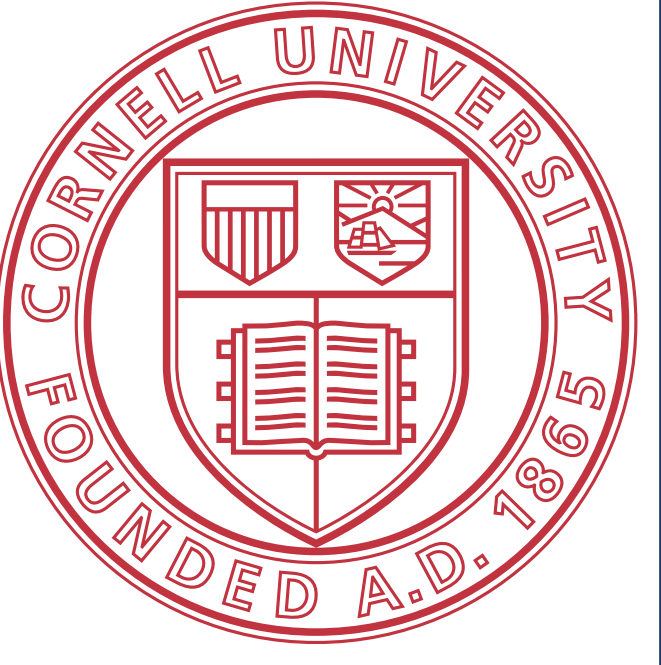




# Practical Haplotype Graph (PHG) to call genotypes from skim sequences to aid in genomic selection



Lynn Johnson<sup>1</sup>, Dan C. Ilut<sup>1</sup>, Zack Miller<sup>1</sup>, Terry M. Casstevens<sup>1</sup>, Peter J. Bradbury<sup>1,2</sup>, Punna Ramu<sup>1</sup>, Cinta M. Romay<sup>1</sup>, Edward S. Buckler<sup>1,2,\*</sup>

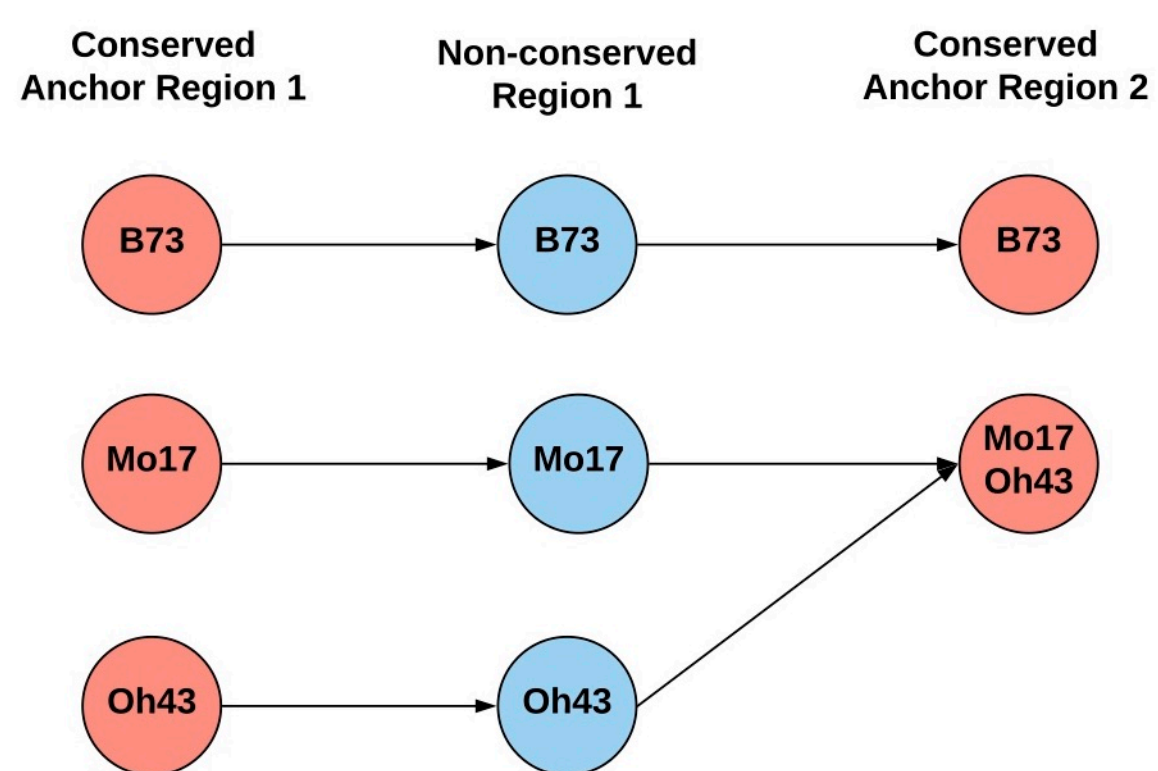
<sup>1</sup>Institute of Genomic Diversity, Cornell University, Ithaca, NY, USA. <sup>2</sup>US Department of Agriculture – Agriculture Research Service (USDA-ARS). \* Correspondence should be addressed to E.S.B. ([esb33@cornell.edu](mailto:esb33@cornell.edu))

## Introduction: Why a graph?

Biology Produces a consistent pattern of a genome

- Conserved genes (and other elements)
- Non-conserved intergenic regions of tremendous variation
- Architecture similar across many species

Maize Example:



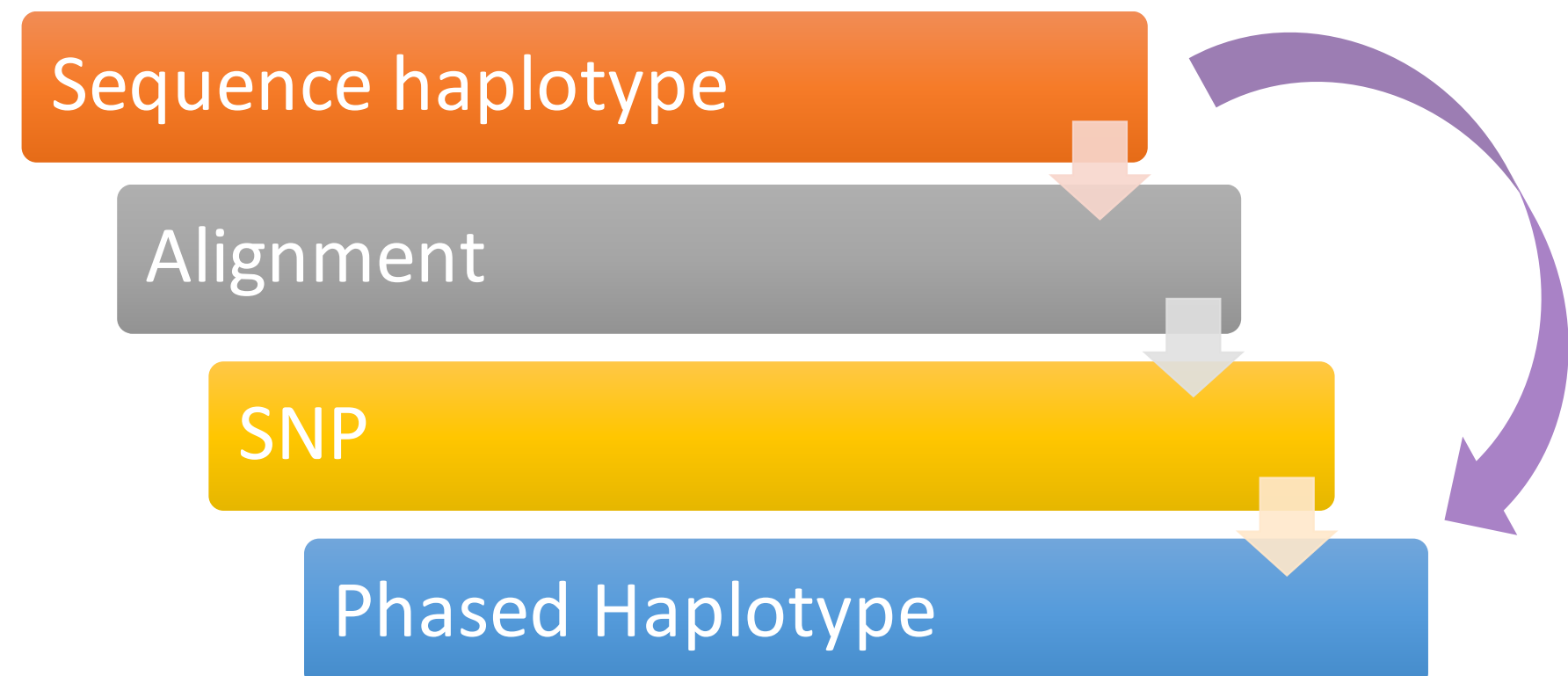
- Representing genomic diversity and complexity is challenging
- Graphs can compactly represent sequence from multiple genomes
- Aligning to a pan-genome is better than aligning to a single reference

## Why the PHG?

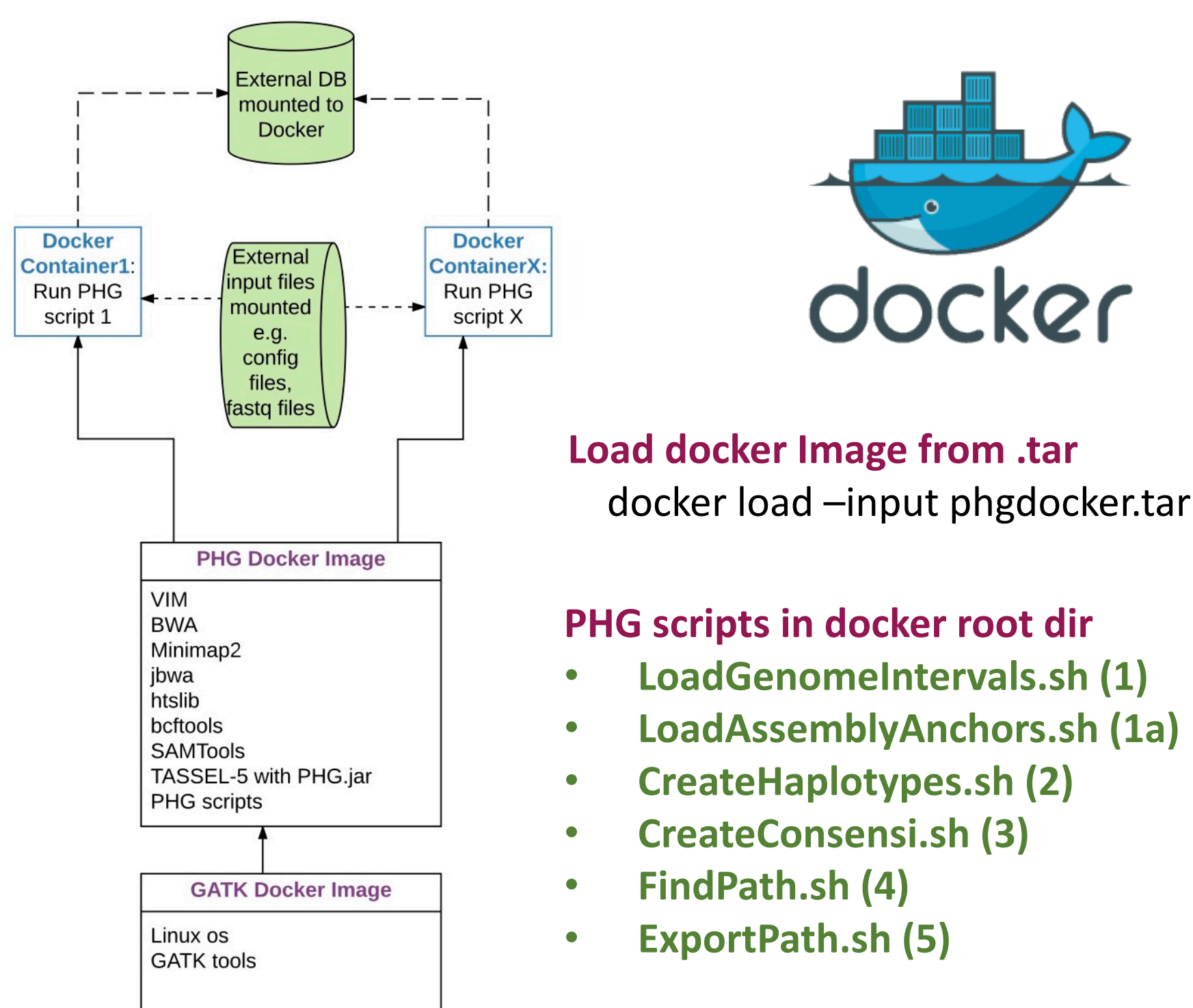
- PHG accepts sequence from multiple technologies (rAmpSeq, Nextera, Nanopore, etc.)
- By pulling taxa into consensus haplotypes, the haplotype graph can be built with low coverage input genomes
- We don't have good assemblies, intergenic regions are horrible
- PHG creates a useful graph even when the data isn't perfect

## PHG Goals

- ❖ Identify the haplotype from low depth sequences
- ❖ Create custom genomes for alignment
- ❖ Call rare haplotypes
- ❖ Compress data



## PHG docker instance



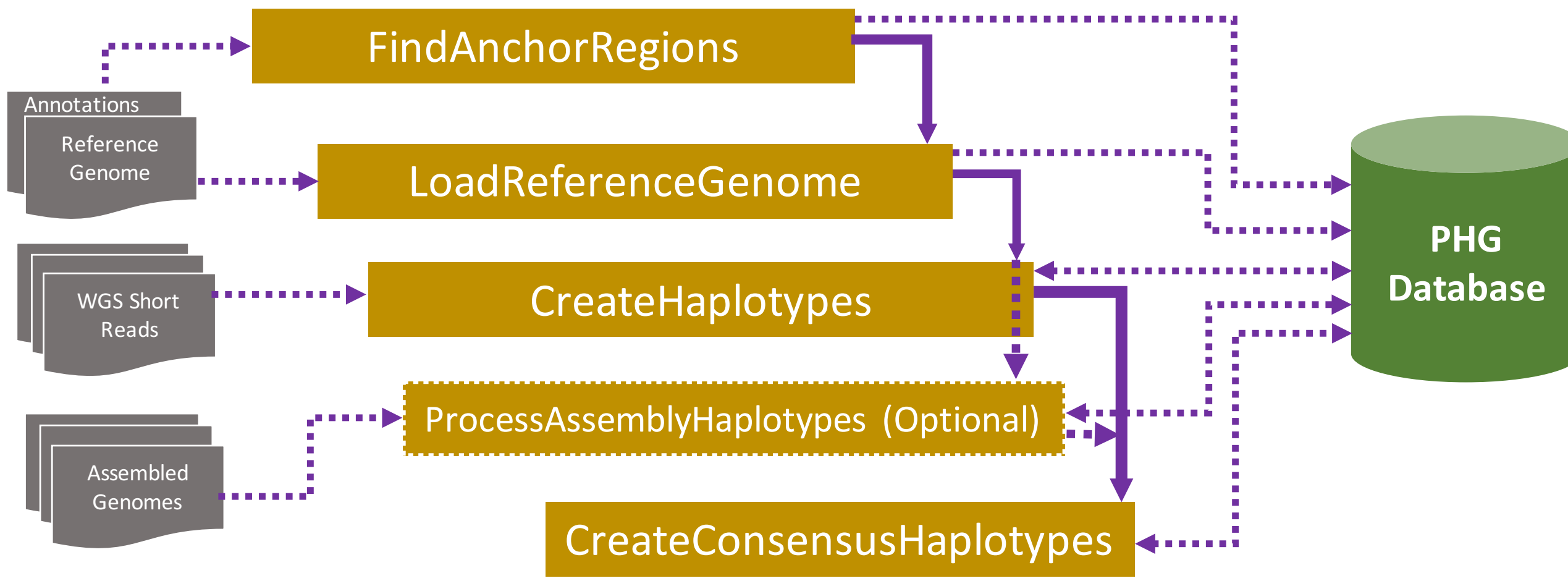
Load docker Image from .tar  
`docker load -input phgdocker.tar`

PHG scripts in docker root dir

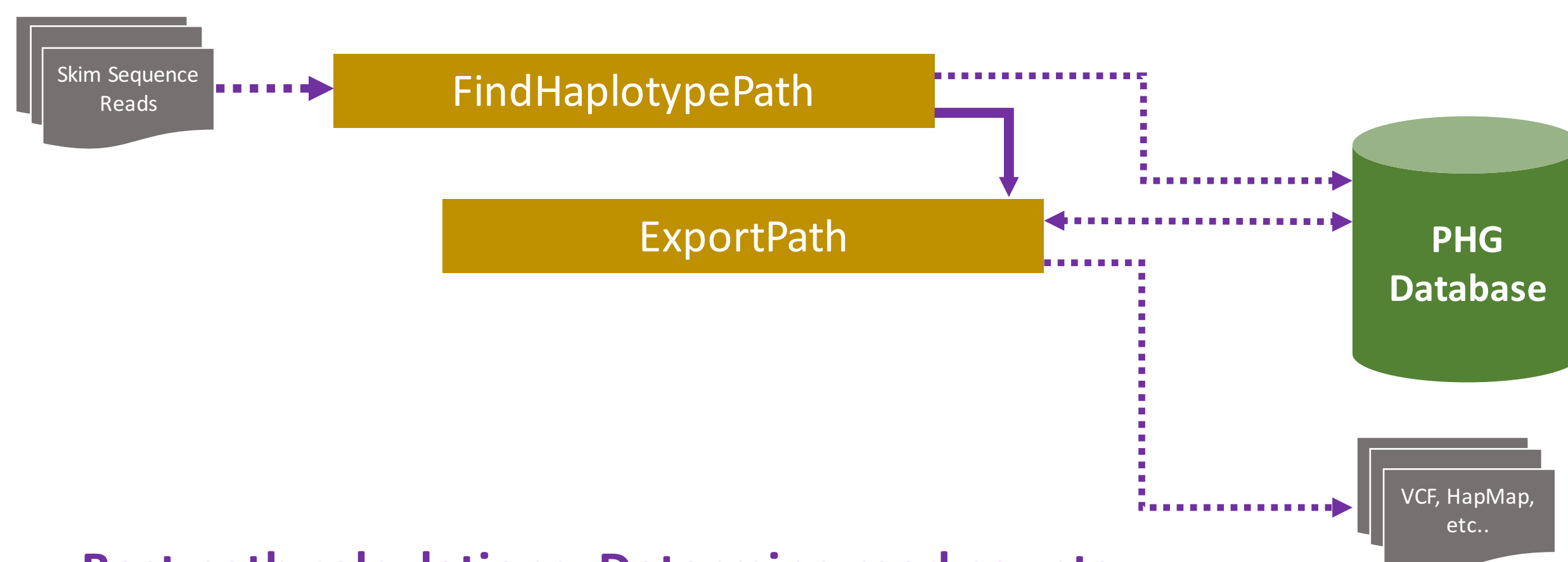
- `LoadGenomeIntervals.sh (1)`
- `LoadAssemblyAnchors.sh (1a)`
- `CreateHaplotypes.sh (2)`
- `CreateConsensi.sh (3)`
- `FindPath.sh (4)`
- `ExportPath.sh (5)`

## Calling genotypes from skim sequences using the PHG

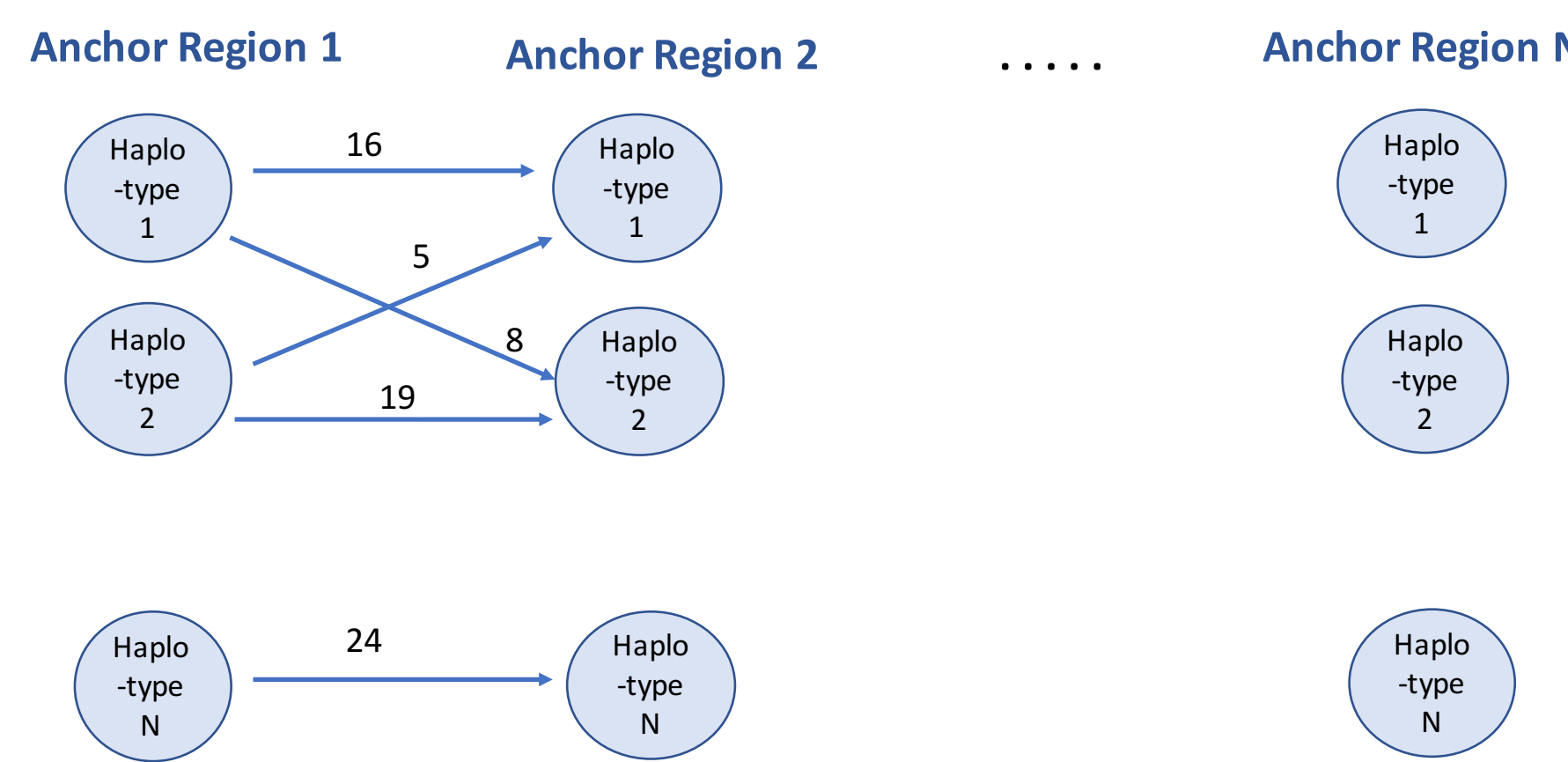
### 1. Populate the database with haplotypes



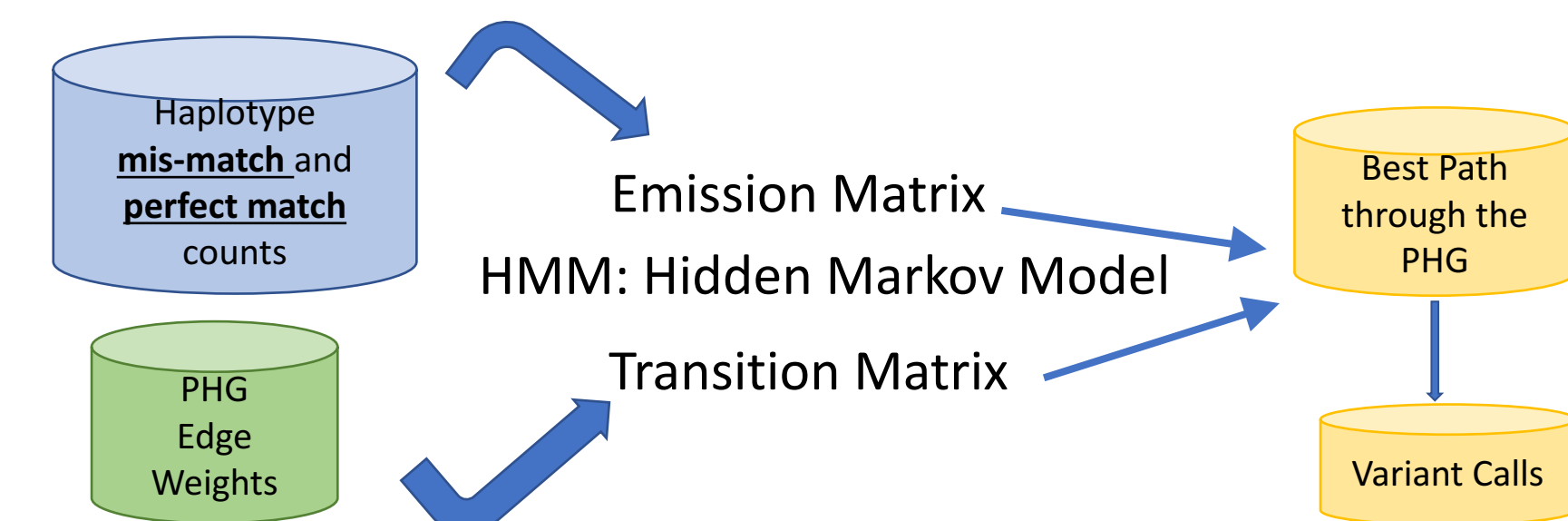
### 2. Infer genotypes from skim sequences



Best path calculations: Determine read counts and edge weights for each haplotype



Use HMM to find the best path through PHG



- **Perfect match count:** Number of sequence reads that match the haplotype sequence exactly.
- **Mis-match count:** Number of reads that map to a haplotype but do not match exactly.

Common pipeline to call genotypes from technology independent skim sequences (GBS, rAmpSeq, Nextera, Nanopore, RNA-seq, DArT-seq, etc.)

## PHG will be available as a docker image

### Running PHG scripts in a docker container

(1) `docker run --name load_phg_container \`  
`-v localMachine:/pathToOutputDir:/tmpDir/outputDir \`  
`-v localMachine:/pathToReference:/tmpDir/refDir \`  
`-t phgdocker:latest \`  
`/LoadGenomeIntervals.sh config.txt ref.fa anchors.bed data.txt`

(2) `docker run --name haplotypes_phg_container \`  
`-v localMachine:/pathToOutputDir:/tmpDir/outputDir \`  
`-v localMachine:/pathToReference:/tmpDir/data/reference/ \`  
`-v localMachine:/pathToWGSBams:/tmpDir/data/bam/DedupBAMS/ \`  
`-v localMachine:/pathToPhgMaizeDB:/tmpDir/output/phgMaizeDB.db \`  
`-t phgdocker:latest \`  
`/CreateHaplotypes.sh /tmpDir/data/config.txt taxon single anchors.bed`

(3) `docker run --name consensus_phg_container \`  
`-v localMachine:/pathToOutputDir:/tmpDir/outputDir \`  
`-v localMachine:/pathToReference:/tmpDir/data/reference/ \`  
`-v localMachine:/pathToPhgMaizeDB:/tmpDir/output/phgMaizeDB.db \`  
`-v localMachine:/pathToGvcfsToLoad:/tmpDir/data/outputs/gvcfs/ \`  
`-v localMachine:/pathToFastasToLoad:/tmpDir/data/fastas/ \`  
`-t phgdocker:latest \`  
`/CreateConsensi.sh /tmpDir/data/config.txt ref.fa ref_version`

(4) `docker run --name findPath_phg_container \`  
`-v localMachine:/pathToOutputDir:/tmpDir/outputDir \`  
`-v localMachine:/pathToReference:/tmpDir/data/reference/ \`  
`-v localMachine:/pathToGBSFastq:/tmpDir/data/fastq/ \`  
`-v localMachine:/pathToConfig:/tmpDir/data/config.txt \`  
`-v localMachine:/pathToFindPathDir:/tmpDir/output/pathDir/ \`  
`-t phgdocker:latest \`  
`/FindPath.sh myTaxa config.txt CONSENSUS ref.fa HAP_METHOD ref_version PATH_METHOD`

For more details: <https://bitbucket.org/bucklerlab/practicalhaplotypegraph/wiki/Home>

## Acknowledgements

This work was generously supported by the USDA-ARS, the Bill & Melinda Gates Foundation (OPP1159867), and the NSF Plant Genome Research Project (IOS#1238014).