

# generate\_simulated\_samples\_byConf.py – Detailed Guide

---

*Automated generator for GSV-benchmarking bash scripts*

---

## 1. Purpose

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

This script 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

## 2. Quick Usage

Follow the GSV\_benchmarking\_tutorial.md document for detailed instructions. In short, you prepare a bunch of files and databases. Then you run the benchmarking script to take that information and create a series of 3 scripts which will make simulated samples and take them through our VARIANT++ classification workflow.

## 3. What the Script Generates

Bash Script	Stage	Key Tasks
<b>script_build_reads.sh</b>	Sample generation	<ul style="list-style-type: none"><li>• Concatenate target / off-target genomes</li><li>• Simulate reads with ISS</li><li>• Merge reads with FLASH</li><li>• Standardise filenames &amp; headers</li><li>• Produce one merged FASTQ per iteration</li></ul>
<b>script_kraken_extract_split.sh</b>	Taxonomic classification	<ul style="list-style-type: none"><li>• Classify merged FASTQs at each <b>--confidence</b></li><li>• Extract target reads with extract_kraken_reads.py</li><li>• Split merged vs unmerged read sets</li></ul>
<b>script_themisto_msweep.sh</b>	GSV quantification	<ul style="list-style-type: none"><li>• Pseudo-align extracted reads with Themisto</li><li>• Run mSWEEP (merged &amp; unmerged) to estimate GSV abundances</li></ul>

---

## 4. Workflow Breakdown

1. **Parse configuration** (**params.txt**): key/value file with paths, iteration counts, GSV list, Kraken confidences, etc.

2. Create directories:

```
<PREFIX>_cat_reads/      # merged FASTQs
<PREFIX>_split_reads/    # unmerged FASTQs
<PREFIX>_merged_reads/   # extracted merged reads
<PREFIX>_results/        # reports and outputs
```

- 3. **Generate mSWEEP annotation files** (`k_#.msweep.txt`) from the metadata TSV.
- 4. **Write `script_build_reads`** – loops over `num_GSV_list` × `num_iters` × `k_columns`.
  - `numGSV == 0` → simulate off-target samples.
  - `numGSV > 0` → select N GSVs, concatenate genomes per GSV, simulate reads.
- 5. **Write `script_kraken_extract_split`** – classify each merged FASTQ for every confidence value, then extract & split reads.
- 6. **Write `script_themisto_msweep`** – pseudo-align reads and run mSWEEP (all k-columns for off-target, single column for on-target).
- 7. **Finish** – print a summary of the three generated bash scripts.

5. Key Parameters (`params.txt`)

Key	Type	Example	Description
<code>target_genome_dir</code>	str	<code>/data/targets/</code>	<code>.fna[.gz]</code> of target species
<code>nontarget_genome_dir</code>	str	<code>/data/offtarget/</code>	Off-target family genomes
<code>input_file</code>	str	<code>ani_clusters.tsv</code>	TSV with genome filenames + GSV codes per <code>k_col</code>
<code>bin_dir</code>	str	<code>/usr/local/bin/</code>	Path to helper scripts ( <code>rename_*</code> , <code>split</code> scripts)
<code>themisto_index</code>	str	<code>/indices/themisto/</code>	Prefix of existing Themisto index
<code>krakendb</code>	str	<code>/databases/kraken2/</code>	Kraken2 database directory
<code>kraken_confidence</code>	list / int	<code>[0,0.1,0.5]</code>	Confidence thresholds to test
<code>num_iters</code>	int	<code>100</code>	Iterations per <code>numGSV</code> setting
<code>num_GSV_list</code>	list[int]	<code>[0,1,3,5,7,9]</code>	GSV counts to test
<code>num_reads_options</code>	list[int]	<code>[10000,20000]</code>	Possible read counts per sample
<code>threads</code>	int	<code>48</code>	CPU threads for all tools
<code>tmp_build</code>	str	<code>\$TMPDIR</code>	Scratch directory (can use environment variable)

6. Expected Outputs

```
<PREFIX>_cat_reads/      # merged FASTQs for Kraken2
<PREFIX>_split_reads/     # unmerged FASTQs for Themisto/mSWEEP
<PREFIX>_merged_reads/    # merged extracted reads by confidence
<PREFIX>_results/         # Kraken reports, Themisto outputs, mSWEEP tables
script_build_reads_<PREFIX>.sh
script_kraken_extract_split_<PREFIX>.sh
script_themisto_msweep_<PREFIX>.sh
```

---

## 7. Dependencies

- Python  $\geq 3.6$
- **ISS** (read simulator)
- **FLASH** (read merger)
- **Kraken2**  $\geq 2.1$
- **extract\_kraken\_reads.py** (KrakenTools)
- **Themisto**  $\geq 2.5$
- **mSWEEP**  $\geq 1.9$
- Helper scripts in `bin_dir` (`rename_sample_files.py`, `rename_headers.py`, `split_extracted_reads_by_conf.py`)

---

## 8. Troubleshooting Tips

- **Empty output?** Verify `num_GSV_list` values exist in your TSV.
  - **File-not-found errors:** Check all paths in `params.txt`.
  - **Kraken2 memory issues:** Add `--memory-mapping` to reduce memory, but request a lot more time.
  - **mSWEEP annotation errors:** Ensure each `k_#.msweep.txt` has one entry per genome.
-