

Advanced Molecular Detection Southeast Region Bioinformatics

Outline



Updates



Agenda



Bactopia



FastANI



GAMMA

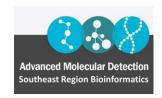


Questions

Updates – ABiL Trainings

ABiL courses

- Online courses will be available shortly for the attendees
 - Pathogen Phylogenomics
 - Quality Assessment of Sequencing Data
- If additional attendees decide to sign up later, that is not a problem, as the courses will still be available, and they can attend on a rolling basis under the contract.

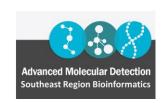


Agenda

September 4 rescheduled to September 11 – Bactopia Tools: HICAP and HpsuisSero **September 18** – Bactopia Tools: Kleborate and Legsta

Future Trainings

- ONT & FL's Flisochar pipeline
- StaPH-B Toolkit Programs/Pipelines
- GISAID flagged SARS-CoV-2
- R Training Series
- Dryad pipeline
- ...and more



FastANI

- Fast Whole-Genome Similarity (ANI) Estimation
- ANI is defined as mean nucleotide identity of orthologous gene pairs shared between two microbial genomes
- FastANI supports pairwise comparison of both complete and draft genome assemblies

ParBLiSS/FastANI: Fast Whole-Genome Similarity (ANI) Estimation (github.com)



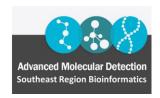
Installation

Available as a module on HPG

module load fastani/1.1

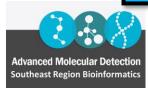
Can be installed through conda

conda create –yp /blue/bphl-<state>/<user>/conda_envs/fastani/ conda activate /blue/bphl-<state>/<user>/conda_envs/fastani/ conda install –c conda-forge –c bioconda fastani



Usage

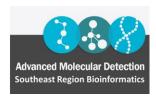
```
thsalikilakshmi@login1:/blue/bphl-florida/thsalikilakshmi/data/HAI/20220727_jax_220708_PLN_WLK_MS_test/assemblies —
[thsalikilakshmi@loginl assemblies]$ fastANI -h
fastANI is a fast alignment-free implementation for computing whole-genome
Average Nucleotide Identity (ANI) between genomes
Example usage:
$ fastANI -q genomel.fa -r genome2.fa -o output.txt
$ fastANI -q genomel.fa --rl genome list.txt -o output.txt
Available options
-h, --help
    Print this help page
-r <value>, --ref <value>
    reference genome (fasta/fastq)[.gz]
--refList <value>, --rl <value>
    a file containing list of reference genome files, one genome per line
-q <value>, --query <value>
    query genome (fasta/fastq)[.gz]
--ql <value>, --queryList <value>
    a file containing list of query genome files, one genome per line
-k <value>, --kmer <value>
    kmer size <= 16 [default : 16]
```



Input

- To compute ANI, use a query genome and a reference genome
- Here we computed ANI between *Escherichia coli* and *Shigella flexneri* genomes
- *E. coli* is provided as a query genome, *S. flexneri* is the reference genome and –o flag is used to give fastani.out as output file

```
[thsalikilakshmi@login2 assemblies]$ fastANI -q /blue/bphl-florida/thsalikilakshmi/data/HAI/20
220727_jax_220708_PLN_WLK_MS_test/assemblies/JBE22000247.fasta -r /blue/bphl-florida/thsalikil
akshmi/data/jbi/20220825_jax_220629_PLN_WLK_MS/assemblies/JBI22000647.fasta -o fastani.out > r
esults_fastani
```



Output

```
[thsalikilakshmi@login2 assemblies]$ fastANI -q /blue/bphl-florida/thsalikilakshmi/data/HAI/20
220727 jax 220708 PLN WLK MS test/assemblies/JBE22000247.fasta -r /blue/bphl-florida/thsalikil
akshmi/data/jbi/20220825 jax 220629 PLN WLK MS/assemblies/JBI22000647.fasta -o fastani.out > r
esults fastani
>>>>>>>
Reference = [/blue/bphl-florida/thsalikilakshmi/data/jbi/20220825 jax 220629 PLN WLK MS/assemb
lies/JBI22000647.fastal
Query = [/blue/bphl-florida/thsalikilakshmi/data/HAI/20220727 jax 220708 PLN WLK MS test/assem
blies/JBE22000247.fastal
Kmer size = 16
Fragment length = 3000
Threads = 1
ANI output file = fastani.out
>>>>>>
INFO [thread 0], skch::Sketch::build, minimizers picked from reference = 346824
INFO [thread 0], skch::Sketch::index, unique minimizers = 343053
INFO [thread 0], skch::Sketch::computeFreqHist, Frequency histogram of minimizers = (1, 340077
 ... (42, 1)
INFO [thread 0], skch::Sketch::computeFregHist, With threshold 0.001%, ignore minimizers occur
ring >= 29 times during lookup.
INFO [thread 0], skch::main, Time spent sketching the reference : 0.317155 sec
INFO [thread 0], skch::main, Time spent mapping fragments in query #1 : 0.81592 sec
INFO [thread 0], skch::main, Time spent post mapping: 0.000211501 sec
```



Result

- Below output implies that the ANI estimate between *E. coli* and *S. flexneri* genomes is 97.4216
- Out of the total 1595 sequence fragments from *E. coli*, 1154 were aligned as orthologous matches

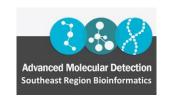
JBI22000647.fasta 97.4216 1154 1595



GAMMA

- Gene Allele Mutation Microbial Assessment
- GAMMA is a command line tool that finds gene matches in microbial genomic data using protein coding (rather than nucleotide) identity, and then translates and annotates the match by providing the type (i.e., mutant, truncation, etc.) and a translated description (i.e., Y190S mutant, truncation at residue 110, etc.)
- GAMMA is helpful in both identifying and explaining how unique alleles differ from their closest known matches

rastanton/GAMMA: Gene Allele Mutation Microbial Assessment (github.com)



Installation

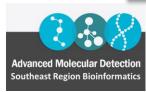
Can be installed through conda

conda create –yp /blue/bphl-<state>/<user>/conda_envs/gamma/ conda activate /blue/bphl-<state>/<user>/conda_envs/gamma/ conda install –c conda-forge –c bioconda gamma



Usage

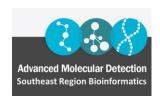
```
# thsalikilakshmi@login2:/blue/bphl-florida/thsalikilakshmi/data/HAI/20220727_jax_220708_PLN_WLK_MS_test/ass...
Preparing transaction: done
Verifying transaction: done
Executing transaction: done
(/blue/bphl-florida/thsalikilakshmi/training/conda envs/gamma) [thsalikilakshmi@login2 assembl
ies]$ GAMMA.py --help
usage: GAMMA.py [-h] [-a] [-e] [-f] [-g] [-n] [-l] [-i PERCENT IDENTITY]
                input fasta database output
This scripts makes annotated gene calls from matches in an assembly using a gene database
positional arguments:
  input fasta
                         input fasta
  database
                         input database
  output
                         output name
options:
                        show this help message and exit
  -h, --help
  -a, --all
                         include all gene matches, even overlaps
                         writes out all protein mutations
  -e, --extended
                         write fasta of gene matches
  -f, --fasta
  -g, --gff
                         write gene matches as gff file
  -n, --name
                        writes name in front of each gene match line
  -1, --headless
                         removes the header from the output gamma file
  -i PERCENT IDENTITY, --percent identity PERCENT IDENTITY
                         minimum nucleotide identity for blat search (default = 90)
(/blue/bphl-florida/thsalikilakshmi/training/conda envs/gamma) [thsalikilakshmi@login2 assembl
ies]$
```



Input

GAMMA.py my_genome.fasta gene_db.fasta output_name [optional arguments]

- Input for GAMMA is a genome or assembly in .fasta format and a multifasta database of the coding sequences of genes
- GAMMA was tested using
 - AR gene databases from AMRFINDERPLUS
 - (Index of /pathogen/Antimicrobial_resistance/AMRFinderPlus/database (nih.gov))
 - ARG_ANNOT
 - (backup.mediterraneeinfection.com/arkotheque/client/ihumed/_depot_arko/articles/2041/arg-annot-v4-aamay2018_doc.fasta)
 - RESFINDER
 - (genomicepidemiology / resfinder_db Bitbucket)

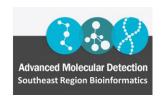


Input

The sample GAMMA input shown below was generated from running GAMMA on a *S. flexneri* using a combination of all the ResFinder databases

genomicepidemiology / resfinder db — Bitbucket

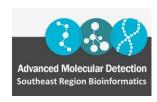
```
(/blue/bphl-florida/thsalikilakshmi/training/conda_envs/gamma) [thsalikilakshmi@login2 assembl ies]$ GAMMA.py /blue/bphl-florida/thsalikilakshmi/data/jbi/20220825_jax_220629_PLN_WLK_MS/asse mblies/JBI22000647.fasta /blue/bphl-florida/thsalikilakshmi/data/HAI/20220727_jax_220708_PLN_WLK_MS_test/assemblies/resfinder_database/resfinder_db/all.fsa output
```



Results

The default output of GAMMA is a tab-delimited file with a .gamma extension with 15 columns

1	Α	В	С	D	Е	F	G	Н	1	J	K	L	М	N	0
1	Gene	Contig	Start	Stop	Match_Ty	Descriptio	Codon_Ch	BP_Chang	Transvers	Codon_Pe	BP_Percer	Percent_L	Match_Lei	Target_Le	Strand
2	catA1_1_V	169	5200	5860	Native	No coding	0	1	0	1	0.9985	1	660	660	-
3	blaTEM-18	191	74	935	Native	No coding	0	0	0	1	1	1	861	861	+
4	sul2_2_AY	191	3355	4171	Native	No coding	0	0	0	1	1	1	816	816	-
5	aph(6)-Id_	191	1655	2492	Native	No coding	0	0	0	1	1	1	837	837	-
6	aph(3")-Ib	191	2491	3295	Native	No coding	0	0	0	1	1	1	804	804	-
7	tet(B)_2_/	205	807	2013	Native	No coding	0	0	0	1	1	1	1206	1206	-
8	mph(A)_1	214	2257	3163	Native	No coding	0	0	0	1	1	1	906	906	+
9	blaOXA-1	244	984	1815	Native	No coding	0	0	0	1	1	1	831	831	-
10	aadA1_3	244	80	872	Mutant	A4V,	1	2	1	0.9962	0.9975	1	792	792	-
11	dfrA14_1_	294	445	919	Native	No coding	0	0	0	1	1	1	474	474	-





Advanced Molecular Detection Southeast Region Bioinformatics

Questions?

bphl-sebioinformatics@flhealth.gov

Lakshmi Thsaliki, MS

Bioinformatician

Lakshmi.Thsaliki@flhealth.gov

Molly Mitchell, PhD

Bioinformatician

Molly.Mitchell@flhealth.gov