File	Description	Links	Programs commonly
extension			used in/produced by
.fa	Fasta files are files that contain a nucleotide or protein sequence, generally without quality scores. This will typically be a high confidence nucleotide sequence or reference assembly. The name of each sequence included is preceded by the > symbol, followed by the sequence.	https://zhanglab.ccmb .med.umich.edu/FAS TA/	Many!
.fq/.fastq (also .fq.gz and .fastq.gz)	Read data from next- generation sequencing. There are four lines per read, including information such as the read name, the nucleotide sequences read and relevant quality scores	https://support.illumi na.com/bulletins/201 6/04/fastq-files- explained.html	Mapping programs, de novo assembly programs, among other
.sam/.bam	Format for reads mapped to a reference genome including important information such as the read, mapping location, and mapping quality. The bam file is binary while the sam file is uncompressed.	https://genome.sph.u mich.edu/wiki/SAM	bwa, STAR, stampy, and others
.vcf/.bcf	Variant file containing information on variant base pairs, mapping quality, depth, and a large number of base quality metrics relative to the reference genome used for mapping. The vcf file is uncompressed which the bcf file is the binary version. Note: the specifics of the vcf format will vary from program to program	https://www.ebi.ac.uk /training/online/cours e/human-genetic- variation- introduction/exercise- title/want-know-how- we-did-it	samtools/beftools, GATK, and others

C	G: :1		. 1 /1 () 1
.g.vcf	Similar to a vcf but also		samtools/bcftools,
	containing information		GATK, and others
	for invariant sites. Note:		
	the specifics of the vcf		
	format will vary from		
	program to program		
ped	Compact file format	http://www.gwaspi.or	plink
map	containing pedigree	g/?page id=145	F
шир	(optional), phenotype	g/:puge_id 113	
	(optional), and genotype		
	data that is used by the		
	plink program and		
	several other programs.		
	The ped file must be		
	paired with a map file		
	which gives the locations		
	of SNPs		
bed	bed files are useful file	https://www.genomat	Can be included to
	formats for storing	ix.de/online help/hel	incorporate known
	interval information for a	p regionminer/bedfor	features (e.g. repetitive
	particular genome	mat help.html	regions, exons) in many
	assembly. The first	mat_nerpt	programs
	column contains the		programs
	linkage		
	_		
	group/scaffold/chromoso		
	me name, the second		
	column contains the		
	feature start position, the		
	third column contains the		
	feature end position.		
	Additional columns after		
	these first three are		
	allowed by most		
	programs		
gff/gtf	Annotation files that	http://genomewiki.uc	Can be included as input
	contain information on	sc.edu/index.php/Gen	with many programs,
	specific features in the	es_in_gtf_or_gff_for	usually will be provided
	genome. A gtf is a	mat	by a database (i.e.
	special case of a gff file	1111111	ensembl, ncbi, etc)
	that contains information		chischiot, neot, etc)
	on gene structure (i.e.		
ingur	exon locations).		gagtle
insnp	A file in the format:		seqtk
	chromosome, position,		
	basepair current, basepair		
	to change to for use with		

the seqtk program for	
updating fasta files	

Lab-specific file formats:

Ancestry-tsv format: posterior probabilities of ancestry for each individual and ancestry informative site in the genome

genotype format: hard-calls for ancestry in the same format as Ancestry-tsv files (2: homozygous parent 2 - malinche, 1: heterozygous ancestry, 0: homozygous parent 1 - birchmanni). We typically use a posterior probability cutoff of 0.9.

hybrid index file format: file including the estimated proportion of the genome from *malinche* and estimated ancestry heterozygosity