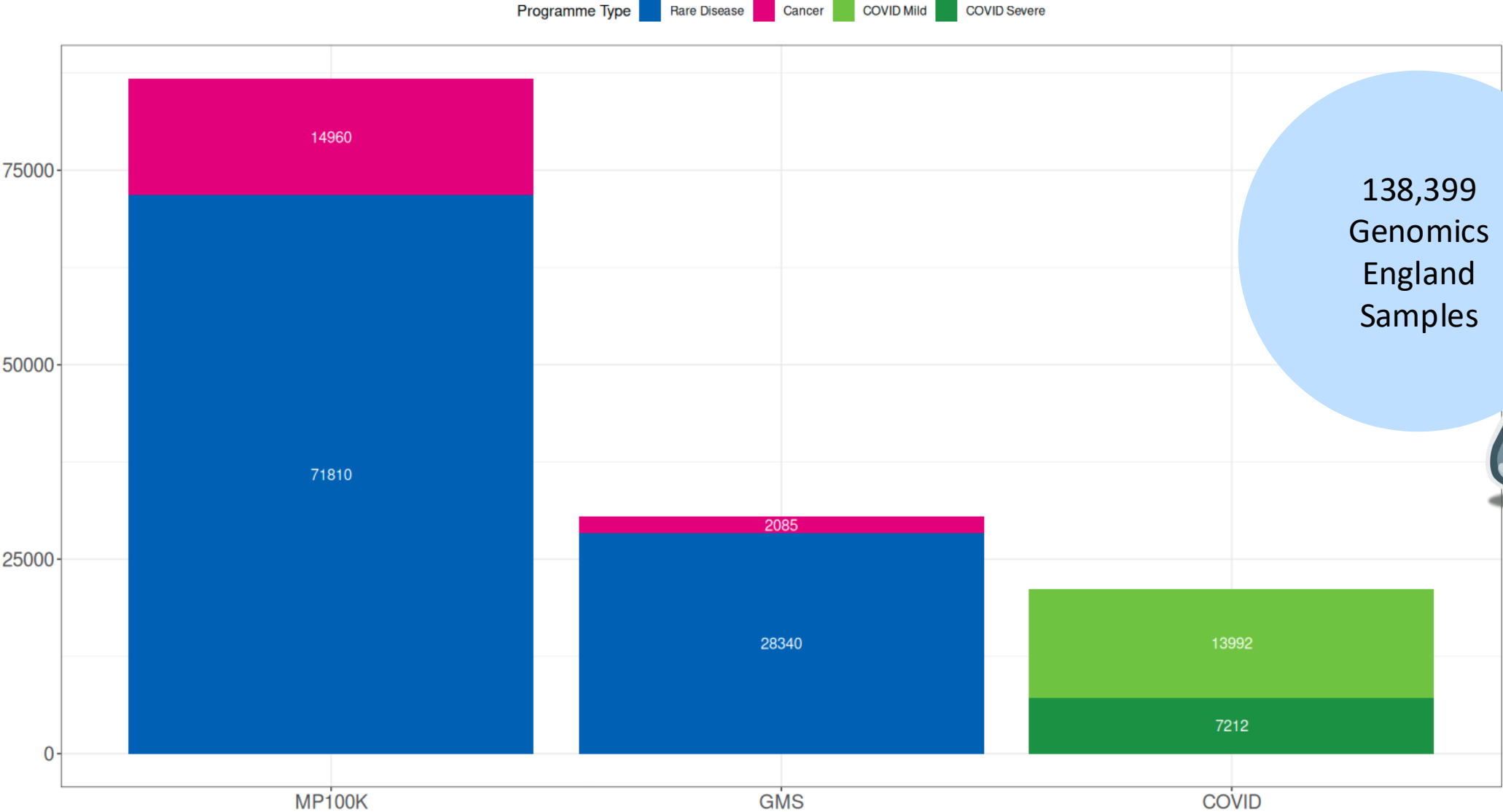
6 rows | 1-7 of 102 columns

6. Aggregate variant files

Aggregate VCFs



AggV3 – germline samples



138,399
Genomics
England
Samples

+7

AggV2 – germline samples

Rare disease



GRCh37 (hg19)
GRCh38 (hg38)

Cancer



Germline GRCh38 (hg38)
Somatic GRCh38 (hg38)

somAgg

Rare disease

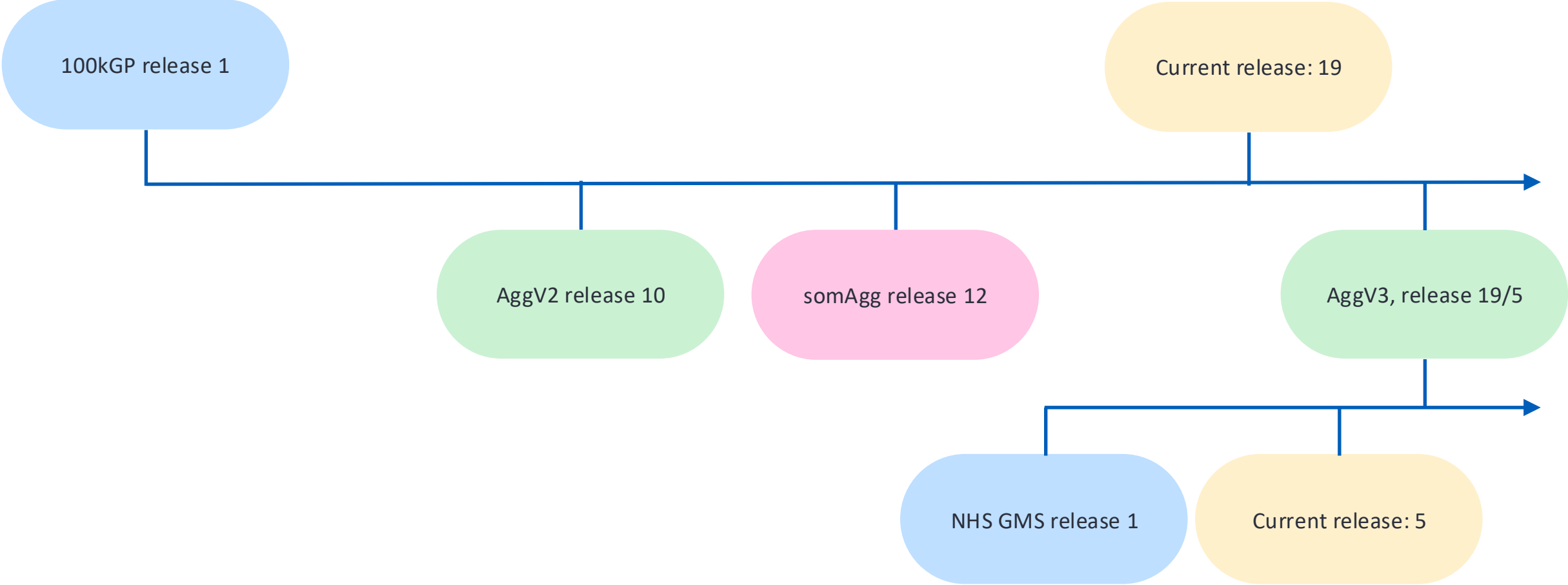
GRCh37 (hg19)
GRCh38 (hg38)

Cancer

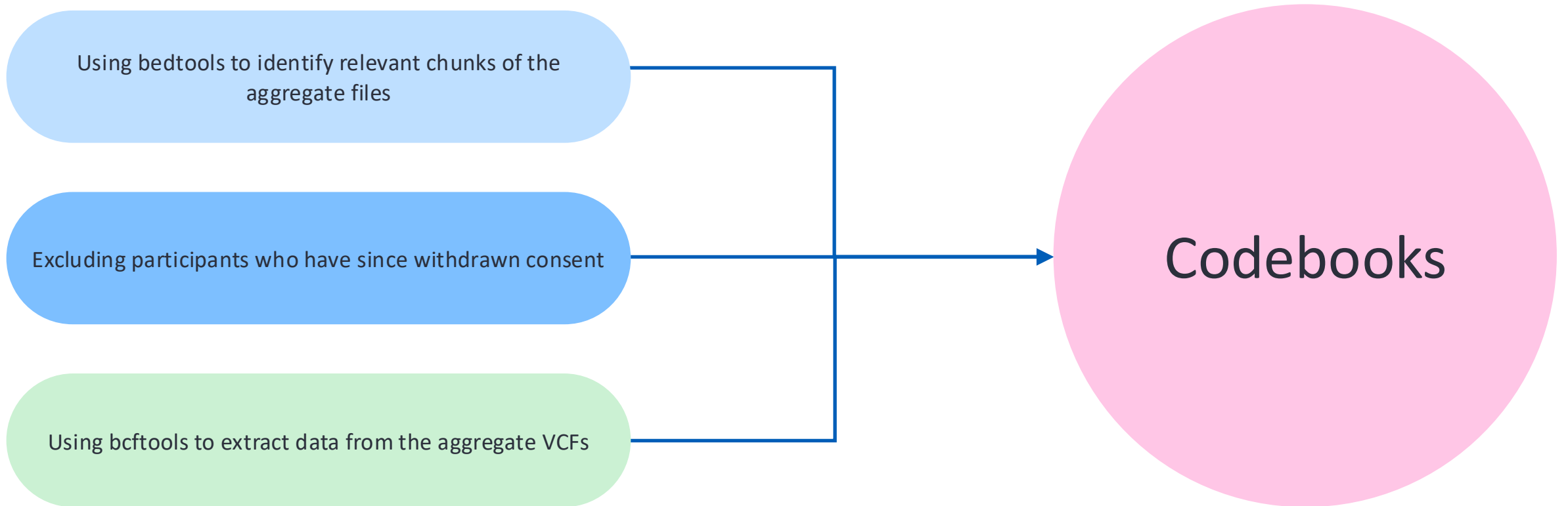


Germline GRCh38 (hg38)
Somatic GRCh38 (hg38)

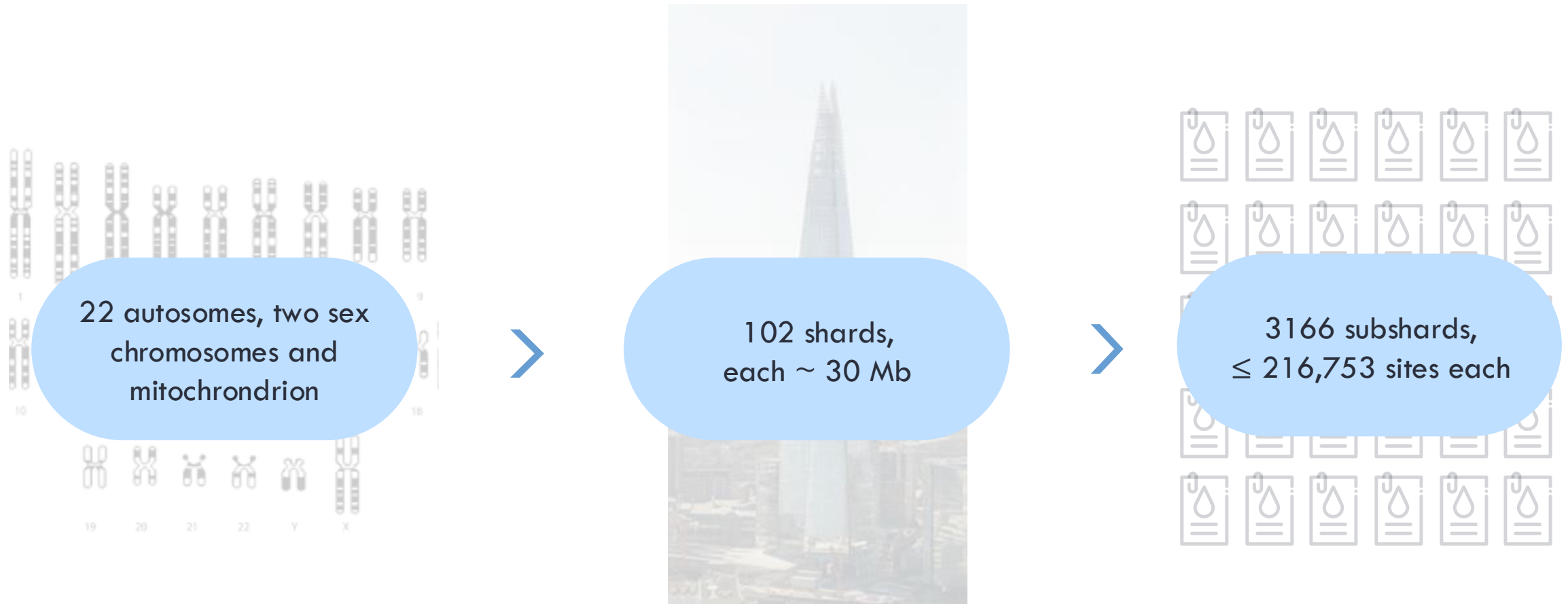
Aggregate VCFs



Aggregate VCFs



AggV3 shards



Same shards used across the genotype, functional annotation and site QC VCFs

AggV2/somAgg 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



AggV3 vs AggV2



More samples

138,399 compared to 78,195, including 100kGP, NHS GMS and COVID-19



Chunks vs shards

Consistent sizing of shards means it's easier to estimate compute needs of large tasks/pipelines



CloudOS

AggV3 only available on CloudOS with no plans to change. Time to get ready for the future Cloud-based RE



Iterative build

We can add more genomes to it as we receive them (>40k at a time to make it viable).



ABRAGEN 3.7.8

All genomes aligned and called using the same pipeline - better comparison within RE and to other similar projects



AggGIAB

Smaller aggregate for testing your workflows

Using bcftools demo

Aggregated Variant Calls (AggV3)

AggV3 is a set of multi-sample VCFs, bringing together short variants in germline genomes from [100kGP](#), [NHS GMS](#) and Covid-19 participants. AggV3 was prepared with by Illumina DRAGEN's Iterative GVCF Genotyper using genomes aligned using the DRAGEN 3.7.8 pipeline. Due to the size of the data, there are actually multiple VCFs, each representing a segment of the genome, known as "shards" and "subshards".

AggV3 contains information on participants who have since withdrawn consent from research. You cannot use them in any new analyses. It is extremely important to remove these samples from your analyses and only use samples included in the latest data release.

The latest updated list of samples for consented participants can be found in an S3 bucket within CloudOS (`s3://512426816668-ge1-data-resources/dragen3.7.8/AggV3_resources/samples/consented_individuals/2026-01-23/aggv3_consented_samples.txt`). When working within interactive sessions, you will need to mount this file to your session before you can use it. For batch analysis, you can provide the file as a parameter by clicking the button next to the `paramValue` textbox and navigating to the file within the File Explorer interface.

As AggV3 is a cross-programme dataset, you may need to update the list of consented individuals yourself at a later stage. For the 100,000 Genomes Project and [NHS-GMS](#) samples, please refer to the latest data release and filtering the `participant` table for `Consenting` in the `programme_consent_status` column. For the COVID19 participants, the list of samples can be used that are part of the latest available release.

To filter the aggregate to these samples, all `bctools` commands should include the flag `-S` `<path_to_consented_participants_list>`.

Submit a ticket to the [Genomics England Service desk](#) if you are unsure of how to filter the dataset for any other use.

Table of contents

- What data is in AggV3?
- How many variant sites are there in AggV3?
- How are the VCFs split up?
- How can I access AggV3?
- Where can I learn more?

- Aggregated Variant Calls (AggV3)
- Data generation, structure and locations
- DRAGEN 3.7.8
- Detailed methods as provided by Illumina
- AggV3 samples
- AggGIAB: A small aggregate with public data to test your workflows
- AggV3 functional annotation
- AggV3 sample quality metrics
- AggV3 site QC
- Processing of multiallelic VCFs
- AggV3 code book
 - AggV3 code book - identifying the correct subshard
 - AggV3 shard lookup tool
 - AggV3 code book - genotype queries
 - AggV3 code book -

Lifebit FedPaaS

pro.cloud-os.prod.a

RUN AND DEBUG: RUN

Open a file which can be debugged or run.

Run and Debug

To customize Run and Debug create a launch.json file.

PROBLEMS

(cloudos) ## Pack

envir

added

- b

The fol

pac

bed

The fol

bedto

Proceed

Downloa

Prepari

Verifying transaction: done

Executing transaction: done

(cloudos) `vscode@db4ccbcb6c9:~/session_data$`

Session data

Add data ▾

Save

Last saved on 02/02/2026 13:09:16

Data Items

No items to display

Welcome my_regions.bed × dragen.vcf.gz

my_regions.bed

1	chr1	230710048	230710048	rs699
---	------	-----------	-----------	-------

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL PORTS

bash + ▾ 🗑 ... | 🗑 ×

```
cloudos) vscode@db4ccb6c9:~/session_data$ bedtools intersect -wo -a my_regions.bed -b filesystems/genomic_data/biallelic_shards.bed
230710048 230710048 rs699 chr1 230054378 231051307 chr1:230054379-231051307 7 23 s3://357851407625-germline-aggre
~/3/data/euw2-dragen-igg-20250430075006-msvcf-version-1/data/shard-msvcf/shard-7/subshard-23/postproc/vcf/dragen.vcf.gz s3://357851407625-germline-aggregate-v3/
2-dragen-igg-20250430075006-msvcf-version-1/data/shard-msvcf/shard-7/subshard-23/postproc/vcf/dragen.vcf.gz.tbi 0
cloudos) vscode@db4ccb6c9:~/session_data$ ^C
cloudos) vscode@db4ccb6c9:~/session_data$
```



Data item successfully added. They should be available in the session shortly.

7. When/why you would use each method



Platform

Tiering and
exomiser
tables



IVA



Small/
Structural
variant
workflows



Querying
the
aggregates



Search by

Tiering and
exomiser
tables



IVA



Small/
Structural
variant
workflows



Querying
the
aggregates



Variants included

Tiering and
exomiser
tables

Variants that have passed filters

IVA

All

Small/
Structural
variant
workflows

All

Querying
the
aggregates

AggV3 – all variants
AggV2 – from GRCh38-aligned
genomes from release 12
somAgg – from GRCh38-aligned
genomes from release 8

Datasets

Tiering and
exomiser
tables

100k and NHS GMS

IVA and
Cohort
browser
omics

100k

Small/
Structural
variant
workflows

100k

Querying
the
aggregates

AggV2 and somAgg: 100k
AggV3: 100k, NHS GMS and Covid

Genome assembly

Tiering and
exomiser
tables

GRCh37 and GRCh38
Assembly as a separate column

Small/
Structural
variant
workflows

GRCh37 and GRCh38 queries simultaneously

IVA

GRCh37 and GRCh38 in separate databases

Querying
the
aggregates

AggV3: genomes realigned to GRCh38
AggV2 and somAgg: GRCh38 only

Underlying VCFs

Tiering and
exomiser
tables

Rare disease: Platypus
Cancer: Strelka

IVA

Rare disease: Platypus
Cancer: Strelka

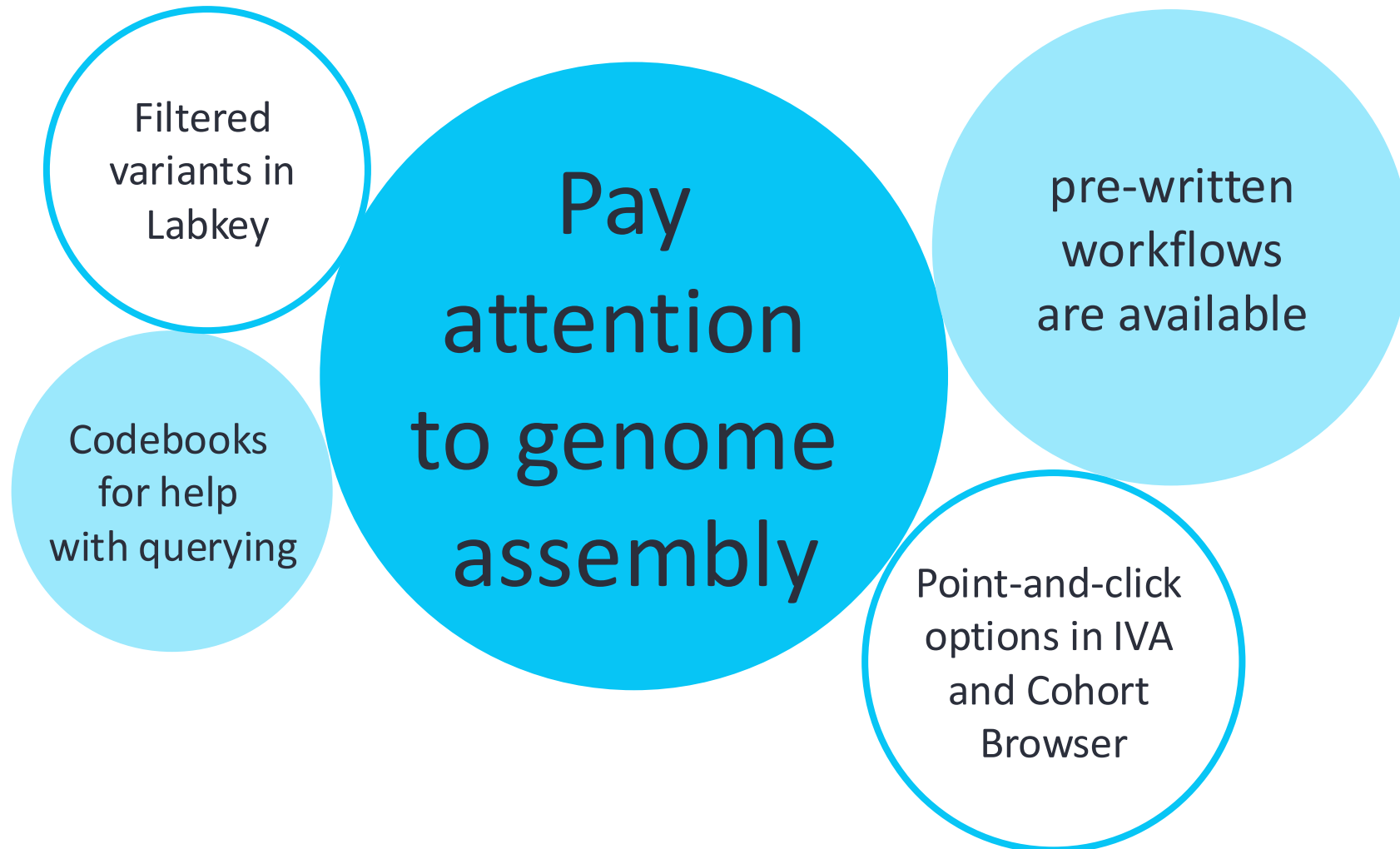
Small/
Structural
variant
workflows

Strelka

Querying
the
aggregates

AggV3: Dragen 3.7.8
AggV2 and somAgg: Strelka

Key takeaways



8. 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>

Questions



All your
microphones
are muted

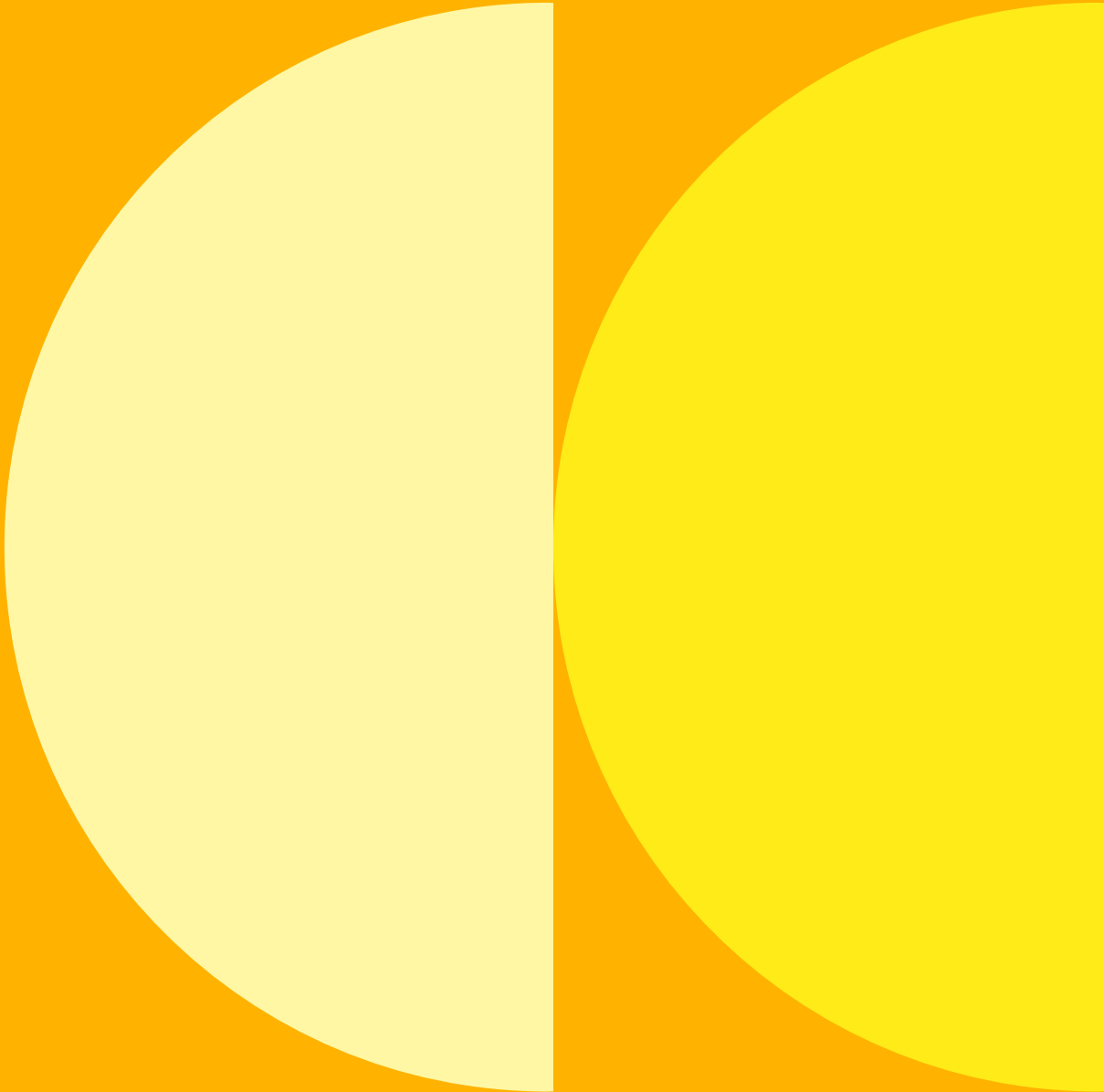


Use the Zoom
Q&A to ask
questions



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

Feedback



Thank you

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