

Annotate your genome with PGAP

Requirements

To run the PGAP pipeline you will need:

- Python (version 3.6 or higher),
- the ability to run Docker (see <https://docs.docker.com/install/> if it is not already installed), Singularity, or Podman
- about 100GB of storage for the supplemental data and working space,
- and 2GB-4GB of memory available per CPU used by your container.
- The CPU must have SSE 4.2 support (released in 2008).

Quick Start

Download the file using either

```
$ curl -OL https://github.com/ncbi/pgap/raw/prod/scripts/pgap.py
```

or

```
$ wget -O pgap.py https://github.com/ncbi/pgap/raw/prod/scripts/pgap.py
```

depending upon which utility your system has installed. If one does not work, try the other.

Install the pipeline. By default it will install in `$HOME/.pgap`, but this location can be changed by setting environmental variable `PGAP_INPUT_DIR`.

```
$ chmod +x pgap.py
$ ./pgap.py --update # required files are downloaded and ex
```

```
tracted
```

Run the pipeline on the *Mycoplasma genitalium* genome provided with the installation:

```
$ ./pgap.py -r -o mg37_results -g $HOME/.pgap/test_genomes/MG37/ASM2732v1.annotation.nucleotide.1.fasta -s 'Mycoplasma genitalium'
```

Output will be located in the `mg37_results` subdirectory as specified by the `-o` flag.

Bring Your Own Data

To run this pipeline using your own genomes, you will need, at a minimum, the multifasta file for the genome, and the associated organism name (genus or genus species).

```
$ ./pgap.py -r -o <results> -g <fasta> -s '<organism_name>'
```

After successful completion the output directory will contain the following files:

- `.fasta` - nucleotide FASTA file that you supplied
- `ani-tax-report.txt` - ANI report in text format
- `ani-tax-report.xml` - ANI report in XML format for machine processing
- `annot-gb.ent` - ASN.1 file in Seq-entry genbank format
- `annot.faa` - FASTA file with all proteins
- `annot.fna` - FASTA file with all nucleotides (note the file might be slightly normalized compared to your input nucleotide FASTA file)
- `annot.gbk` - annotations in flatfile format
- `annot.gff` - annotations in GFF3 format
- `annot.sqn` - ASN.1 file in Seq-submit format

- annot_cds_from_genomic.fna - FASTA file with nucleotide sequences of all coding regions
- annot_translated_cds.faa - FASTA file with translated sequences of all coding regions
- annot_with_genomic_fasta.gff - file combining annotations in GFF format and nucleotide sequence in FASTA format used in some third party applications, like Roary
- checkm.txt - CheckM output for this genome
- cwltool.log - CWL tool log that could be instrumental for post mortem analysis of failures
- fastaval.xml - XML file with validation results for input FASTA file

Useful options

Command	Description
<code>-g <path>, --genome <path></code>	Path to genomic fasta
<code>-s 'organism', --organism 'organism'</code>	Genus, or genus species
<code>-r, --report-usage-true</code>	Report anonymized usage metadata to NCBI
<code>-n, --report-usage-false</code>	Do not report anonymized usage metadata to NCBI
<code>-o <path>, --output <path></code>	Output directory to be created, which may include a full path
<code>--ignore-all-errors</code>	Ignore errors from quality control analysis, in order to obtain a draft annotation
<code>--no-internet</code>	Disable internet access for all programs in pipeline
<code>-D <path>, --docker <path></code>	Docker-compatible executable (e.g. docker, podman, singularity), which may include a full path like /usr/bin/docker
<code>--taxcheck</code>	Also calculate the Average Nucleotide Identity to type assemblies
<code>--taxcheck-only</code>	Only calculate the Average Nucleotide Identity to type assemblies, do not run

Command	Description
	PGAP
<code>--auto-correct-tax</code>	Override the organism provided in the input YAML file, if the taxcheck predicts a different organism with high confidence. Use in combination with the <code>--taxcheck</code> flag
<code>-d, --debug</code>	Debug mode. Retain intermediate files needed for <u>investigating failures</u>