

Appendix 2 *ancestryinfer* pipeline

Output of *mixnmatch* simulations *or* your own data can be input into the *ancestryinfer* pipeline to run local ancestry inference following (Corbett-Detig & Nielsen, 2017).

Install

Option 1 – install dependencies:

```
git clone https://github.com/Schumerlab/ancestryinfer.git
```

To install dependencies, follow instructions outlined in:

```
installation_instructions.txt
```

Test that the install and pipeline are working:

```
cd ancestryinfer  
  
perl Ancestry_HMM_parallel_v5.pl  
hmm_configuration_file_nonparallel.cfg
```

Option 2 – load docker file for dependencies:

```
docker pull schumer/mixnmatch-ancestryinfer-image:mixnmatch-ancestryinfer-docker
```

```
docker run -it mixnmatch-ancestryinfer-image bash
```

Test that the install and pipeline are working:

```
cd ancestryinfer  
  
perl Ancestry_HMM_parallel_v5.pl  
hmm_configuration_file_nonparallel.cfg
```

Setting parameters in the configuration file

There are example configuration files available on github:

```
hmm_configuration_file_parallel.cfg  
hmm_configuration_file_nonparallel.cfg
```

Parameter descriptions:

Parameter	Description	Example	Include if
<code>genome1=</code>	User provided fasta file for species 1	<code>genome1=xiphophorus_birchmanni_10x_12Sep2018_yDAA6.fasta</code>	Always
<code>genome2=</code>	User provided fasta file for species 2	<code>genome2=Xmalinche_dovetail_assembly.fa</code>	Always
<code>read_type=</code>	Indicated whether data is paired end or single end	<code>read_type=PE</code>	Always
<code>read_list=</code>	Provide list (including full paths) to the reads to be analyzed	<code>read_list=combined_all_call_hybrids_read_list</code> example list format for paired end data (single end file should contain one line per individual): <code>./reads/CALL1_read1.fq.gz ./reads/CALL1_read2.fq.gz</code> <code>./reads/CALL2_read1.fq.gz ./reads/CALL2_read2.fq.gz</code> <code>./reads/CALL3_read1.fq.gz ./reads/CALL3_read2.fq.gz</code>	Always
<code>read_length=</code>	Provide expected read length	<code>read_length=150</code>	Always
<code>prop_genome_genome1_parent=</code>	Expected proportion of the genome derived from the parent species listed under genome1	<code>prop_genome_genome1_parent=0.5</code>	If not provided, AncestryHMM will attempt to estimate (may increase run time)
<code>number_indiv_per_job=</code>	Parallelize jobs such that each job processes this number of individuals. Low numbers mean high parallelization and high number mean low parallelization.	<code>number_indiv_per_job=1</code>	Always
<code>program_path=</code>	Path to the program install folder	<code>program_path=/home/groups/schumer/shared_bin/Ancestry_HMM_pipeline</code>	If not provided the program

			will assume necessary scripts and programs are in the working directory
<code>provide_AIMs=</code>	Coordinates and identities of ancestry informative sites that distinguish the two parent species	<code>provide_AIMs=Xbirmanni10xgenome_ancestry_informative_sites_filterF1</code> Example list format: <pre> ScyDAA6-2-HRSCAF-26 58345 T C ScyDAA6-2-HRSCAF-26 58976 T A ScyDAA6-2-HRSCAF-26 59896 T C ScyDAA6-2-HRSCAF-26 60164 G A ScyDAA6-2-HRSCAF-26 63105 G A ScyDAA6-2-HRSCAF-26 65532 G A ScyDAA6-2-HRSCAF-26 66290 C A ScyDAA6-2-HRSCAF-26 68233 T C ScyDAA6-2-HRSCAF-26 70398 G A ScyDAA6-2-HRSCAF-26 73869 G A </pre>	Required unless provided genomes are on the same coordinate system and can be auto detected
<code>provide_counts=</code>	Counts of parental allele frequencies at ancestry informative sites (and recombination rates between adjacent sites if available)	<code>provide_counts=Xbirmanni10xgenome_Xmalinche_observed_parental_counts_filterF1</code> Example format: <pre> ScyDAA6-2-HRSCAF-26 163722 129 3 0 54 0.00000078 ScyDAA6-2-HRSCAF-26 166158 135 5 0 54 0.00001374 ScyDAA6-2-HRSCAF-26 166535 6 0 0 6 0.00000754 </pre> Columns are: <pre> Chromosome site allele1_count_parent1 allele2_count_parent1 allele1_count_parent2 allele2_count_parent2 recombination_rate </pre>	If not provided, the program will assume that provided ancestry informative sites are fixed between species
<code>per_site_error=</code>	Per-site error parameter for HMM (i.e. due to sequencing error, contamination, etc)	<code>per_site_error=0.02</code>	Always
<code>gen_initial_admix=</code>	Estimated generation of initial admixture	<code>gen_initial_admix=20</code>	If not provided, AncestryHMM will attempt to estimate (may increase run time)

<code>focal_chrom_list=</code>	Provide a list of chromosomes to run (other chromosomes will not be run)	<code>focal_chrom_list=mychrs.txt</code> Example: <code>ScyDAA6-2-HRSCAF-26</code> <code>ScyDAA6-7-HRSCAF-50</code>	Not required
<code>rec_M_per_bp=</code>	Estimated recombination rate in Morgans/bp	<code>rec_M_per_bp=0.00000002</code>	Always; Use an estimate for a related species if not available
<code>max_alignments=</code>	Limit analysis to a maximum number of alignments (for computational speed)	<code>max_alignments=2000000</code>	Optional
<code>retain_intermediate_files=</code>	Keep all intermediate files. Warning: setting this to 1 results in a high space footprint for a large run; only recommended for troubleshooting.	<code>retain_intermediate_files=0</code>	Options are 1 to keep or 0 to delete.
<code>posterior_thresh=</code>	Posterior probability threshold to use for identifying ancestry transition intervals	<code>posterior_thresh=0.9</code>	Recommended 0.8-1
<code>job_submit_command=</code>	Option to run sequentially if using Docker image for dependencies or from a desktop computer. Set bash to run sequentially and sbatch to run in	<code>job_submit_command=bash</code> or <code>job_submit_command=sbatch</code>	Always required

	parallel on a slurm cluter		
<code>slurm_command_map=</code> <code>slurm_command_variant_call=</code> <code>slurm_command_hmm=</code>	If running on a slurm cluster, provide cluster specific parameters for queues, time & memory	<code>slurm_command_map=#!/bin/sh #SBATCH --ntasks=1 #SBATCH --cpus-per-task=1 #SBATCH -p schumer --mem=64000 #SBATCH --time=02:30:00</code> <code>slurm_command_variant_call=#!/bin/sh #SBATCH --ntasks=1 #SBATCH --cpus-per-task=1 #SBATCH -p schumer --mem=64000 #SBATCH --time=05:00:00</code> <code>slurm_command_hmm=#!/bin/sh #SBATCH --ntasks=1 #SBATCH --cpus-per-task=1 #SBATCH -p schumer --mem=64000 #SBATCH --time=03:00:00</code>	Required if running on a cluster

Examples

Several example files are available with the git repository including example configuration files

Running the pipeline

After setting the parameters in the configuration file and loading required dependencies, simply run:

```
perl mixnmatch/ simulate_admixed_genomes_v6.pl
hybrid_simulation_configuration.cfg
```

where path is the path to your simulator install