



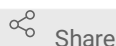
Jul 29, 2022

# BASIC PROTOCOL 4: Pan-genome Copy Number Variant Calling

In 1 collection

miriam.goldman<sup>1,2</sup>, chunyu.zhao<sup>3,4</sup><sup>1</sup>Data Science and Biotechnology, Gladstone Institutes, San Francisco, CA, USA,;<sup>2</sup>Biomedical Informatics, University of California San Francisco, San Francisco, CA;<sup>3</sup>Data Science, Chan Zuckerberg Biohub, San Francisco, CA, USA,;<sup>4</sup>Data Science and Biotechnology, Gladstone Institutes, San Francisco, CA, USA

1 Works for me



Share

[dx.doi.org/10.17504/protocols.io.n92ldzke7v5b/v1](https://dx.doi.org/10.17504/protocols.io.n92ldzke7v5b/v1)

miriam.goldman

## ABSTRACT

This protocol describes the CNV module of MIDAS2, which takes as input metagenomic sequencing reads from a set of samples and generates files with CNV genotypes for each sample for all detected species. There are two steps for population CNV calling: (1) single-sample quantification of copy number for each gene in the pangenome of each species with the midas2 run\_genes command and (2) population CNV calling with the midas2 merge\_genes command. Basic Protocols 1 (Species) and 2 (MIDASDB) should be run before this protocol.

## DOI

[dx.doi.org/10.17504/protocols.io.n92ldzke7v5b/v1](https://dx.doi.org/10.17504/protocols.io.n92ldzke7v5b/v1)

## PROTOCOL CITATION

miriam.goldman , chunyu.zhao 2022. BASIC PROTOCOL 4: Pan-genome Copy Number Variant Calling. **protocols.io**  
<https://dx.doi.org/10.17504/protocols.io.n92ldzke7v5b/v1>

## COLLECTIONS ⓘ

 **MIDAS 2 Protocol**

## LICENSE

\_\_\_\_\_ This is an open access protocol distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited

CREATED

Jul 28, 2022

LAST MODIFIED

Jul 29, 2022

PROTOCOL INTEGER ID

67835

PARENT PROTOCOLS

Part of collection

[MIDAS 2 Protocol](#)

- 1 Perform species prescreening as described in Basic Protocol 1.



BASIC PROTOCOL 1: Species Prescreening  
by **miriam.goldman**

PREVIEW

RUN



- 2 Download MIDASDB as described in Basic Protocol 2.



BASIC PROTOCOL 2: Download MIDAS Reference  
Database  
by **miriam.goldman**

PREVIEW

RUN



- 3 Execute the run\_genes command for each sample

1. Conceptually, a typical invocation of the run\_genes command proceeds by four steps:
2. Select the list of species abundant enough for accurate metagenotyping based on the species profiling results and user-defined species selection criterion. Taking SRR172902 as an example, run\_genes expect to find the species profiling results at midas2\_output/SRR172902/species/species\_profile.tsv.
3. Compile the pangenomes for these species and build a sample-customized pan-genome bowtie2 index.
4. Align reads to this index with bowtie2.
5. For each gene in the pan-genome, normalize gene coverage by the mean coverage of all that species' SCGs to estimate copy number per cell [11].
6. Output a read mapping summary and CNV estimates for each species.

```
midas2 run_genes
```

```
for sample_name in SRR172902 SRR172903  
do  
  midas2 run_genes \  
  --sample_name ${sample_name} \  
  -1 reads/${sample_name}.fastq.gz \  
  --midasdb_name uhgg --midasdb_dir midasdb_uhgg \  
  --species_list 100122,100277 \  
  --select_by median_marker_coverage,unique_fraction_covered \  
  --select_threshold=0,0.6 \  
  --num_cores 8 midas2_output  
done
```

The number of CPUs used is specified via `--num_cores 8`.

- 4 Prepare sample manifest file for merging purpose. We can use the same `list_of_samples.tsv` generated by step 6 in Basic Protocol 1.



BASIC PROTOCOL 1: Species Prescreening  
by miriam.goldman

PREVIEW

RUN



- 5 Upon the completion of `run_genes` for all the samples listed in the `list_of_samples.tsv`, MIDAS2 merges the CNV profiles across samples with the `merge_genes` command.

```
midas2 merge_genes --samples_list list_of_samples.tsv \  
--midasdb_name uhgg --midasdb_dir midasdb_uhgg \  
--min_copy 0.5 \  
--num_cores 2 midas2_output/merge
```

6 Population pangenome CNV analysis has finished successfully when all the following output files are created under the directory midas2\_output/merge/genes/ without any error message

- genes\_summary.tsv: merged single-sample CNV summary containing information such as mean\_coverage.
- <species\_id>/<species\_id>.genes\_copynum.tsv.lz4: gene-by-sample matrix of copy-number estimates
- <species\_id>/<species\_id>.genes\_preabs q.tsv.lz4: gene-by-sample matrix of gene presence / absence
- <species\_id>/<species\_id>.genes\_depth.tsv.lz4: gene-by-sample read coverage matrix
- <species\_id>/<species\_id>.genes\_reads.tsv.lz4: gene-by-sample read counts matrix