



G00x - Generalizable Germline-Targeting Clinical Trial Pipeline

- C de f r generating main text figures 2 thr ugh 7.
- C de f r generating the spreadsheet Data S10. "S urce data".

Second, the repository presents the G00x pipeline designed to facilitate the analysis of germline-targeting vaccine clinical trials. It is an all-in-one pipeline and analysis platform that parses, validates, calculates B cell frequencies, runs cellranger and combines all analyses into a pliable data frame. The pipeline supports datasets from both G002 and G003 trials (related to germline targeting to elicit VRC01-class antibodies against HIV) and ensures that all data is processed and validated to maintain the integrity of the clinical trial results. The pipeline is readily modifiable to enable handling of other types of germline-targeting trials, with other types of analyses, and can also be modified to enable data storage/data flow pathways different from that used in G002 and G003.

Data Access

If you don't want to run the pipeline, you can access the important data via the following links.

- The annotated, filtered and paired antibody sequences:
 - G002 [final_df.feather](#)
 - G003 [final_df.feather](#)
- Merged summary file with all frequencies reported in this study:
 - G002 [fl_w_and_sequencing_lng_names.csv](#)
 - G003 [fl_w_and_sequencing_lng_names.csv](#)

G002

To run testing, you will need a [g002](#) directory in the root of the project. You can get this with the sync command:

```
# this will sync the entire contents of the files into g002 directory
# Beware, this file is 1.5TB
aws sync --delete s3 s3://iavig002public/g002/ ./g002/
```

However, we have also setup so you don't need to download the entire directory. You can run the following to get a subset of the data.

```
# get the sorting directory
aws s3 cp --recursive s3://iavig002public/g002/G002/sorting ./g002/G002/sorting

# get the sequencing directory excluding large bcl and fastq files
aws s3 cp --recursive s3://iavig002public/g002/G002/sequencing ./g002/G002/sequencing --exclude *working_directory/*
--exclude *.fastq.gz --exclude *.tif --exclude *.cbcl --exclude *.imf1 --exclude *.filter --exclude *.bin --exclude
*Logs/* --exclude *_stdout --exclude *_stderr --exclude *Autofocus/* --exclude *Intensities/*
```

G003

To run testing, you will need a [g003](#) directory in the root of the project. You can get this with the sync command

```
# this will sync the entire contents of the files into g003 directory
# Beware, this file is 1.5TB
aws sync --delete s3 s3://iavig003public/g003/ ./g003/
```

However, we have also setup so you don't need to download the entire directory. You can run the following to get a subset of the data.

```
# get the sorting directory
aws s3 cp --recursive s3://iavig003public/g003/G003/sorting ./g003/G003/sorting

# get the sequencing directory excluding large bcl and fastq files
aws s3 cp --recursive s3://iavig003public/g003/G003/sequencing ./g003/G003/sequencing --exclude *working_directory/*
--exclude *.fastq.gz --exclude *.tif --exclude *.cbcl --exclude *.imf1 --exclude *.filter --exclude *.bin --exclude
*Logs/* --exclude *_stdout --exclude *_stderr --exclude *Autofocus/* --exclude *Intensities/*
```

Pipeline

Installation pre-requisites

While not necessary, we highly recommend using the [conda](#) open-source package and environment manager. For the purposes of this repository, only a minimal installer for anaconda is necessary (Miniconda).

[Miniconda command line installers](#)

Installation

This installation assumes that `git` and `conda` are in your path.

```
# clone the repository
git clone https://github.com/SchiefLab/G00x.git

# change directory
cd G00x

# this will create a conda environment called g00x and install the package
./install.sh
```

Validation

The most important part of this clinical trial working, is that everything is validated. That means that the data is validated, the code is validated, and the results are validated. This is a very important part of the process and should not be skipped.

FACS/Sorting validation

G002

These are instructions on how to validate the file structure containing the FACS/Sorting data before uploading to the box.

```
#Run the validator for flow from the command line
g00x g002 validate flow my_path/to/box/G002/

# If you used the AWS sync command above, you can use the command:
g00x g002 validate flow ./g002/G002/sorting/G002
```

Your folder structure should look like this.

```
# if you used the above commands to sync the data. The folder structure will look like this

g002/G002/sorting
└── G002
    ├── Prescreens
    │   ├── Prescreen_RunDate220825_UploadDate221021
    │   ├── Prescreen_RunDate220826_UploadDate221021
    ...
    └── Sorts
        ├── Sort_RunDate220927_UploadDate221013
        ├── Sort_RunDate220928_UploadDate221014
        ├── Sort_RunDate220929_UploadDate221014
        ├── Sort_RunDate220930_UploadDate221014
```

G003

```
#Run the validator for flow from the command line
g00x g003 validate flow my_path/to/G003/

# If you used the AWS sync command above, you can use the command:
g00x g003 validate flow ./g003/G003/sorting/G003
```

For G003, all the data is available in S3 bucket. The structure should look as below:

```
g003/G003/sorting
└── G003
    ├── Prescreens
    │   ├── Sort_RunDate230807_UploadDate230807
    │   ├── Sort_RunDate230808_UploadDate230808
    ...
    └── Sorts
        ├── Sort_RunDate230807_UploadDate230807
        ├── Sort_RunDate230808_UploadDate230808
        ├── Sort_RunDate230809_UploadDate230809
        ├── Sort_RunDate231003_UploadDate231027
```

Sequencing files validation

G002

In order to validate the sequencing files, it must be merged with the flow data. Thus, you will need both Glue and Box access. Validation can be accomplished by validating the merge command

```
g00x g002 sequencing -f /path/to/flow -s /path/to/sequencing -o output_file

# if you have the AWS structure
g00x g002 merge -s ./g002/G002/sequencing/G002 -f ./g002/G002/sorting/G002 -o my_merged_file
```

The sequencing folder structure should be as follows:

```
G002
└── run0002
    └── 221006_VH00497_31_AAAVKCLHV
        └── sample_manifest.csv
└── run0003
    └── 221019_VH00497_32_AAANGGVM5
        └── sample_manifest.csv
└── run0004
    ├── 221101_VL00414_3_AACFYLCM5
    └── 221103_VL00414_4_AAATJGYM5
        └── sample_manifest.csv
```

G003

In G003, all the data are stored in the S3 bucket; thus, you only need access.

```
g00x g003 sequencing -f /path/to/flow -s /path/to/sequencing -o output_file

# if you have the AWS structure
g00x g003 merge -s ./g003/G003/sequencing/G003 -f ./g003/G003/sorting/G003 -o my_merged_file
```

```
G003
└── run0001
    └── 230818_NB552490_0059_AHL7GFBGXM
        └── sequencing_manifest.csv
└── run0002
    └── 230821_NB552490_0060_AHL7GGBGXM
        └── sequencing_manifest.csv
└── run0003
    └── 231108_NB552490_0065_AHL7CCBGXM
        └── sequencing_manifest.csv
```

Using the above command, you will generate a file called `my_merged_file.csv` which will have the following fields.

```
ptid
group
weeks
visit_id
probe_set
sample_type
run_date
sort_pool
hashtag
run_dir_path
pool_number
sorted_date
vdj_sequencing_replicate
cso_sequencing_replicate
vdj_library_replicate
cso_library_replicate
bio_replicate
vdj_index
```


field	description
value	what is the value of the gate or frequency
hashtag	if this sample has been sorted, what is the hashtag
sort_p_l	the sorting position, P01 - P10

BCR sequence analysis

G002

```
g00x g002 pipeline demultiplex -f ./g002/G002/sorting/G002/ -s ./g002/G002/sequencing/G002/ -o ./g002/G002/output/demultiplex

g00x g002 pipeline vdj -d ./g002/G002/output/demultiplex.feather -o ./g002/G002/output/vdj

g00x g002 pipeline cso -d ./g002/G002/output/demultiplex.feather -o ./g002/G002/output/cso

# Merge VDJ and CSO dataframes, run the output through SADIE for
g00x g002 pipeline airr -k 3 -c ./g002/G002/output/cso.feather -v ./g002/G002/output/vdj.feather -o ./g002/G002/output/final_df

# Output merged.feather not used to generate figures, but is a sanity check.
g00x g002 validate merge -s ./g002/G002/sequencing/G002/ -f ./g002/G002/sorting/G002 -o ./g002/G002/output/merged

# Output flow_and_sequencing.feather used to generate figures; will contain counts, frequencies, and metadata from
# flow and sequencing data.
g00x g002 analysis report -s ./g002/G002/output/final_df.feather -f ./g002/G002/output/flow_output.feather -o ./g002/G002/output/flow_and_sequencing
```

G003

```
g00x g003 pipeline demultiplex -f ./g003/G003/sorting/G003/ -s ./g003/G003/sequencing/G003/ -o ./g003/G003/output/demultiplex

g00x g003 pipeline vdj -d ./g003/G003/output/demultiplex.feather -o ./g003/G003/output/vdj

g00x g003 pipeline cso -d ./g003/G003/output/demultiplex.feather -o ./g003/G003/output/cso

# Merge VDJ and CSO dataframes, run the output through SADIE for
g00x g003 pipeline airr -k 3 -c ./g003/G003/output/cso.feather -v ./g003/G003/output/vdj.feather -o ./g003/G003/output/final_df

# Output merged.feather not used to generate figures, but is a sanity check.
g00x g003 validate merge -s ./g003/G003/sequencing/G003/ -f ./g003/G003/sorting/G003 -o ./g003/G003/output/merged

# Output flow_and_sequencing.feather used to generate figures; will contain counts, frequencies, and metadata from
# flow and sequencing data.
g00x g003 analysis report -s ./g003/G003/output/final_df.feather -f ./g003/G003/output/flow_output.feather -o ./g003/G003/output/flow_and_sequencing
```

Development

Testing

If you'd like to help development. Please use pre-commit before committing any files.

```
# runs pre-commit for syntax,formatting and static type checking
poetry run pre-commit run --all-files
```

Now you can run the tests via poetry

```
poetry run pytest -sv --log-cli-level DEBUG tests
```

Figures and Tables

Main Figures

Main figures as they are displayed in the G002 and G003 paper. Figure 1 is not included as it is manually generated.

```
g00x plots fig2
g00x plots fig3
g00x plots fig4
g00x plots fig5
g00x plots fig6
g00x plots fig7
```

Supplementary Figures

Supplementary figures are in order.

```
g00x plots supppfigures --all #TODO: Update later
```

Supplementary C mparis n Tables

```
g00x plots comparison-tables --all
g00x plots tables --all
```

Issues

Please submit any issues to the [issues page](#) and we are happy to help.

License

L PL

Copyright © Jordan R. Willis, TrySinc mb, Caleb Kibet, VISC, and IAVI