

GSV benchmarking - step by step guide

Overview

Genomic Sequence Variant (**GSV**) benchmarking evaluates how well different classifier settings distinguish closely related genomes.

This workflow clusters reference genomes from a target species into GSVs, simulates metagenomic samples from those references, and tests multiple Kraken2 and mSWEEP configurations **to maximise recall and precision**.

Off-target genomes from other species in the same family are processed in parallel to measure **false-positive rates**, allowing the pipeline to pinpoint the best combination of:

- Number of GSV clusters
- Kraken2 classifier confidence thresholds
- Count filtration to account for false positive classification

Table of Contents

- [Step 1: On-target genomes](#)
 - [1.1 Download all reference genomes](#)
 - [1.2 Run ANI on all genomes](#)
 - [1.3 Use hclust\(\) and cutree\(\) to create GSVs](#)
- [Step 2: Off-target genomes](#)
 - [2.1 Identify genomes for same family as target species](#)
 - [2.2 Remove accessions](#)
- [Step 3: Benchmarking](#)
 - [3.1 Preparing resources](#)
 - [3.1.1 Install VARIANT++ conda environment](#)
 - [3.1.2 Make themisto database](#)
 - [3.1.3 Confirm ANI clusters are in same order as themisto db](#)
 - [3.1.4 Download kraken database](#)
 - [3.2 Benchmarking with on target genomes \(recall/precision\) and off target genomes \(false positive rate\)](#)
 - [3.2.1 Make scripts](#)
 - [3.2.2 Run script to build simulated reads](#)
 - [3.2.3 Run script to run kraken and extract reads](#)
 - [3.2.4 Run script to classify GSVs in simulated samples](#)
- [Step 4: Analyzing benchmarking results](#)

Step 1: On-target genomes

1.1 Download all reference genomes

In progress

Example script to download genomes from accession list:

```
#!/bin/bash
#SBATCH -t 04:00:00 -n 1 --cpus-per-task=1 --mem=10G -o get_ncbi_genomes.out
module load WebProxy

datasets download genome accession --inputfile accessions_mh_2435.txt --include
genome --filename Mh_genbank.zip
```

1.2 Run ANI on all genomes

In progress

1.3 Use hclust() and cutree() to create GSVs

In progress

Step 2: Off-target genomes

1.1 Identify genomes for same family as target species

In progress

1.2 Remove accessions

In progress

Get a list of accessions for all genomes in the same family as your target genome, then use the python script "INSERT_HERE" on that metadata sheet to remove accessions for genomes from your target species.

Step 3: Bechmarking

What you'll need:

1. VARIANT++ github code
 - This includes the benchmarking code and anaconda recipe to install dependencies
2. ANI annotations with clusters from 2-30 (or more)
3. Target genomes
4. Off-target genomes
5. Themisto database
6. Kraken database

3.1 Preparing resources

3.1.1 Install VARIANT++ conda environment

We'll need to download the VARIANT++ github directory which will include all the necessary scripts and recipe for the conda environment.

I tend to use miniconda installed locally, but you can also just look for the latest anaconda module and load it directly. Once you have conda available, run this:

```
git clone https://github.com/Microbial-Ecology-Group/VARIANTplusplus.git

cd VARIANTplusplus/

conda env create -f envs/VARIANT++_env.yaml

conda activate VARIANT++_env
# if that doesn't work, use "source" instead of "conda" for the command above.
```

There's one software tool, themisto, that comes in the VARIANT++ repository, but unfortunately we have another step to make it available for our benchmarking analysis. We'll get the absolute path to the themisto executable you just downloaded and add that to your \$PATH so that it comes up automatically.

```
# Navigate to directory from the VARIANT++ dir
cd bin/

# Get path
pwd

# Add that directory to your $PATH variable
# Replace "/path/to/your/wd/" with your path
echo 'export PATH="/path/to/your/wd/:$PATH"' >> ~/.bashrc

# Now run the bash.rc file
source ~/.bashrc
```

Additionally, to make it easy to run the benchmarking scripts from anywhere on the server, we can add the absolute path to the "bin/Benchmarking_code" directory.

If you just did the step above, you should still be in the "bin/" directory. We'll move into the "Benchmarking_code" directory and repeat the steps.

```
# Navigate to directory from the VARIANT++ "bin/" dir
cd Benchmarking_code/

# Get path
pwd
```

```
# Add that directory to your $PATH variable
# Replace "/path/to/your/wd/" with your path
echo 'export PATH="/path/to/your/wd/:$PATH"' >> ~/.bashrc

# Now run the bash.rc file
source ~/.bashrc
```

At this point, navigate to wherever you want to run the benchmarking. We can run everything from this location.

3.1.2 Make themisto database

We'll use your target genomes to create a themisto database.

```
ls -d /scratch/user/u.ed124096/Mh_benchmarking/Mh_2435_genomes/*.fna >
2025_themisto_input_Mh_sequences.txt

themisto build -k 21 -i 2025_themisto_input_Mh_sequences.txt -o
2025_themisto_index_no --temp-dir tmp -t 2 --mem-gigas 4
```

3.1.3 Confirm ANI clusters are in same order as themisto db

Refer back to the file we made in R with all genomes in the first column and subsequent columns named "k_#". Open it with excel and make sure you sort all the rows so that they match the order of the genomes in the file you just made during the creation of the themisto database, in this case

[2025_themisto_input_Mh_sequences.txt](#). Save it with the fixed order, we'll be pointing to this below.

3.1.4 Download kraken database

The coreNT kraken database. It's about 240 GB and can be downloaded like this:

```
# Download database
wget https://genome-idc.s3.amazonaws.com/kraken/k2_core_nt_20241228.tar.gz

# Make directory for db contents
mkdir -p k2_core_nt_20241228

# unzip it
tar -xvzf k2_core_nt_20241228.tar.gz -C k2_core_nt_20241228/
```

3.2 Benchmarking with on target genomes (recall/precision) and off target genomes (false positive rate)

3.2.1 Make scripts

1. Modify the `params.txt` file to point to the corresponding directories to your genomes, the themisto database, ANI annotatinos, and kraken database.
 - Only fill out the variables as present in the template, don't add any quotes, spaces, or comments after the variables.
 - I recommend running a test with only a few simulations, low number of genomes per GSV, and maybe reduce the kraken confidence scores that's included. Once you can get through all steps, you can come back and modify the `params.txt` file again with the full simulation paramters.
2. Run the `generate_simulated_samples_byConf.py` script which will make three seperate scripts that we'll run in portions.

Command to run in terminal (don't need sbatch)

```
python generate_simulated_samples_byConf.py --benchmarking_params params.txt
```

This will output three scripts. They'll have different names depending on what you used for the `output_name_prefix` parameter in the `params.txt` file.

- `script_build_reads_test_benchmarking.sh`
- `script_kraken_extract_split_test_benchmarking.sh`
- `script_themisto_msweep_test_benchmarking.sh`

3.2.2 Run script to build simulated reads

The scripts we just made are way too long to run at once, so we'll use the script `split_script_iteration_markers.py`. This will break up the processes evenly across any given number of scripts. It will make a ".sh" command file and a corresponding ".sbatch" file that references the previous ".sh" script. So, each subset of commands will have a bash script, sbatch script, and a log file with similar names.

For each of the three main scripts, we'll change the parameters we provide the script, `split_script_iteration_markers.py`, which will change how the resulting sbatch scripts are made.

The main parameters of importance is the time `--time` and the memory `--mem` for each script that is created and this will change based on what tools we are running, how many subsets of scripts we made, and how long it takes to run your genomes.

In this python command, we'll point to the first script, "script_build_reads_test_benchmarking.sh", then specify the number of scripts to make, the prefix for the scripts that will be made.

```
python split_script_iteration_markers.py --input_file
script_build_reads_test_benchmarking.sh --num_output_files 20 --name_prefix
buildReads_byConf --mem 30G --time 24:00:00
```

Read the output text, you might notice if you used too high a number of some scripts have 0 lines or "0 marker lines". This means that there weren't enough code chunks to split into that many scripts. These scripts are now in a new directory, `scripts/`.

If, for example, only 13 of the 20 scripts had commands then we can use the following command to submit those scripts using sbatch. make sure to change the prefix to the left of the word "part" to match the prefix you used above.

Make sure you have the conda environment activated.

```
for i in {1..13}; do
  sbatch "scripts/buildReads_byConf_part_${i}.sbatch"
done
```

Waiting for scripts to finish running.

You'll now have to wait for those scripts to finish running. Remember to do a good testing run, so you have an idea of how long it will take per simulated. Maybe take chunks from one of the scripts and try running all the commands to see how long you'll need.

Once everything is done running, you can consider erasing the contents in the temporary directory specified by the `tmp_build` paramater, in this case "temp_conf". This can end up taking a lot of space so keep an eye on your storage space. Make sure you only delete the contents and not the directory itself. The next steps will also use that temporary directory.

You can do this at the end of each step, just make sure everything completed succesfully by looking at the end of the logs in the `scripts/` directory.

What if some of my scripts didn't completely finish?

If you look at the logs and see that they were cutoff due to time, you should be able to identify the last sample that was succesfully run by that script and erase everything up to that set of commands. Then you can resubmit, but try to calculate how much time would actually be needed before re-submitting.

3.2.3 Run script to run kraken and extract reads

This next section requires alot more resources because we'll be using kraken2 and the core nt database. Loading that large database requires at least 250G of memory. We could add "--memory-mapping" as a kraken2 flag to reduce the required memory, but it ends up taking way longer, so better to just request the full amount.

We'll modify the `split_script_iteration_markers.py` command again, to update the sbatch header so that it includes the 250G in memory. Remember to change the time based on your testing.

We'll run the script to split up the next series of commands from the second script, `script_kraken_extract_split_test_benchmarking.sh`.

```
python split_script_iteration_markers.py --input_file
script_kraken_extract_split_test_benchmarking.sh --num_output_files 20 --
name_prefix kraken_byConf --mem 250G --time 12:00:00
```

Submit them as before, remembering to update the prefix.

```
for i in {1..20}; do
  sbatch "scripts/kraken_byConf_part_${i}.sbatch"
done
```

Remember to check your temporary directory in case you want to empty it, but make sure you don't delete the folder itself.

3.2.4 Run script to classify GSVs in simulated samples

This final step will be similar, but requires much less memory. So update the command for `split_script_iteration_markers.py` again to look something like this, pointing to the last script, `script_themisto_msweep_test_benchmarking.sh`:

```
python split_script_iteration_markers.py --input_file
script_themisto_msweep_test_benchmarking.sh --num_output_files 20 --name_prefix
classify_byConf --mem 10G --time 12:00:00
```

Submit them as before, remembering to update the prefix.

```
for i in {1..20}; do
  sbatch "scripts/classify_byConf_part_${i}.sbatch"
done
```

Step 4: Analyzing benchmarking results

In progress

4.1 Parse msweep results without a filter

4.2 Analyze off target results and false positive rate in R

4.2.1 Pick filtration rates for each GSV

4.3 Evaluate classification performance

4.3.1 Pick optimal GSV number