

# *Analyzing Clinical and Genomic Oncological Data with {genieBPC} and {gnomeR}*

---

R\Medicine Demo

June 7, 2023

Sammi Brown

Jessica Lavery

Karissa Whiting



# Agenda



Projects GENIE & GENIE BPC



Clinico-Genomic Data Processing Pipeline



Case study



Clinical data processing with {genieBPC}



Genomic data processing with {gnomeR}

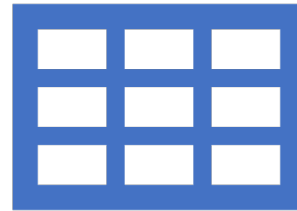


Conclusion

# Projects GENIE & GENIE BPC



The goal of **Project GENIE Biopharma Collaborative (BPC)** is to augment the existing registry genomic data from AACR Project GENIE with enhanced clinical (phenomic) data to support clinical-genomics analyses.



Phenomic data are curated using the PRISMM curation model to capture detailed information on cancer diagnosis, drug regimens, disease status from radiology reports, pathology reports and medical oncologist assessments, structured in several datasets with over 700 feature variables.



Analyses using linked clinico-genomic databases – including GENIE BPC – will help to drive advancements in precision oncology in identifying the genomic alterations and drug therapies that optimize clinical outcomes.

# Genomic data included in GENIE

Researchers receive genomic data in different formats and types

The AACR Project GENIE data repository is comprised of one type of genomic data called **tumor DNA sequencing assays**

- Collected from tumor samples via biopsy/resection
- Compare DNA sequence in cancer cells with that in normal cells

Sequencing assays can be broad or targeted

- Broad regions: whole genome/whole exome sequencing
- Targeted regions: gene panels
  - GENIE data consists of data from targeted gene panels from high-throughput (huge amounts of data) sequencing assays, also referred to as next-generation sequencing (NGS)

# GENIE BPC Data

- Data are publicly released by cancer cohort: non-small cell lung (NSCLC), colorectal (CRC), breast, pancreas, prostate, bladder
- New versions of data are released periodically to include additional patients and variables and to incorporate data corrections
- .csv and .txt data files are available for download from Sage Bionetworks' Synapse data sharing platform
- Downloading each file individually poses challenges for efficient and reproducible workflows

# {genieBPC} & {gnomeR} R Packages



**The {genieBPC} package is a pipeline to programmatically access the data corresponding to each release from Synapse to support reproducibility, and to create datasets linking clinical and genomic data for analysis.**



**Created and developed by**

Samantha Brown  
Michael Curry  
Hannah Fuchs  
Jessica Lavery  
Axel Martin  
Dan Sjoberg  
Karissa Whiting



**The {gnomeR} package provides a consistent framework for genetic data wrangling, processing, visualization and analysis.**



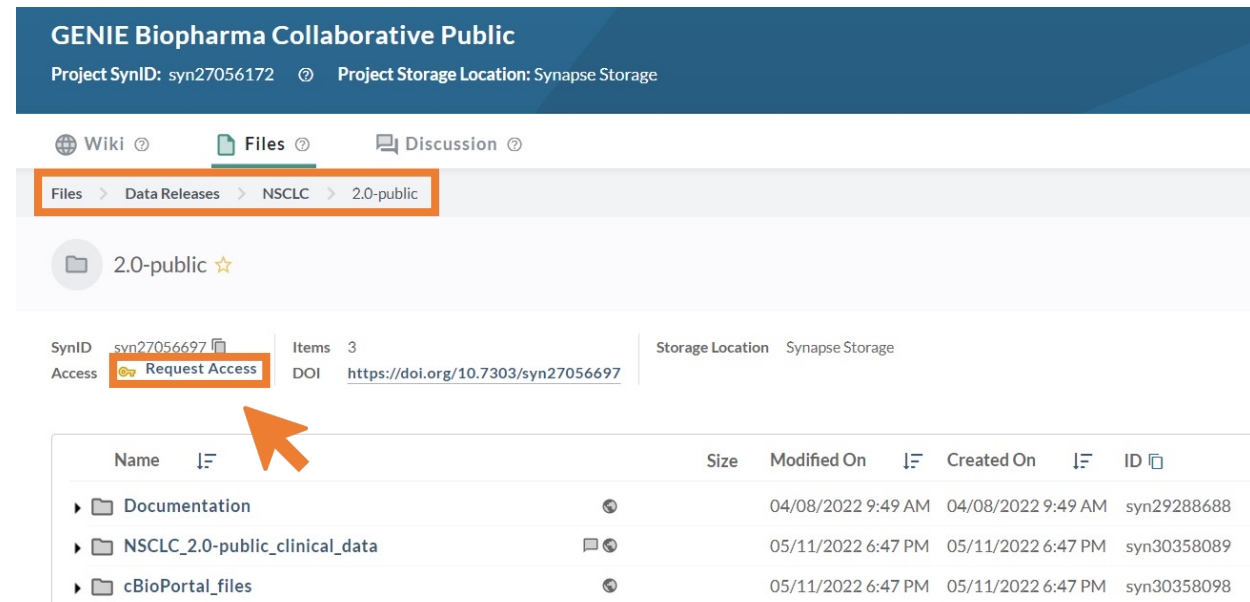
**Created and developed by**

Arshi Arora  
Michael Curry  
Hannah Fuchs  
Axel Martin  
Karissa Whiting

# Register for Synapse Account

## Instructions:

1. Register for a [‘Synapse’ account](https://www.synapse.org/#). Be sure to create a username and password. **Do NOT connect via your Google account.**
  - a) <https://www.synapse.org/#>
2. Accept the **Synapse account terms of use**.
3. Navigate to GENIE Biopharma Collaborative Public page
  - a) <https://www.synapse.org/#!Synapse:syn27056172/wiki/616601>
4. In the Files folder, navigate to Data Releases → NSCLC → 2.0-public
5. Select *Request Access*, review the **terms of data use** and click *Accept*



GENIE Biopharma Collaborative Public  
Project SynID: syn27056172 Project Storage Location: Synapse Storage

Wiki Files Discussion

Files > Data Releases > NSCLC > 2.0-public

2.0-public ☆

SynID syn27056697 Items 3  
Access Request Access DOI <https://doi.org/10.7303/syn27056697> Storage Location Synapse Storage

Name	Size	Modified On	Created On	ID
Documentation		04/08/2022 9:49 AM	04/08/2022 9:49 AM	syn29288688
NSCLC_2.0-public_clinical_data		05/11/2022 6:47 PM	05/11/2022 6:47 PM	syn30358089
cBioPortal_files		05/11/2022 6:47 PM	05/11/2022 6:47 PM	syn30358098

# Installation Instructions

Installing {genieBPC}:

```
install.packages("GENIE-BPC/genieBPC")
```

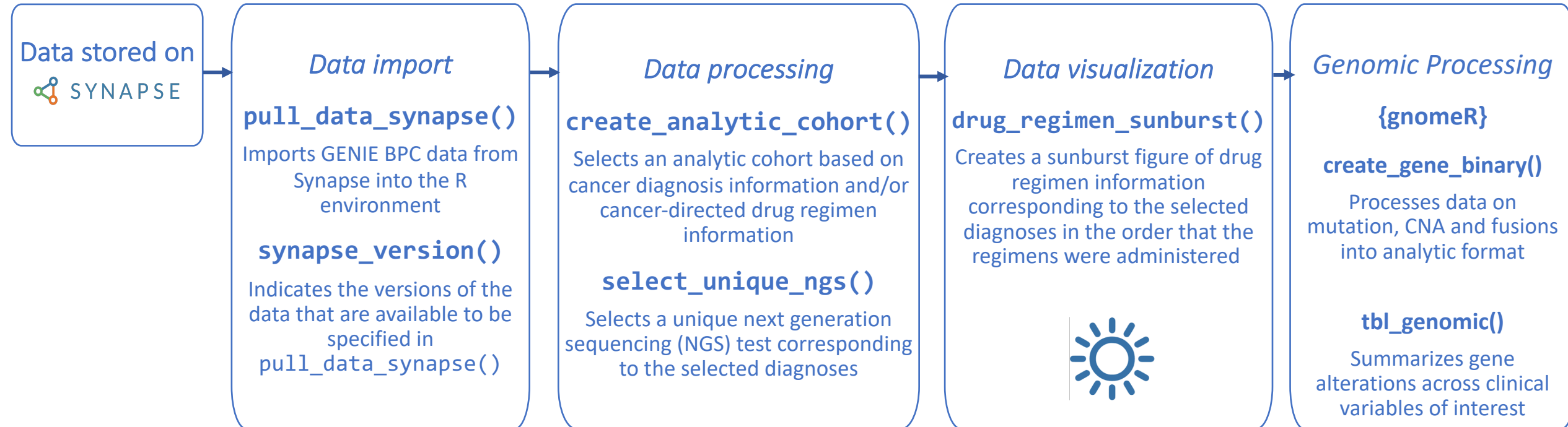
Installing {gnomeR}:

```
# install.packages("devtools")  
devtools::install_github("MSKCC-Epi-Bio/gnomeR")
```

- These instructions are also included in the Demo.R script on our GitHub repository: [https://github.com/GENIE-BPC/intro\\_to\\_genieBPC\\_and\\_gnomeR](https://github.com/GENIE-BPC/intro_to_genieBPC_and_gnomeR)
- Further R package details are available on the {genieBPC} [GitHub repo](#) & [website](#) and the {gnomeR} [GitHub repo](#) & [website](#)
- Note: Both R packages require R version  $\geq 3.6$



# Clinico-Genomic Data Processing Pipeline





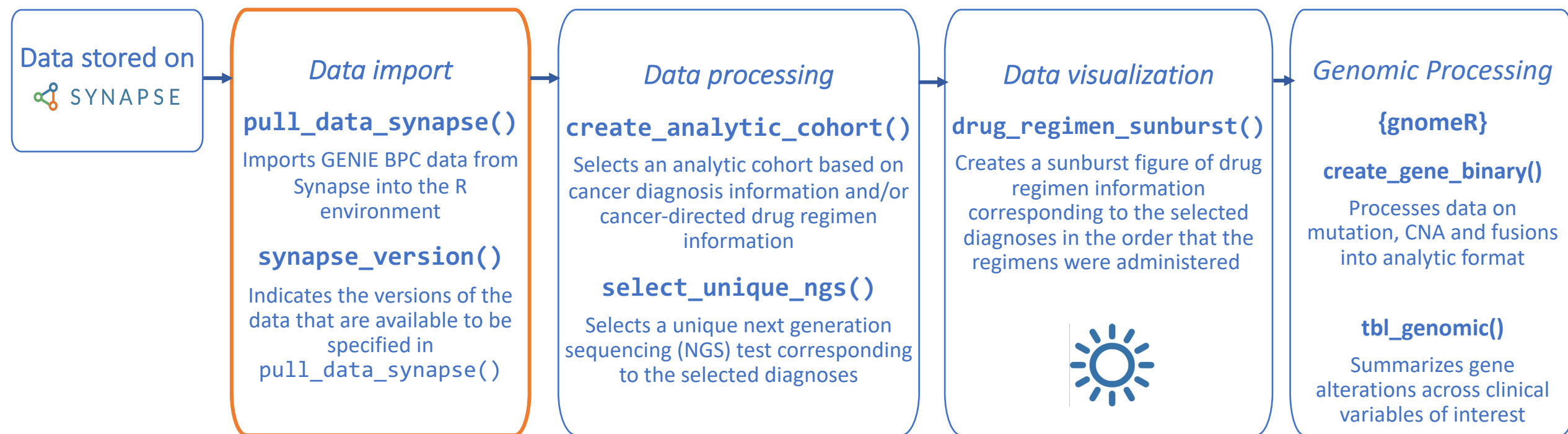
# Case Study

---

Create a cohort of patients who were diagnosed with Stage IV adenocarcinoma non-small cell lung cancer (NSCLC) and received Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab as their first cancer-directed drug regimen after diagnosis.

Follow along using the Demo.R script on our GitHub repository: [https://github.com/GENIE-BPC/intro\\_to\\_genieBPC\\_and\\_gnomeR](https://github.com/GENIE-BPC/intro_to_genieBPC_and_gnomeR)

# Clinico-Genomic Data Processing Pipeline



# Set Synapse Credentials

To pull data from Synapse, users must create a Synapse account and store their Synapse credentials in the R environment. The `set_synapse_credentials()` function will store credentials during each R session:

```
set_synapse_credentials(username = 'your_username',  
                        password = 'your_password')
```

# `synapse_version()`

- Helper function that returns a table of GENIE BPC data releases that are currently available
- `synapse_version()` has one input: `most_recent = TRUE/FALSE`
  - Calling `genieBPC::synapse_version(most_recent = TRUE)` will return a table with each cancer cohort and its latest data release version
  - Calling `genieBPC::synapse_version(most_recent = FALSE)` will return a table with all cancer cohorts and data releases available

# synapse\_version()

- Helper function that returns a table of GENIE BPC data releases that are currently available
- `synapse_version()` has one input: `most_recent = TRUE/FALSE`
  - Calling `genieBPC::synapse_version(most_recent = TRUE)` will return a table with each cancer cohort and its latest data release version
  - Calling `genieBPC::synapse_version(most_recent = FALSE)` will return a table with all cancer cohorts and data releases available

## `synapse_version(most_recent = TRUE)`

cohort	version	release_date	all_versions
BLADDER	v1.1-consortium	November 2022	Most Recent Versions
BrCa	v1.2-consortium	October 2022	Most Recent Versions
CRC	v1.2-consortium	August 2021	Most Recent Versions
CRC	v2.0-public	October 2022	Most Recent Versions
NSCLC	v2.1-consortium	August 2021	Most Recent Versions
NSCLC	v2.0-public	May 2022	Most Recent Versions
PANC	v1.2-consortium	January 2023	Most Recent Versions
Prostate	v1.2-consortium	January 2023	Most Recent Versions

# pull\_data\_synapse()

- Pull GENIE BPC clinical and genomic data directly from Synapse into R
- Can specify cancer type (``cohort``) and version of data (``version``)
  - Version of the data is updated periodically on Synapse with re-releases (new variables available, additional QA, etc.)
- Returns a nested list of data frames for each cancer site for the accompanying version

Argument	Description	Acceptable Values
cohort	<ul style="list-style-type: none"><li>• GENIE BPC Project cancer</li><li>• Currently, NSCLC and CRC are the only two publicly available datasets</li></ul>	<ul style="list-style-type: none"><li>• NSCLC</li><li>• CRC</li><li>• BrCa</li><li>• PANC</li><li>• Prostate</li><li>• BLADDER</li></ul>
version	Version of the data (e.g v1.1-consortium, v2.0-public)	<ul style="list-style-type: none"><li>• Values can be found in <code>synapse_version()</code></li></ul>

Demo: Run **pull\_data\_synapse()** for case study





# Demo: Run **pull\_data\_synapse()** for case study

```
library(genieBPC)
```



# Demo: Run **pull\_data\_synapse()** for case study

```
library(genieBPC)  
set_synapse_credentials()
```



# Demo: Run **pull\_data\_synapse()** for case study

```
library(genieBPC)
set_synapse_credentials()
nscclc_synapse_data <- pull_data_synapse(cohort = "NSCLC", version = "v2.0-public")
```



# Demo: Run **pull\_data\_synapse()** for case study

```
library(genieBPC)
```

```
set_synapse_credentials()
```

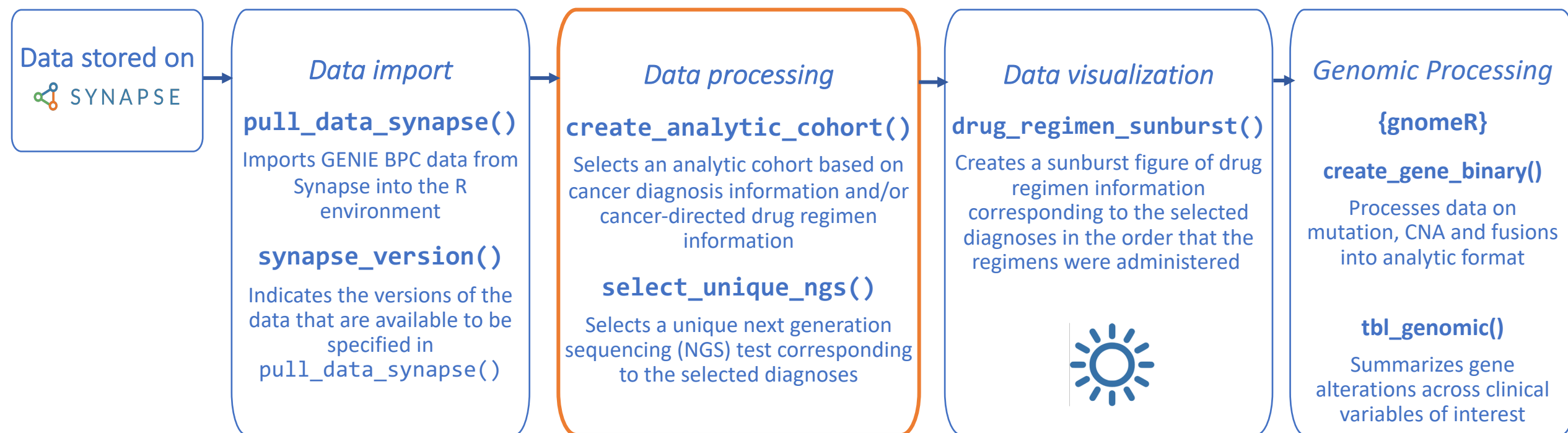
```
nsc1c_synapse_data <- pull_data_synapse(cohort = "NSCLC", version = "v2.0-public")
```

Calling **nsc1c\_synapse\_data\$NSCLC\_v2.0** returns a list of datasets in **nsc1c\_synapse\_data**:

- pt\_char
- ca\_dx\_index
- ca\_dx\_non\_index
- ca\_drugs
- prissmm\_pathology
- prissmm\_imaging
- prissmm\_md
- cpt
- mutations\_extended
- cna
- fusions



# Clinico-Genomic Data Processing Pipeline



# create\_analytic\_cohort()



Create a cohort from the GENIE BPC data

Cancer diagnosis information such as cancer cohort, treating institution, histology, and stage at diagnosis  
Cancer-directed regimen information including regimen name and regimen order.



This function returns all clinical and genomic data for the selected patients

# create\_analytic\_cohort()

Argument	Description	Acceptable Values
data_synapse	List returned from pull_data_synapse()	<ul style="list-style-type: none"><li>Name of object in global environment that was returned from pull_data_synapse()</li></ul>

# create\_analytic\_cohort()

Argument	Description	Acceptable Values
data_synapse	List returned from pull_data_synapse()	<ul style="list-style-type: none"><li>Name of object in global environment that was returned from pull_data_synapse()</li></ul>
index_ca_seq	Index cancer sequence. Default is 1, indicating the patient's first index cancer. This is the cancer that met the eligibility criteria for the project and was selected at random for PRISSMM phenomic data curation.	<ul style="list-style-type: none"><li>Numeric (1+)</li></ul>



# create\_analytic\_cohort()

Argument	Description	Acceptable Values
data_synapse	List returned from pull_data_synapse()	<ul style="list-style-type: none"><li>• Name of object in global environment that was returned from pull_data_synapse()</li></ul>
index_ca_seq	Index cancer sequence. Default is 1, indicating the patient's first index cancer. This is the cancer that met the eligibility criteria for the project and was selected at random for PRISSMM phenomic data curation.	<ul style="list-style-type: none"><li>• Numeric (1+)</li></ul>
institution	GENIE BPC participating institution. Default selection is all institutions. <i>Note that not all institutions curated data for all cancer sites.</i>	<ul style="list-style-type: none"><li>• DFCI</li><li>• MSK</li><li>• UHN</li><li>• VICC</li></ul>

# create\_analytic\_cohort()

Argument	Description	Acceptable Values
stage_dx	Stage at diagnosis. Default selection is all stages.	<ul style="list-style-type: none"><li>• Stage I</li><li>• Stage II</li><li>• Stage III</li><li>• Stage I-III NOS</li><li>• Stage IV</li></ul>

# create\_analytic\_cohort()

Argument	Description	Acceptable Values
stage_dx	Stage at diagnosis. Default selection is all stages.	<ul style="list-style-type: none"><li>• Stage I</li><li>• Stage II</li><li>• Stage III</li><li>• Stage I-III NOS</li><li>• Stage IV</li></ul>
histology	Cancer histology. Default selection is all histologies.	<ul style="list-style-type: none"><li>• Adenocarcinoma</li><li>• Squamous cell</li><li>• Sarcoma</li><li>• Small cell carcinoma</li><li>• Other histologies/mixed tumor</li></ul>

# create\_analytic\_cohort()

Argument	Description	Acceptable Values
regimen_drugs	Vector with names of drugs in cancer-directed regimen, separated by a comma. For example, to specify a regimen consisting of Carboplatin and Pemetrexed Disodium, specify <code>regimen_drugs = "Carboplatin, Pemetrexed Disodium"</code> .	Acceptable values are found in the <code>drug_names_by_cohort</code> dataset provided with this package.
regimen_type	Indicates whether the regimen(s) specified in <code>regimen_drugs</code> indicates the exact regimen to return, or if regimens containing the drugs listed in <code>regimen_drugs</code> should be returned.	<ul style="list-style-type: none"><li>• Exact</li><li>• Containing</li></ul>

# Example: regimen\_drugs and regimen\_type

regimen_drugs	regimen_type	Example regimens returned
Carboplatin	Exact	<ul style="list-style-type: none"><li>• Carboplatin</li></ul>
Carboplatin	Containing	<ul style="list-style-type: none"><li>• Carboplatin</li><li>• Carboplatin, Cisplatin</li><li>• Carboplatin, Paclitaxel</li><li>• Carboplatin, Pemetrexed Disodium</li><li>• etc.</li></ul>

# create\_analytic\_cohort()

Argument	Description	Acceptable Values
regimen_order	Order of cancer-directed regimen. If multiple drugs are specified, regimen_order indicates the regimen order for all drugs; different values of regimen_order cannot be specified for different drug regimens.	<ul style="list-style-type: none"><li>Numeric (1+)</li></ul>
regimen_order_type	Specifies whether the 'regimen_order' parameter refers to the order of receipt of the drug regimen within the cancer diagnosis (across all other drug regimens; "within cancer") or the order of receipt of the drug regimen within the times that that drug regimen was administered (e.g. the first time carboplatin pemetrexed was received, out of all times that the patient received carboplatin pemetrexed; "within regimen").	<ul style="list-style-type: none"><li>Within cancer</li><li>Within regimen</li></ul>

# create\_analytic\_cohort()

Argument	Description	Acceptable Values
return_summary	Specifies whether summary tables are returned using {gtsummary}. Default is FALSE.	<ul style="list-style-type: none"><li>• TRUE</li><li>• FALSE</li></ul>

# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

Case Study: Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis





# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

```
nslc_cohort <- create_analytic_cohort(
```

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis



# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

```
nsclc_cohort <- create_analytic_cohort(  
  data_synapse = nsclc_synapse_data$NSCLC_v2.0,
```

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis



# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis

```
nsccl_cohort <- create_analytic_cohort(  
  data_synapse = nsccl_synapse_data$NSCLC_v2.0,  
  stage_dx = c("Stage IV"),
```



# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis

```
nsccl_cohort <- create_analytic_cohort(  
  data_synapse = nsccl_synapse_data$NSCLC_v2.0,  
  stage_dx = c("Stage IV"),  
  histology = "Adenocarcinoma",
```



# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis

```
nsclc_cohort <- create_analytic_cohort(  
  data_synapse = nsclc_synapse_data$NSCLC_v2.0,  
  stage_dx = c("Stage IV"),  
  histology = "Adenocarcinoma",  
  regimen_drugs = c("Carboplatin, Pemetrexed Disodium",
```



# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis

```
nsclc_cohort <- create_analytic_cohort(  
  data_synapse = nsclc_synapse_data$NSCLC_v2.0,  
  stage_dx = c("Stage IV"),  
  histology = "Adenocarcinoma",  
  regimen_drugs = c("Carboplatin, Pemetrexed Disodium",  
                    "Cisplatin, Pemetrexed Disodium",
```



# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis

```
nsclc_cohort <- create_analytic_cohort(  
  data_synapse = nsclc_synapse_data$NSCLC_v2.0,  
  stage_dx = c("Stage IV"),  
  histology = "Adenocarcinoma",  
  regimen_drugs = c("Carboplatin, Pemetrexed Disodium",  
                    "Cisplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Carboplatin, Pemetrexed Disodium",
```



# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis

```
nsclc_cohort <- create_analytic_cohort(  
  data_synapse = nsclc_synapse_data$NSCLC_v2.0,  
  stage_dx = c("Stage IV"),  
  histology = "Adenocarcinoma",  
  regimen_drugs = c("Carboplatin, Pemetrexed Disodium",  
                    "Cisplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Carboplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Cisplatin, Pemetrexed Disodium"),
```





# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis

```
nsclc_cohort <- create_analytic_cohort(  
  data_synapse = nsclc_synapse_data$NSCLC_v2.0,  
  stage_dx = c("Stage IV"),  
  histology = "Adenocarcinoma",  
  regimen_drugs = c("Carboplatin, Pemetrexed Disodium",  
                    "Cisplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Carboplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Cisplatin, Pemetrexed Disodium"),  
  regimen_type = "Exact",
```



# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis

```
nsccl_cohort <- create_analytic_cohort(  
  data_synapse = nsccl_synapse_data$NSCLC_v2.0,  
  stage_dx = c("Stage IV"),  
  histology = "Adenocarcinoma",  
  regimen_drugs = c("Carboplatin, Pemetrexed Disodium",  
                    "Cisplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Carboplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Cisplatin, Pemetrexed Disodium"),  
  regimen_type = "Exact",  
  regimen_order = 1,
```



# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis

```
nsccl_cohort <- create_analytic_cohort(  
  data_synapse = nsccl_synapse_data$NSCLC_v2.0,  
  stage_dx = c("Stage IV"),  
  histology = "Adenocarcinoma",  
  regimen_drugs = c("Carboplatin, Pemetrexed Disodium",  
                    "Cisplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Carboplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Cisplatin, Pemetrexed Disodium"),  
  regimen_type = "Exact",  
  regimen_order = 1,  
  regimen_order_type = "within cancer",
```



# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis

```
nsclc_cohort <- create_analytic_cohort(  
  data_synapse = nsclc_synapse_data$NSCLC_v2.0,  
  stage_dx = c("Stage IV"),  
  histology = "Adenocarcinoma",  
  regimen_drugs = c("Carboplatin, Pemetrexed Disodium",  
                    "Cisplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Carboplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Cisplatin, Pemetrexed Disodium"),  
  regimen_type = "Exact",  
  regimen_order = 1,  
  regimen_order_type = "within cancer",  
  return_summary = TRUE
```



# Demo: **create\_analytic\_cohort()** for case study using NSCLC 2.0-public data

**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis

```
nsclc_cohort <- create_analytic_cohort(  
  data_synapse = nsclc_synapse_data$NSCLC_v2.0,  
  stage_dx = c("Stage IV"),  
  histology = "Adenocarcinoma",  
  regimen_drugs = c("Carboplatin, Pemetrexed Disodium",  
                    "Cisplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Carboplatin, Pemetrexed Disodium",  
                    "Bevacizumab, Cisplatin, Pemetrexed Disodium"),  
  regimen_type = "Exact",  
  regimen_order = 1,  
  regimen_order_type = "within cancer",  
  return_summary = TRUE  
)
```



nsccl\_cohort  
\$tbl\_overall\_  
summary

Characteristic	N = 241 patients <sup>†</sup>
Number of diagnoses per patient in cohort_ca_dx data frame	
1	241 (100%)
Number of regimens per patient in cohort_ca_drugs data frame	
1	241 (100%)
Number of CPTs per patient in cohort_ngs data frame	
1	222 (92%)
2	18 (7.5%)
4	1 (0.4%)
<sup>†</sup> n (%)	

nsccl\_cohort  
\$tbl\_cohort

Characteristic	N = 241 Diagnoses <sup>1</sup>
<b>Cohort (cohort)</b>	
NSCLC	241 (100%)
<b>Institution (institution)</b>	
DFCI	92 (38%)
MSK	118 (49%)
VICC	31 (13%)
<b>Stage at diagnosis (stage_dx)</b>	
Stage IV	241 (100%)
<b>Histology (ca_hist_adeno_squamous)</b>	
Adenocarcinoma	241 (100%)
<sup>1</sup> n (%)	

nsccl\_cohort  
\$tbl\_drugs

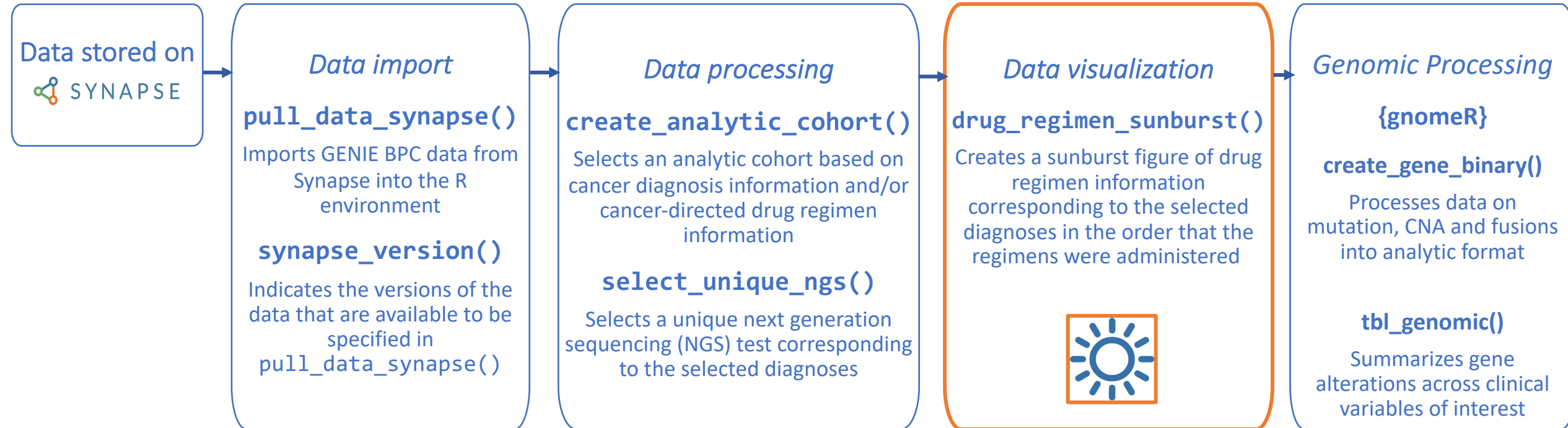
Characteristic	N = 241 Regimens <sup>1</sup>
<b>Cohort (cohort)</b>	
NSCLC	241 (100%)
<b>Institution (institution)</b>	
DFCI	92 (38%)
MSK	118 (49%)
VICC	31 (13%)
<b>Drugs in regimen (regimen_drugs)</b>	
Bevacizumab, Carboplatin, Pemetrexed Disodium	52 (22%)
Bevacizumab, Cisplatin, Pemetrexed Disodium	27 (11%)
Carboplatin, Pemetrexed Disodium	124 (51%)
Cisplatin, Pemetrexed Disodium	38 (16%)
<sup>1</sup> n (%)	



nsccl\_cohort  
\$tbl\_ngs

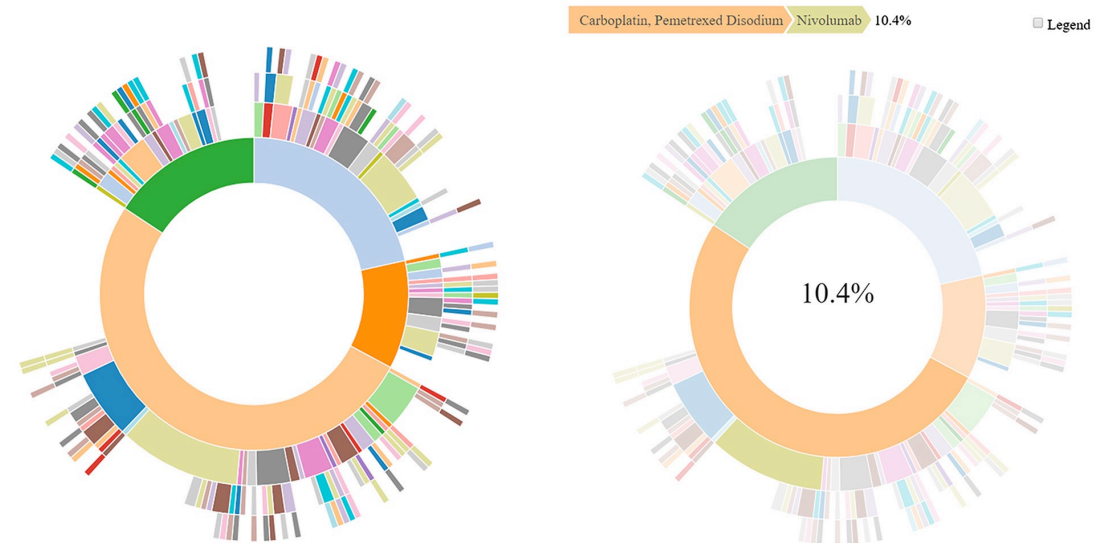
Characteristic	N = 262 Cancer Panel Tests <sup>1</sup>
<b>Cohort (cohort)</b>	
NSCLC	262 (100%)
<b>Institution (institution)</b>	
DFCI	99 (38%)
MSK	126 (48%)
VICC	37 (14%)
<b>OncoTree code (cpt_oncotree_code)</b>	
LCLC	1 (0.4%)
LUAD	253 (97%)
LUAS	1 (0.4%)
LUSC	1 (0.4%)
NSCLC	4 (1.5%)
NSCLCPD	2 (0.8%)
<b>Sequence assay ID (cpt_seq_assay_id)</b>	
DFCI-ONCOPANEL-1	1 (0.4%)
DFCI-ONCOPANEL-2	57 (22%)
DFCI-ONCOPANEL-3	41 (16%)
MSK-IMPACT341	3 (1.1%)
MSK-IMPACT410	61 (23%)
MSK-IMPACT468	62 (24%)
VICC-01-SOLIDTUMOR	26 (9.9%)
VICC-01-T5A	1 (0.4%)
VICC-01-T7	10 (3.8%)
<sup>1</sup> n (%)	

# Clinico-Genomic Data Processing Pipeline



# drug\_regimen\_sunburst()

- Visualize the complete treatment course for selected cancer diagnoses
- Each ring corresponds to a regimen (i.e., innermost ring is first regimen, second innermost ring is second regimen, etc.)
- Interactive figure: Can hover to see regimen names and percent of patients receiving that regimen



# drug\_regimen\_sunburst()

Argument	Description	Acceptable Values
data_synapse	List returned from pull_data_synapse()	<ul style="list-style-type: none"><li>Name of object in global environment that was returned from pull_data_synapse()</li></ul>
data_cohort	The list returned from the create_analytic_cohort() function call	<ul style="list-style-type: none"><li>Name of object in global environment that was returned from create_analytic_cohort()</li></ul>
max_n_regimens	The maximum number of regimens displayed in the sunburst plot	<ul style="list-style-type: none"><li>Integer &gt;0</li></ul>

# Demo: **drug\_regimen\_sunburst()** for case study using NSCLC 2.0-public data

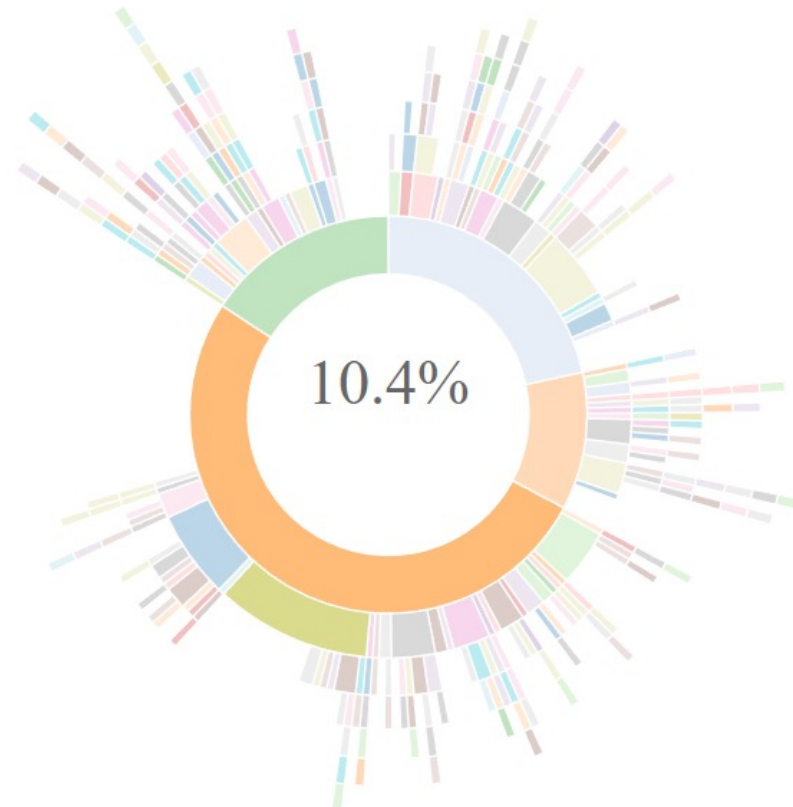
**Case Study:** Create a cohort of patients who were diagnosed with **Stage IV adenocarcinoma NSCLC** and received **Carboplatin and Pemetrexed +/- Bevacizumab or Cisplatin and Pemetrexed +/- Bevacizumab** as their **first** cancer-directed drug regimen after diagnosis

```
nsccl_sunburst <- drug_regimen_sunburst(  
  data_synapse = nsccl_synapse_data$NSCLC_v2.0,  
  data_cohort = nsccl_cohort)
```



```
nscic_sunburst$  
sunburst_plot
```

Carboplatin, Pemetrexed Disodium Nivolumab 10.4%

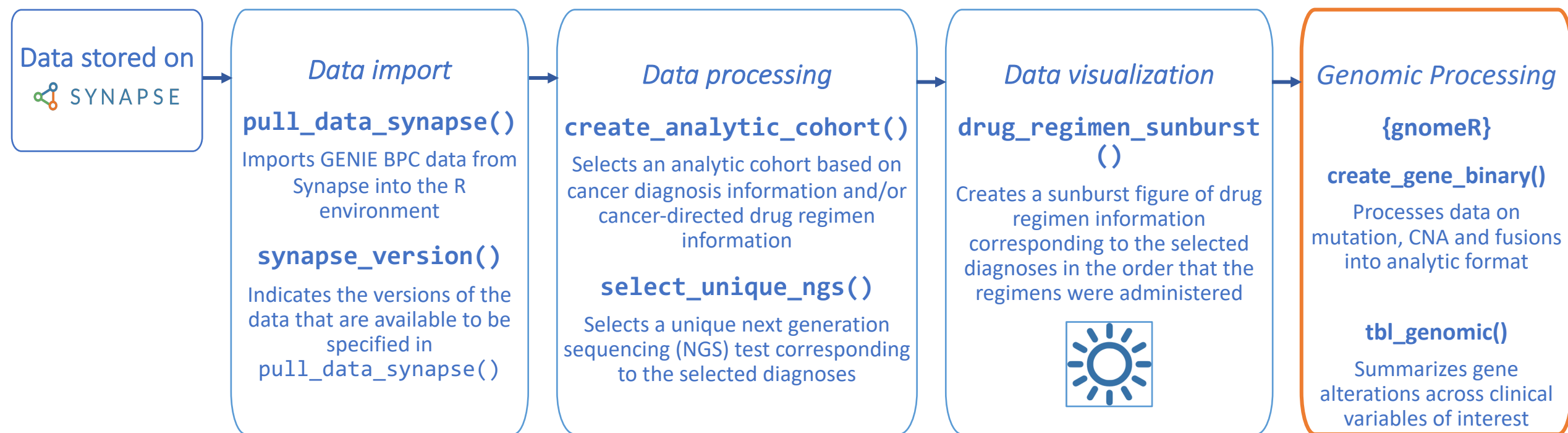




# Genomic Data Processing With {gnomeR}

---

# Clinico-Genomic Data Processing Pipeline





# Case Study

Using the cohort of patients who were diagnosed with Stage IV adenocarcinoma NSCLC, we will:

1. **Process data** into a analysis-ready matrix of gene alteration events
2. **Summarize genomic alteration frequencies** and analyze differences between males and females

```
1 nslc_cohort <- create_analytic_cohort(  
2   data_synapse =  
3     nslc_synapse_data$NSCLC_v2.0,  
4   stage_dx = c("Stage IV"),  
5   histology = "Adenocarcinoma",  
6   regimen_drugs =  
7     c("Carboplatin, Pemetrexed Disodium",  
8       "Cisplatin, Pemetrexed Disodium",  
9       "Bevacizumab, Carboplatin, Pemetrexed Disodium",  
10      "Bevacizumab, Cisplatin, Pemetrexed Disodium"),  
11   regimen_type = "Exact",  
12   regimen_order = 1,  
13   regimen_order_type = "within cancer",  
14   return_summary = TRUE  
15 )
```

Characteristic	N = 241 <sup>1</sup>
naaccr_sex_code	
Female	145 (60%)
Male	96 (40%)
<sup>1</sup> n (%)	

# Overview of Genomic Data

We will be processing and analyzing data on:

## 1) Mutations

```
1 mutations <- nsclc_synapse_data$NSCLC_v2.0$mutations_extended
```

## 2) Discrete Copy Number Alterations

```
1 cna <- nsclc_synapse_data$NSCLC_v2.0$cna
```

## 3) Fusions

```
1 fusions <- nsclc_synapse_data$NSCLC_v2.0$fusions
```

# Processing Data

# Issues When Processing Multi-Institutional Genomic Data

i

i

i

i

# Issues When Processing Multi-Institutional Genomic Data



## 1. Multiple Samples Per Patient



# Issues When Processing Multi-Institutional Genomic Data



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses



# Issues When Processing Multi-Institutional Genomic Data



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses



## 2. Data Formats & Gene Standards Often Inconsistent



# Issues When Processing Multi-Institutional Genomic Data



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses



## 2. Data Formats & Gene Standards Often Inconsistent

- Column names, data formats and gene names may differ between studies or even within studies





# Issues When Processing Multi-Institutional Genomic Data



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses



## 2. Data Formats & Gene Standards Often Inconsistent

- Column names, data formats and gene names may differ between studies or even within studies



## 3. Cohort Inclusion



# Issues When Processing Multi-Institutional Genomic Data



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses



## 2. Data Formats & Gene Standards Often Inconsistent

- Column names, data formats and gene names may differ between studies or even within studies



## 3. Cohort Inclusion

- Samples with no alterations may be dropped when pulling data



# Issues When Processing Multi-Institutional Genomic Data



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses



## 2. Data Formats & Gene Standards Often Inconsistent

- Column names, data formats and gene names may differ between studies or even within studies



## 3. Cohort Inclusion

- Samples with no alterations may be dropped when pulling data



## 4. Multi-Institutional Studies Use Several Gene Panels

# Issues When Processing Multi-Institutional Genomic Data



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses



## 2. Data Formats & Gene Standards Often Inconsistent

- Column names, data formats and gene names may differ between studies or even within studies



## 3. Cohort Inclusion

- Samples with no alterations may be dropped when pulling data



## 4. Multi-Institutional Studies Use Several Gene Panels

- Samples may be sequenced using different panels, therefore the non-overlapping genes have to be annotated as missing

# Select One Sample Per Patient



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses

Argument	Description
data_cohort	Output object of the create_analytic_cohort function.
oncotree_code	Character vector specifying which sample OncoTree codes to keep. See 'cpt_oncotree_code' column of data_cohort.
sample_type	Character specifying which type of genomic sample to prioritize. Options are 'Primary', 'Local', and 'Metastasis'. Default is to not select a NGS sample based on the sample type.
min_max_time	Character specifying if the first or last genomic sample recorded should be kept. Options are 'min' (first) and 'max' (last).

# Select One Sample Per Patient



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses
- Patients can have many NGS reports

Argument	Description
data_cohort	Output object of the create_analytic_cohort function.
oncotree_code	Character vector specifying which sample OncoTree codes to keep. See 'cpt_oncotree_code' column of data_cohort.
sample_type	Character specifying which type of genomic sample to prioritize. Options are 'Primary', 'Local', and 'Metastasis'. Default is to not select a NGS sample based on the sample type.
min_max_time	Character specifying if the first or last genomic sample recorded should be kept. Options are 'min' (first) and 'max' (last).

# Select One Sample Per Patient



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses
- Patients can have many NGS reports
- We can use `select_unique_ngs()` to select 1 sample per patient

Argument	Description
<code>data_cohort</code>	Output object of the <code>create_analytic_cohort</code> function.
<code>oncotree_code</code>	Character vector specifying which sample OncoTree codes to keep. See 'cpt_oncotree_code' column of <code>data_cohort</code> .
<code>sample_type</code>	Character specifying which type of genomic sample to prioritize. Options are 'Primary', 'Local', and 'Metastasis'. Default is to not select a NGS sample based on the sample type.
<code>min_max_time</code>	Character specifying if the first or last genomic sample recorded should be kept. Options are 'min' (first) and 'max' (last).

# Select One Sample Per Patient



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses
- Patients can have many NGS reports
- We can use `select_unique_ngs()` to select 1 sample per patient
- This function prioritizes characteristics of interest (e.g. sample type)

Argument	Description
<code>data_cohort</code>	Output object of the <code>create_analytic_cohort</code> function.
<code>oncotree_code</code>	Character vector specifying which sample OncoTree codes to keep. See 'cpt_oncotree_code' column of <code>data_cohort</code> .
<code>sample_type</code>	Character specifying which type of genomic sample to prioritize. Options are 'Primary', 'Local', and 'Metastasis'. Default is to not select a NGS sample based on the sample type.
<code>min_max_time</code>	Character specifying if the first or last genomic sample recorded should be kept. Options are 'min' (first) and 'max' (last).



# Select One Sample Per Patient



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses
- Patients can have many NGS reports
- We can use `select_unique_ngs()` to select 1 sample per patient
- This function prioritizes characteristics of interest (e.g. sample type)
- If patient only has one report, it will be returned regardless of criteria

Argument	Description
<code>data_cohort</code>	Output object of the <code>create_analytic_cohort</code> function.
<code>oncotree_code</code>	Character vector specifying which sample OncoTree codes to keep. See 'cpt_oncotree_code' column of <code>data_cohort</code> .
<code>sample_type</code>	Character specifying which type of genomic sample to prioritize. Options are 'Primary', 'Local', and 'Metastasis'. Default is to not select a NGS sample based on the sample type.
<code>min_max_time</code>	Character specifying if the first or last genomic sample recorded should be kept. Options are 'min' (first) and 'max' (last).

# Select One Sample Per Patient



## 1. Multiple Samples Per Patient

- We often need to select a single sample per patient for analyses

```
1 nrow(nsclc_cohort$cohort_ngs)
```

[1] 262

```
1 nsclc_samp <- select_unique_ngs(  
2   data_cohort = nsclc_cohort$cohort_ngs,  
3   oncotree_code = "LUAD",  
4   sample_type = "Metastasis",  
5   min_max_time = "max"  
6 )  
7  
8 nrow(nsclc_samp)
```

[1] 241



# Format Data in Analysis-ready Matrix

sample_id	CREBBP	GLI2	KRAS	MAP3K1	PIK3C2B	PBRM1	MYC.Del
GENIE-DFCI-004022-1313	0	0	0	0	0	NA	0
GENIE-DFCI-000013-8840	0	0	1	0	0	0	0
GENIE-DFCI-000136-6004	0	0	1	0	0	0	0
GENIE-DFCI-000215-8010	0	0	0	0	0	0	0
GENIE-DFCI-000381-9526	0	0	0	0	0	0	0

# Format Data in Analysis-ready Matrix

- Next we want to get our genomic data in an **analysis-friendly format**.

sample_id	CREBBP	GLI2	KRAS	MAP3K1	PIK3C2B	PBRM1	MYC.Del
GENIE-DFCI-004022-1313	0	0	0	0	0	NA	0
GENIE-DFCI-000013-8840	0	0	1	0	0	0	0
GENIE-DFCI-000136-6004	0	0	1	0	0	0	0
GENIE-DFCI-000215-8010	0	0	0	0	0	0	0
GENIE-DFCI-000381-9526	0	0	0	0	0	0	0

# Format Data in Analysis-ready Matrix

- Next we want to get our genomic data in an **analysis-friendly format**.
- `create_gene_binary()` from {gnomeR} will give us a data frame of  **$n$  patients  $\times$   $p$  alterations**.

sample_id	CREBBP	GLI2	KRAS	MAP3K1	PIK3C2B	PBRM1	MYC.Del
GENIE-DFCI-004022-1313	0	0	0	0	0	NA	0
GENIE-DFCI-000013-8840	0	0	1	0	0	0	0
GENIE-DFCI-000136-6004	0	0	1	0	0	0	0
GENIE-DFCI-000215-8010	0	0	0	0	0	0	0
GENIE-DFCI-000381-9526	0	0	0	0	0	0	0

# Format Data in Analysis-ready Matrix

- Next we want to get our genomic data in an **analysis-friendly format**.
- `create_gene_binary()` from {gnomeR} will give us a data frame of ***n* patients x *p* alterations**.
- Alteration columns are denoted by the **gene name** if mutation (e.g. **TP53**) or **gene name + .Amp, .fus, .Del** (**TP53.Del**) for other alterations types.

sample_id	CREBBP	GLI2	KRAS	MAP3K1	PIK3C2B	PBRM1	MYC.Del
GENIE-DFCI-004022-1313	0	0	0	0	0	NA	0
GENIE-DFCI-000013-8840	0	0	1	0	0	0	0
GENIE-DFCI-000136-6004	0	0	1	0	0	0	0
GENIE-DFCI-000215-8010	0	0	0	0	0	0	0
GENIE-DFCI-000381-9526	0	0	0	0	0	0	0

# Format Data in Analysis-ready Matrix

- Next we want to get our genomic data in an **analysis-friendly format**.
- `create_gene_binary()` from {gnomeR} will give us a data frame of ***n* patients x *p* alterations**.
- Alteration columns are denoted by the **gene name** if mutation (e.g. **TP53**) or **gene name + .Amp, .fus, .Del** (**TP53.Del**) for other alterations types.
- Each cell will have **0** if no alteration, **1** if altered, or **NA** if that gene was not tested in that patient.

sample_id	CREBBP	GLI2	KRAS	MAP3K1	PIK3C2B	PBRM1	MYC.Del
GENIE-DFCI-004022-1313	0	0	0	0	0	NA	0
GENIE-DFCI-000013-8840	0	0	1	0	0	0	0
GENIE-DFCI-000136-6004	0	0	1	0	0	0	0
GENIE-DFCI-000215-8010	0	0	0	0	0	0	0
GENIE-DFCI-000381-9526	0	0	0	0	0	0	0

# Get Data in Standardized Format



## 2. Data Formats & Gene Standards Often Inconsistent

- Column names, data formats and gene names may differ between studies or even within studies!

**{gnomeR}** functions are designed to work with standard alteration data formats used in common platforms like cBioPortal.

GENIE CNA and fusion data are slightly different than this standard, therefore we need to reformat it using:

- `gnomeR::reformat_fusions()`
- `gnomeR::pivot_cna_longer()`



# Get Data in Standardized Format

Format fusions to follow the cBioPortal standard:

Before:

1 fusions

Hugo_Symbol	Tumor_Sample_Barcode	Fusion
TP53	GENIE-MSK-P-0004827-T01-IM5	TP53-intragenic
XPO1	GENIE-MSK-P-0004827-T01-IM5	XPO1-USP34 fusion
USP34	GENIE-MSK-P-0004827-T01-IM5	XPO1-USP34 fusion

After:

1 reformat\_fusion(fusions)

sample_id	site_1_hugo_symbol	site_2_hugo_symbol	fusion
GENIE-MSK-P-0004827-T01-IM5	TP53	NA	TP53
GENIE-MSK-P-0004827-T01-IM5	USP34	XPO1	USP34-XPO1

# Get Data in Standardized Format

Format CNA to follow the cBioPortal standard:

**Before:**

```
1 cna
```

Hugo_Symbol	GENIE.DFCI.003908.234520	GENIE.DFCI.002183.6917
APC	1	0
ARID1B	-1	0
BCL2	-1	0
TP53	0	0

**After:**

```
1 gnomeR::pivot_cna_longer(cna)
```

hugo_symbol	sample_id	alteration
APC	GENIE-DFCI-003908-234520	gain
ARID1B	GENIE-DFCI-003908-234520	loss
BCL2	GENIE-DFCI-003908-234520	loss
BCL6	GENIE-DFCI-002183-6917	gain
CARD11	GENIE-DFCI-002183-6917	gain

# | Get Data in Standardized Format

```
1 mutations <- nsclc_synapse_data$NSCLC_v2.0$mutations_extended
2 cna <- nsclc_synapse_data$NSCLC_v2.0$cna
3 fusions <- nsclc_synapse_data$NSCLC_v2.0$fusions
```

```
1 reformat_fusions <- gnomeR::reformat_fusion(fusions)
2
3 nrow(reformat_fusions)
```

[1] 538

```
1 reformat_cna <- gnomeR::pivot_cna_longer(cna)
2
3 nrow(reformat_cna)
```

[1] 38163



# Format Data in Analysis-ready Matrix

## Basic code:

```
1  gnomeR::create_gene_binary(  
2    mutation = mutations,  
3    cna = reformat_cna,  
4    fusion = reformat_fusions)
```

# Format Data in Analysis-ready Matrix

- Data is now in standardized format and can be processed using {gnomeR}  
`create_gene_binary()`

## Basic code:

```
1  gnomeR::create_gene_binary(  
2    mutation = mutations,  
3    cna = reformat_cna,  
4    fusion = reformat_fusions)
```

# Format Data in Analysis-ready Matrix

- Data is now in standardized format and can be processed using {gnomeR}  
`create_gene_binary()`
- We will add additional arguments to `create_gene_binary()` to help address remaining data processing issues

## Basic code:

```
1  gnomeR::create_gene_binary(  
2    mutation = mutations,  
3    cna = reformat_cna,  
4    fusion = reformat_fusions)
```

# Cohort Inclusion



## 3. Cohort Inclusion

- Samples with no alterations may be dropped when pulling raw genomic data

```
1  gnomeR::create_gene_binary(  
2    mutation = mutations,  
3    cna = reformat_cna,  
4    fusion = reformat_fusions  
5    samples = nsclc_samp$cpt_genie_sample_id)
```

The **samples** argument will ensure all study IDs have a row in resulting analysis data, even if they are not present in genomic files

# Not All Patients Tested on Same Panel



## 4. Multi-Institutional Studies Use Several Gene Panels

- Samples may be sequenced using different panels therefore the non overlapping genes have to be annotated as missing

```
1  gnomeR::create_gene_binary(  
2    mutation = mutations,  
3    cna = reformat_cna,  
4    fusion = reformat_fusions,  
5    samples = nsclc_samp$cpt_genie_sample_id,  
6    specify_panel = nsclc_panels)
```

The `specify_panels` argument can insert **NAs** when we know that gene was not tested for a specific set of patients.



# Not All Patients Tested on Same Panel

To use `specify_panels`, we first need to create a data frame indicating which patient IDs were sequenced on which panels.

```
1 nslc_panels <- data.frame(  
2   sample_id = nslc_samp$cpt_genie_sample_id,  
3   panel_id = nslc_samp$cpt_seq_assay_id) %>%  
4   mutate(panel_id = ifelse(!is.na(panel_id),  
5                             panel_id, "no"))  
6  
7 nslc_panels %>% head()
```

sample_id	panel_id
GENIE-DFCI-000013-8840	DFCI-ONCOPANEL-2
GENIE-DFCI-000136-6004	DFCI-ONCOPANEL-2
GENIE-DFCI-000215-8010	DFCI-ONCOPANEL-2
GENIE-DFCI-000381-9526	DFCI-ONCOPANEL-2
GENIE-DFCI-000410-10003	DFCI-ONCOPANEL-2
GENIE-DFCI-000583-11175	DFCI-ONCOPANEL-2

# Not All Patients Tested on Same Panel

## Without Panel Annotation

```
1 binmat1 <- gnomeR::create_gene_binary(  
2   mutation = mutations,  
3   cna = reformat_cna,  
4   fusion = reformat_fusions,  
5   samples = nsclc_samp$cpt_genie_sample_id,  
6   specify_panel = "no")
```

sample_id	panel_id	GLI2	KRAS	PIK3C2B	PBRM1
GENIE-DFCI-004022-1313	DFCI-ONCOPANEL-1	0	0	0	0
GENIE-DFCI-000013-8840	DFCI-ONCOPANEL-2	0	1	0	0
GENIE-MSK-P-0002725-T01-IM3	MSK-IMPACT341	0	1	0	0
GENIE-MSK-P-0017722-T02-IM6	MSK-IMPACT468	0	0	0	0

## With Panel Annotation

```
1 binmat2 <- gnomeR::create_gene_binary(  
2   mutation = mutations,  
3   cna = reformat_cna,  
4   fusion = reformat_fusions,  
5   samples = nsclc_samp$cpt_genie_sample_id,  
6   specify_panel = nsclc_panels)
```

sample_id	panel_id	GLI2	KRAS	PIK3C2B	PBRM1
GENIE-DFCI-004022-1313	DFCI-ONCOPANEL-1	0	0	0	NA
GENIE-DFCI-000013-8840	DFCI-ONCOPANEL-2	0	1	0	0
GENIE-MSK-P-0002725-T01-IM3	MSK-IMPACT341	NA	1	NA	0
GENIE-MSK-P-0017722-T02-IM6	MSK-IMPACT468	NA	0	NA	0

# Ensure Gene Names Are Consistent Across Studies



## 2. Data Formats & Gene Standards Often Inconsistent

- Column names, data formats and gene names may differ between studies or even within studies.

```
1 no_recode <- gnomer::create_gene_binary(  
2   samples = nsclc_samp$cpt_genie_sample_id,  
3   mutation = mutations,  
4   cna = reformat_cna,  
5   fusion = reformat_fusions,  
6   specify_panel = nsclc_panels,  
7   recode_aliases = "no")
```

```
1 recode <- gnomer::create_gene_binary(  
2   samples = nsclc_samp$cpt_genie_sample_id,  
3   mutation = mutations,  
4   cna = reformat_cna,  
5   fusion = reformat_fusions,  
6   specify_panel = nsclc_panels,  
7   recode_aliases = "impact")
```

```
1 setdiff(names(no_recode), names(recode)) %>% head()
```

[1] "MRE11A" "RFWD2" "H3F3A" "FAM46C" "HIST1H3D" "WHSC1L1"

# | Process Data: Final Dataset

Let's run `create_gene_binary()` with the `samples`, `specify_panel` and `recode_aliases` arguments.

First we create `nsclc_panels`:

```
1 nsclc_panels <- data.frame(  
2   sample_id = nsclc_samp$cpt_genie_sample_id,  
3   panel_id = nsclc_samp$cpt_seq_assay_id) %>%  
4   mutate(panel_id = ifelse(!is.na(panel_id),  
5                             panel_id, "no"))
```

Then run `create_gene_binary()`:

```
1 gene_binary <- gnomer::create_gene_binary(  
2   mutation = mutations,  
3   cna = reformat_cna,  
4   fusion = reformat_fusions,  
5   samples = nsclc_samp$cpt_genie_sample_id,  
6   specify_panel = nsclc_panels,  
7   recode_aliases = "impact")
```



# Analyzing Data

# Issues When Analyzing Multi-Institutional Genomic Data



# Issues When Analyzing Multi-Institutional Genomic Data



## 1. Multiple Testing and False Positives



# Issues When Analyzing Multi-Institutional Genomic Data



## 1. Multiple Testing and False Positives

- Many hypothesis tests done simultaneously can lead to high rate of false positive findings.





# Issues When Analyzing Multi-Institutional Genomic Data



## 1. Multiple Testing and False Positives

- Many hypothesis tests done simultaneously can lead to high rate of false positive findings.
- Very low prevalence genes are often not very informative.



# Issues When Analyzing Multi-Institutional Genomic Data



## 1. Multiple Testing and False Positives

- Many hypothesis tests done simultaneously can lead to high rate of false positive findings.
- Very low prevalence genes are often not very informative.
- Choose a threshold (e.g. 1% or 5 %) *a priori* to limit number of hypothesis tests and consider reporting a q-value (adjusted for multiple testing).



# Issues When Analyzing Multi-Institutional Genomic Data



## 1. Multiple Testing and False Positives

- Many hypothesis tests done simultaneously can lead to high rate of false positive findings.
- Very low prevalence genes are often not very informative.
- Choose a threshold (e.g. 1% or 5 %) *a priori* to limit number of hypothesis tests and consider reporting a q-value (adjusted for multiple testing).



## 2. Limited Power To Detect Clinical Associations When Sparse Alterations

# Issues When Analyzing Multi-Institutional Genomic Data



## 1. Multiple Testing and False Positives

- Many hypothesis tests done simultaneously can lead to high rate of false positive findings.
- Very low prevalence genes are often not very informative.
- Choose a threshold (e.g. 1% or 5 %) *a priori* to limit number of hypothesis tests and consider reporting a q-value (adjusted for multiple testing).



## 2. Limited Power To Detect Clinical Associations When Sparse Alterations

- If biologically meaningful, you may want to summarize on gene or pathway level

# Case Study

# Case Study

- We will use the processed binary data frame data (`gene_binary`) to **summarize genomic alterations** overall in the cohort, and **by sex**.

# Case Study

- We will use the processed binary data frame data (`gene_binary`) to **summarize genomic alterations** overall in the cohort, and **by sex**.
- First, we need to join clinical data on `sex` to genomic data

# Add Clinical Variable To Data

```
1 # get patient IDs and sample IDs
2 patient_index <- nscic_cohort$cohort_ngs %>%
3   select(record_id, cpt_genie_sample_id)
4
5 # Join sex data to patient ID index
6 select_clinical <- nscic_cohort$cohort_pt_char %>%
7   select(record_id, naaccr_sex_code) %>%
8   left_join(patient_index)
9
10 # Join all to gene binary data
11 gene_binary <- gene_binary %>%
12   left_join(select_clinical,
13             by = c("sample_id"= "cpt_genie_sample_id")) %>%
14   select(-record_id)
15
16 gene_binary <- gene_binary %>%
17   select(sample_id, naaccr_sex_code, everything())
18
19 gene_binary %>%
20   select(naaccr_sex_code) %>%
21   tbl_summary()
```

Characteristic	N = 241 <sup>1</sup>
naaccr_sex_code	
Female	145 (60%)
Male	96 (40%)
<sup>1</sup> n (%)	





# Subset By a Prevalence Threshold



## 1. Multiple Testing and False Positives

- Use `subset_by_frequency(t)` to **subset genes** above a given threshold
- `t` indicates a **prevalence threshold** between 0 (`t = 0`) and 100% (`t = 1`)
- `other_vars` retains the clinical variable of interest in the resulting data set

# | Subset By 40% Threshold

```
1 ncol(gene_binary)
```

[1] 1403

```
1 nsclc_subset <- gene_binary %>%  
2   subset_by_frequency(t = .4, other_vars = naaccr_sex_code)  
3  
4 ncol(nsclc_subset)
```

[1] 6



# Subset by a Panel

Use `subset_by_panel()` to subset genes in a given targeted panel.

```
1 ncol(gene_binary)
```

[1] 1403

```
1 nsclc_subset_panel <- gene_binary %>%  
2   subset_by_panel(panel_id = 'IMPACT300', other_vars = naaccr_sex_code)  
3  
4 ncol(nsclc_subset_panel)
```

[1] 220

# Summarize Alterations with `tbl_genomic()`

- `tbl_genomic()` is a wrapper function for `gtsummary::tbl_summary()` specifically designed for presenting genomic data
- You can use any `{gtsummary}` function on top of `tbl_genomic()` to customize the table (e.g. `bold_labels()`)

# | Summarize Alterations with tbl\_genomic

Create a simple `tbl_genomic` object, then bold the labels.

```
1 nsc1c_subset %>%  
2   select(-naaccr_sex_code) %>%  
3   tbl_genomic() %>%  
4   bold_labels()
```

Characteristic	N = 241 <sup>1</sup>
TP53	124 (51%)
GBA.Amp	17 (46%)
Unknown	204
KRAS	98 (41%)
JAZF1.Amp	15 (41%)
Unknown	204
<sup>1</sup> n (%)	

# Summarize Genes with tbl\_genomic()

You may want to analyze on the **gene level** instead of the alteration level.

Use `summarize_by_gene()` first, then pass to `tbl_genomic()`:

 **Note:** `summarize_by_gene()` should come before passing to `subset_by_frequency()`

```
1 tbl_gene <- gene_binary %>%
2   select(-naaccr_sex_code) %>%
3   summarize_by_gene() %>%
4   subset_by_frequency(t = .4) %>%
5   tbl_genomic()
6
7 tbl_gene
```

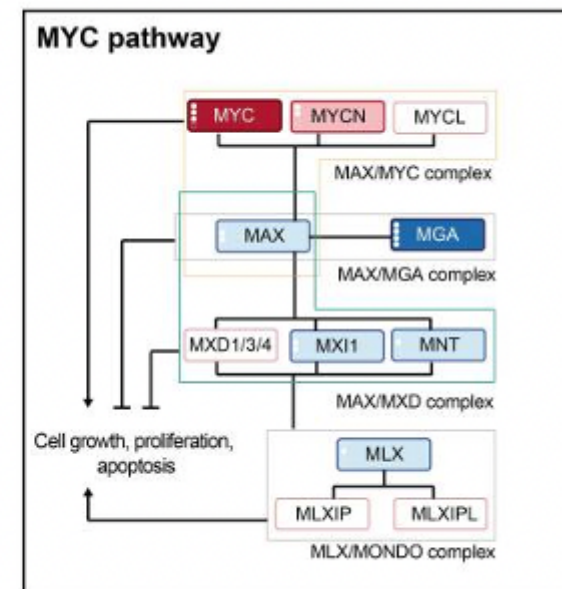
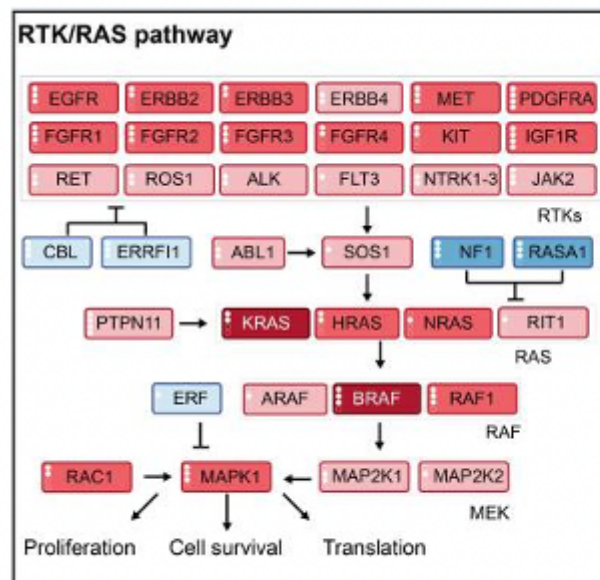
Characteristic	N = 241 <sup>1</sup>
TP53	135 (56%)
GBA	17 (46%)
Unknown	204
PTK2B	17 (46%)
Unknown	204
KRAS	107 (44%)
WRN	38 (41%)
Unknown	149
JAZF1	15 (41%)
Unknown	204
<sup>1</sup> n (%)	

# Summarize Oncogenic Pathways

## 2. Limited Power To Detect Clinical Associations When Sparse Alterations

- If biologically meaningful, you may want to summarize on Alteration, Gene or Pathway Level

Often we want to analyze alterations on the pathway level.



# Summarize Oncogenic Pathways



## 2. Limited Power To Detect Clinical Associations When Sparse Alterations

- If biologically meaningful, you may want to summarize on Alteration, Gene or Pathway Level

{gnomeR} offers several default pathways that can be added with `add_pathways()`

```
1 paths <- gnomeR::pathways %>% names()  
2 paths
```

[1] “RTK/RAS” “Nrf2” “PI3K” “TGFB” “p53”

[6] “Wnt” “Myc” “Cell cycle” “Hippo” “Notch”



# Summarize Oncogenic Pathways

You can also add a custom pathway:

```
1 path_df <- gene_binary %>%  
2   add_pathways(custom_pathways = c("SPOP.mut", "FOXA1.mut"))  
3  
4 path_df %>% select("pathway_custom") %>%  
5   tbl_summary()
```

Characteristic	N = 241 <sup>1</sup>
pathway_custom	2 (0.8%)
<sup>1</sup> n (%)	

Note: You must specify **.mut**, **.Amp**, **.Del** for alterations **custom\_pathways**

# | Summarize Alteration Pathways

```
1 path_df <- gene_binary %>%  
2   select(-naaccr_sex_code) %>%  
3   add_pathways()  
4  
5 path_df %>%  
6   select(contains("pathway")) %>%  
7   tbl_summary() %>%  
8   bold_labels()
```



# Summarize Alteration Pathways

```
1 path_df <- gene_binary %>%
2   select(-naaccr_sex_code) %>%
3   add_pathways()
4
5 path_df %>%
6   select(contains("pathway")) %>%
7   tbl_summary() %>%
8   bold_labels()
```

Characteristic	N = 241 <sup>1</sup>
pathway_RTK/RAS	217 (90%)
pathway_Nrf2	48 (20%)
pathway_PI3K	114 (47%)
pathway_TGFB	34 (14%)
pathway_p53	181 (75%)
pathway_Wnt	51 (21%)
pathway_Myc	70 (29%)
pathway_Cell cycle	90 (37%)
pathway_Hippo	31 (13%)
pathway_Notch	72 (30%)
<sup>1</sup> n (%)	



# Comparing Alteration Frequencies Across Clinical Data

We can easily compare frequencies by sex using the `by` argument:

```
1 tbl_gene <- gene_binary %>%
2   subset_by_frequency(
3     t = .4,
4     other_vars = naaccr_sex_code) %>%
5   tbl_genomic(by = naaccr_sex_code) %>%
6   bold_labels()
7
8 tbl_gene
```

# Comparing Alteration Frequencies Across Clinical Data

We can easily compare frequencies by sex using the `by` argument:

```
1 tbl_gene <- gene_binary %>%
2   subset_by_frequency(
3     t = .4,
4     other_vars = naaccr_sex_code) %>%
5   tbl_genomic(by = naaccr_sex_code) %>%
6   bold_labels()
7
8 tbl_gene
```

Characteristic	Overall, N = 241 <sup>1</sup>	Female, N = 145 <sup>1</sup>	Male, N = 96 <sup>1</sup>
TP53	124 (51%)	72 (50%)	52 (54%)
GBA.Amp	17 (46%)	8 (36%)	9 (60%)
Unknown	204	123	81
KRAS	98 (41%)	66 (46%)	32 (33%)
JAZF1.Amp	15 (41%)	9 (41%)	6 (40%)
Unknown	204	123	81
<sup>1</sup> n (%)			

# Comparing Alteration Frequencies Across Clinical Data

## 1. Multiple Testing and False Positives

We can use {gtsummary}'s `add_p()` and `add_q()` for hypothesis testing

```
1 tbl_gene <- gene_binary %>%
2   subset_by_frequency(
3     t = .4,
4     other_vars = naaccr_sex_code) %>%
5   tbl_genomic(
6     by = naaccr_sex_code) %>%
7   bold_labels() %>%
8   add_p() %>%
9   add_q()
10
11 tbl_gene
```

Characteristic	Overall, N = 241 <sup>1</sup>	Female, N = 145 <sup>1</sup>	Male, N = 96 <sup>1</sup>	p-value <sup>2</sup>	q-value <sup>3</sup>
<b>TP53</b>	124 (51%)	72 (50%)	52 (54%)	0.5	0.7
<b>GBA.Amp</b>	17 (46%)	8 (36%)	9 (60%)	0.2	0.3
Unknown	204	123	81		
<b>KRAS</b>	98 (41%)	66 (46%)	32 (33%)	0.059	0.2
<b>JAZF1.Amp</b>	15 (41%)	9 (41%)	6 (40%)	>0.9	>0.9
Unknown	204	123	81		

<sup>1</sup> n (%)

<sup>2</sup> Pearson's Chi-squared test

<sup>3</sup> False discovery rate correction for multiple testing

# Comparing Alteration Frequencies Across Clinical Data

```
1 tbl_gene <- gene_binary %>%
2   subset_by_frequency(
3     t = .4,
4     other_vars = naaccr_sex_code) %>%
5   tbl_genomic(by = naaccr_sex_code) %>%
6   bold_labels() %>%
7   add_p() %>%
8   add_q()
9
10 tbl_gene
```

Characteristic	Overall, N = 241 <sup>1</sup>	Female, N = 145 <sup>1</sup>	Male, N = 96 <sup>1</sup>	p-value <sup>2</sup>	q-value <sup>3</sup>
<b>TP53</b>	124 (51%)	72 (50%)	52 (54%)	0.5	0.7
<b>GBA.Amp</b>	17 (46%)	8 (36%)	9 (60%)	0.2	0.3
Unknown	204	123	81		
<b>KRAS</b>	98 (41%)	66 (46%)	32 (33%)	0.059	0.2
<b>JAZF1.Amp</b>	15 (41%)	9 (41%)	6 (40%)	>0.9	>0.9
Unknown	204	123	81		

<sup>1</sup> n (%)

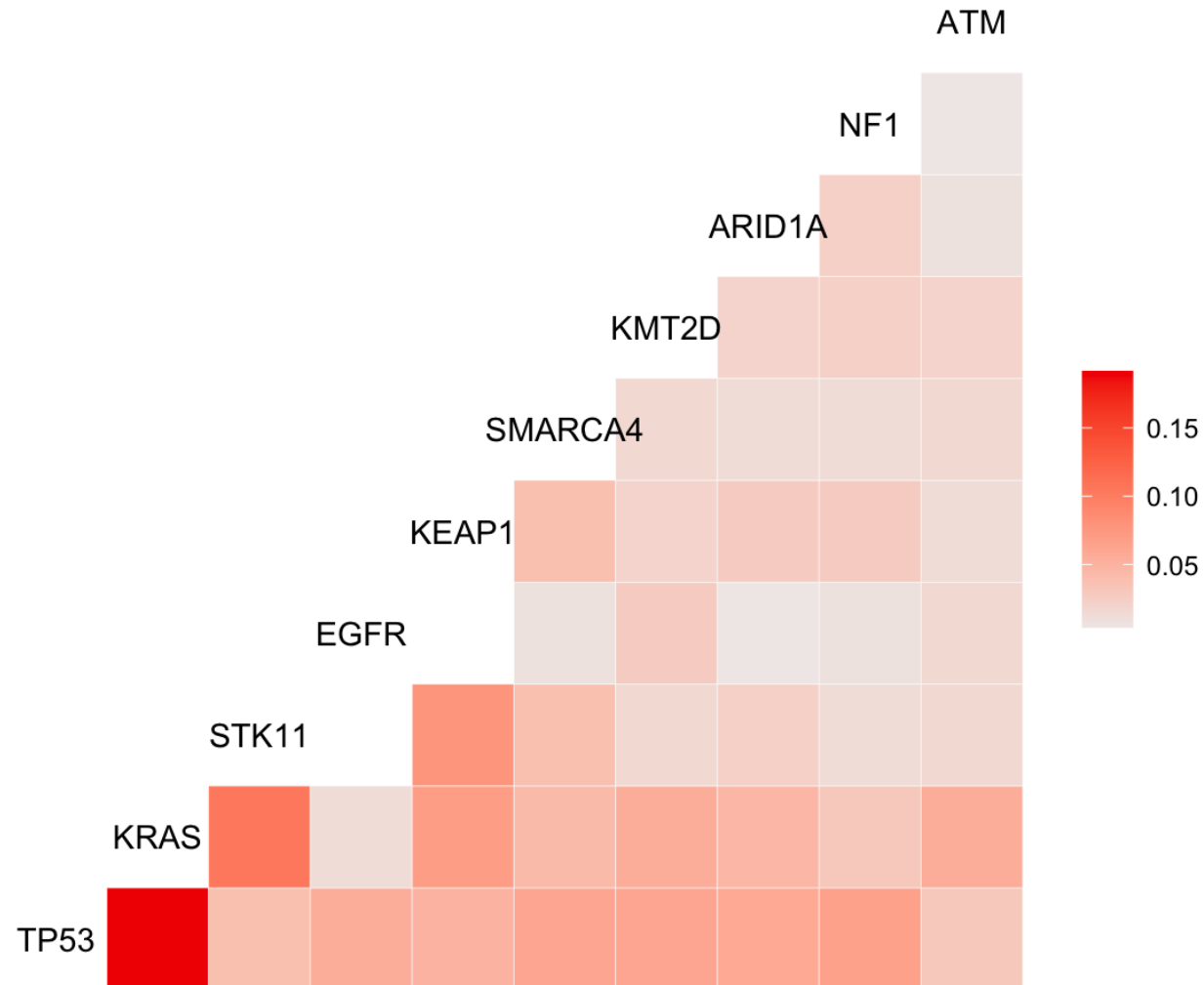
<sup>2</sup> Pearson's Chi-squared test

<sup>3</sup> False discovery rate correction for multiple testing



# Visualize Data

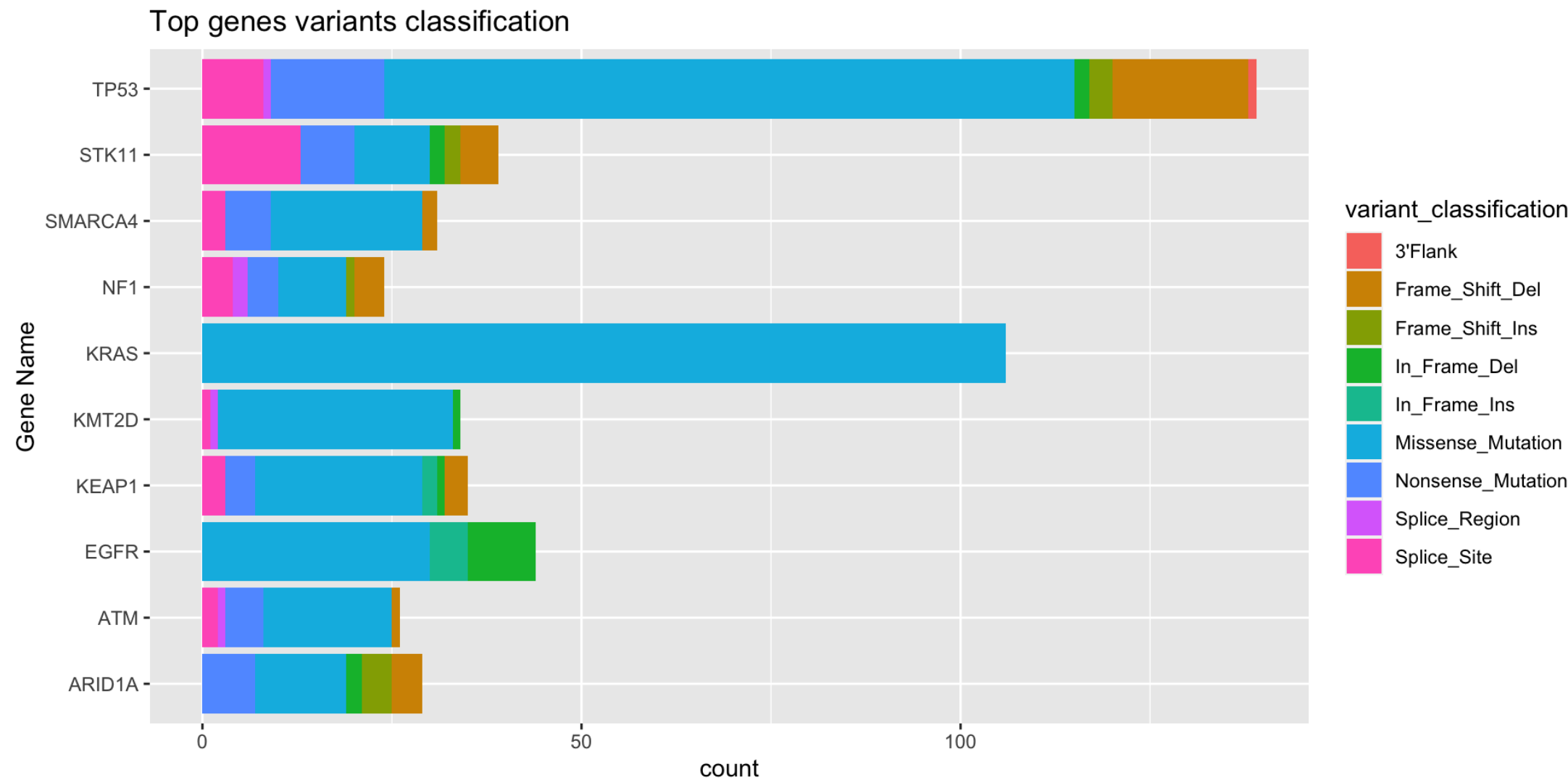
```
1  gnomeR::ggcomut(nscic_cohort$cohort_mutations_extended)
```





# Visualize Data

```
1  gnomer::ggtopgenes(nscic_cohort$cohort_mutations_extended)
```



# Additional Items & Next Steps

- Additional **visualizations** and **color palettes** useful for genomic data are available in {gnomeR} package
- Some data may require additional data checks. See [{gnomeR vignette}](#) for helpful tips on data QA.
- It may be appropriate to **oncoKB annotate** your data and only analyze oncogenic mutations (see [oncoKB.org](https://oncoKB.org) for more information).
- Some projects may utilize CNA Segmentation data. See [{gnomeR documentation}](#) for more information on available tools.

# Conclusion

- The **{genieBPC}** & **{gnomeR}** R packages offer a reproducible pipeline to create cohorts for clinico-genomics analyses
- **{genieBPC}** streamlines data access and clinical data processing from multiple clinical data files of varying structure to create analytic cohorts
- **{gnomeR}** facilitates annotation and analysis of complicated genomic data.



*Note: {gnomeR} can be used for genomic data processing and analyses outside of the GENIE BPC project!*

# Thank You!

Thanks to Hannah Fuchs for contributions to slides and contributing {gnomeR} & {genieBPC}  
authors: Michael Curry, Hannah Fuchs, Axel Martin, Arshi Arora