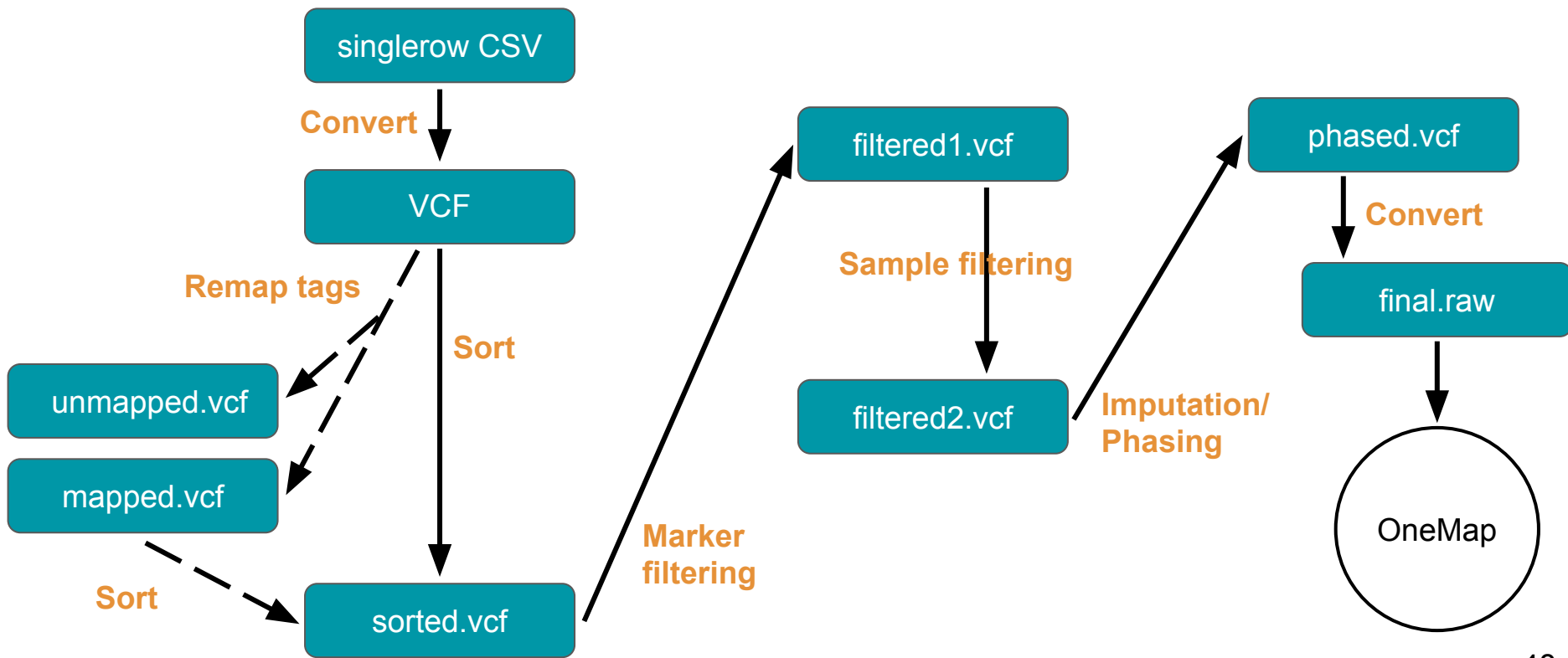# Wednesday, 25 July 2018

- Marker filtering

- Sample Filtering

- Imputation
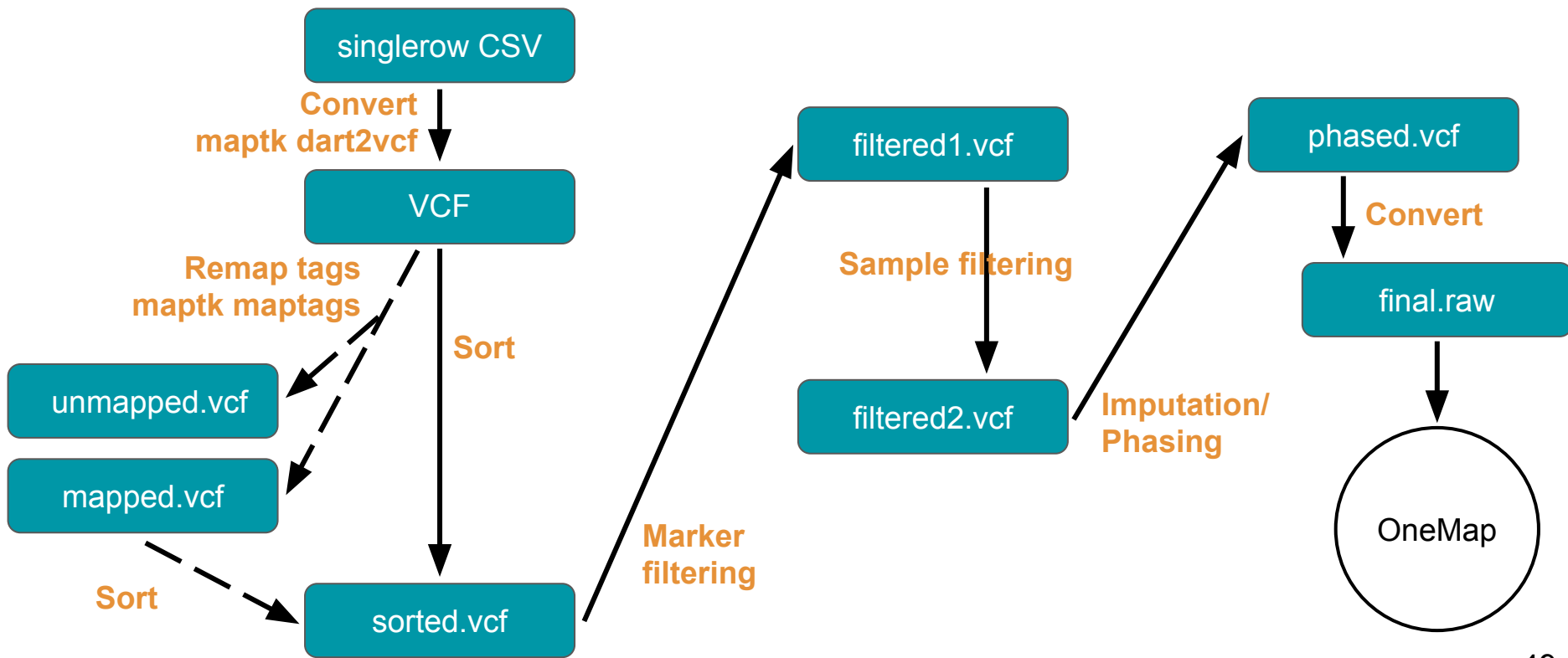
Windows users will need WinSCP today

+ Daily reminder to copy your commands for later reference :-)

# DArTseq to mapping: analysis flow chart
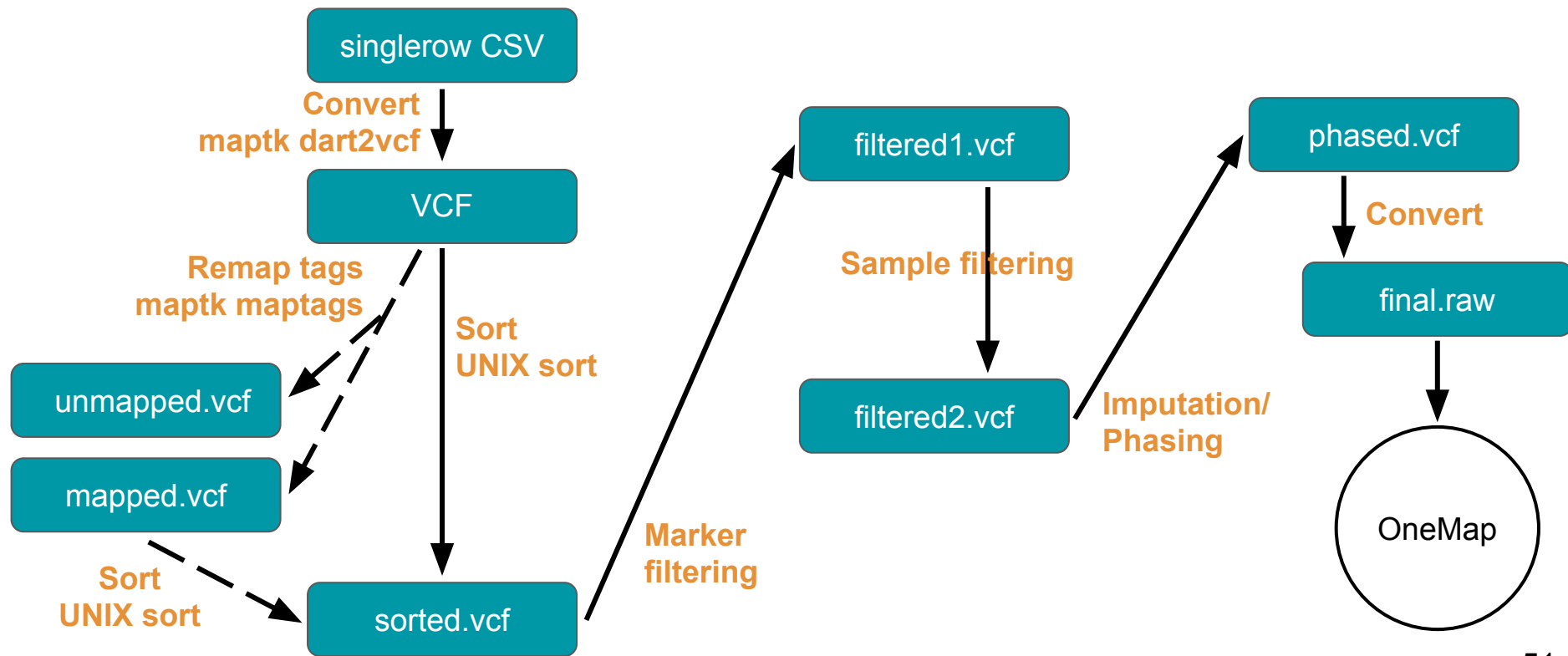


48

# DArTseq to mapping: analysis flow chart

# Preparing the (mapped) VCF file for filtering



Why sort?

# DArTseq to mapping: analysis flow chart

# Filtering loci for segregating markers

# What markers are useful for F1 linkage mapping?

- Useful markers have two or more alleles segregating in a Mendelian manner
- At least one of the parents must be heterozygous



1 Aa

# $X^2$ test for Mendelian segregation

Goodness-of-fit test:
- Expected counts vs. observed counts
- "Test statistics that follow a chi-squared distribution arise from an assumption of independent normally distributed data" – Wikipedia



|   | A | a |
|---|---|---|
| **A** | AA | Aa |
| **a** | Aa | aa |

1 AA : 2 Aa : 1 aa

|   | A | a |
|---|---|---|
| **A** | AA | Aa |
| **A** | AA | Aa |

1 AA : 1 Aa

|   | a | a |
|---|---|---|
| **A** | Aa | Aa |
| **a** | aa | aa |

1 Aa : 1 aa

# MapTK chif1 command



```
[lyons@ip-172-31-31-27 ~]$ maptk chif1
Prototype mismatch: sub main::assert: none vs (&;@) at /usr/local/lib/Exporter.pm line 66.
 at /usr/local/bin/maptk line 40.

Usage:    maptk chif1 [options] <in.vcf>

Options: -o <file>       Write output to FILE [stdout]
         -A <ufloat>     Max. phred allele freq. best-fit P-val [50.0]
         -G <ufloat>     Max. phred genotype freq. best-fit P-val [9999.0]
         -P <P1>[,<P2>]  Parental IDs for the input population (recommended)
         -e <uint>       Phred-scaled genotype-call error (see Notes 3) [30]
         -a              Perform chi-sqr on allele depths (AD) field [GT]
         -F              Exclude sites for which scores cannot be applied
         -p              Force pseudo-testcross markers only
         -S              Silence/disable verbose reporting
         -t              Write output as a table (see Note 4) [VCF]
         -h              This help document
```
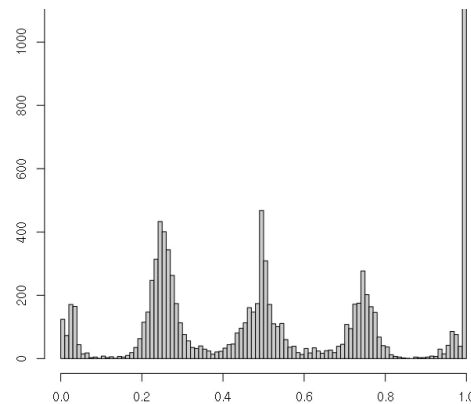
Notes:

1. This script applies a chi-squared goodness-of-fit test for Mendelian genotype frequencies, for a determined allele frequency class (3:1, 1:1, etc). The allele frequency class is selected within a goodness-of-fit tolerance threshold set by the '-A' option. The chi-squared values calculated here have been compared to those calculated by JoinMap4 and there is (almost) complete agreement.

2. Parental genotypes are inferred and stored by the 'P0GT' key in the INFO field. If the IDs of the parents are passed via the '-P' option, and the samples are included in the input VCF, an attempt to orient the inferred genotype calls with respect to parent (determined by the relative order of the parental samples listed in the VCf header) is made. If successful, a 'P0PHASED' key is applied to the INFO field.

3. Requires at least five diploid F1 samples (ten or more recommended).

4. Yates' correction for continuity is applied to sites with less than ten observed chromosomes, and no calculations are performed on sites with less than five observed chromosomes.

5. Value passed to the '-e' argument must be a positive integer. It is recommended to set '-e' to the minimum GQ value use for filtering.

6. When enabling the '-a' flag, and the input data are at low-coverage, it is necessary to include the parental IDs via '-P' to calculate the P-value statistics accurately.

7. The default output format is VCF, the '-t' option outputs a tab-delimited table (list of sub-field values) of relevant statistics a-la `vcftools --get-INFO`.

**Note**: this command requires all sites to have less than 50% missing data.
VCFtools can help you remove sites with >50% missing data.

# MapTK chif1

How do you know if chif1 ran correctly?

# MapTK chif1: test statistics

- **F1AFP**: Phred-scaled P-value for $X^2$ goodness-of-fit test on allele frequencies
- **F1GTP**: Phred-scaled P-value for $X^2$ goodness-of-fit test on genotype frequencies
- **F1X2**: $X^2$ value for goodness-of-fit test on locus allele frequencies
- **P0GT**: The Inferred parental genotypes, ordered alphabetically by Sample ID
- **P0PHASED**: Boolean tag indicating the locus passed the specified $X^2$ test (and, if the option was enabled, the parental genotypes were able to be inferred)

Phred = $-10 \log_{10}(\text{Prob})$

# DArTseq to mapping: analysis flow chart

# Filtering samples for full-sib progeny

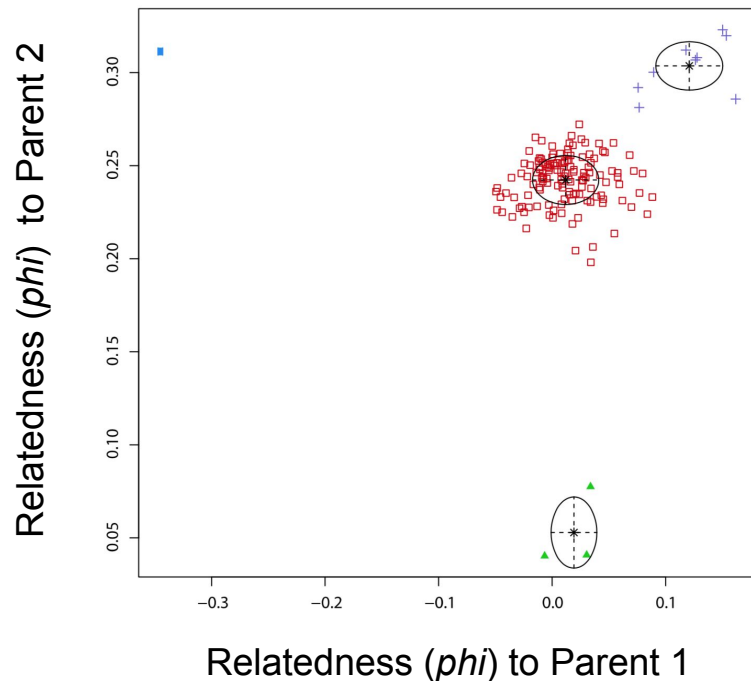# Sources of error in crossing experiments

1. Half-sibs
2. Selfs
3. Volunteer seedlings
4. Human error

# Filtering samples for full-sib progeny

- Calculate relatedness* of putative offspring w.r.t. the parents.
- Relatedness *phi* about half of relatedness expected:

| relationship | Expected | Reported *phi* |
|---|---|---|
| clonal | 1.0 | 0.5 |
| parent–child | 0.5 | 0.25 |
| off-type | 0.0 | 0.0 |



Relatedness (*phi*) to Parent 1

*Manichaikul. 2010. Bioinfo. doi: 10.1093/bioinformatics/btq559
ICGMC. 2015. G3 Journal. doi: 10.1534/g3.114.015008
https://bitbucket.org/rokhsar-lab/gbs-analysis/src/master/

61

# Filtering samples for full-sib progeny

vcftools relatedness2

➔ Remember to check where it wants a file, and with a <u>prefix</u>

maptk mclust

➔ Check for files ending in .dat and .dat.pdf
➔ **Note**: mclust runs R to produce the plots, if the Mclust R package is not installed, R will not print the plots

# Filtering samples for full-sib progeny
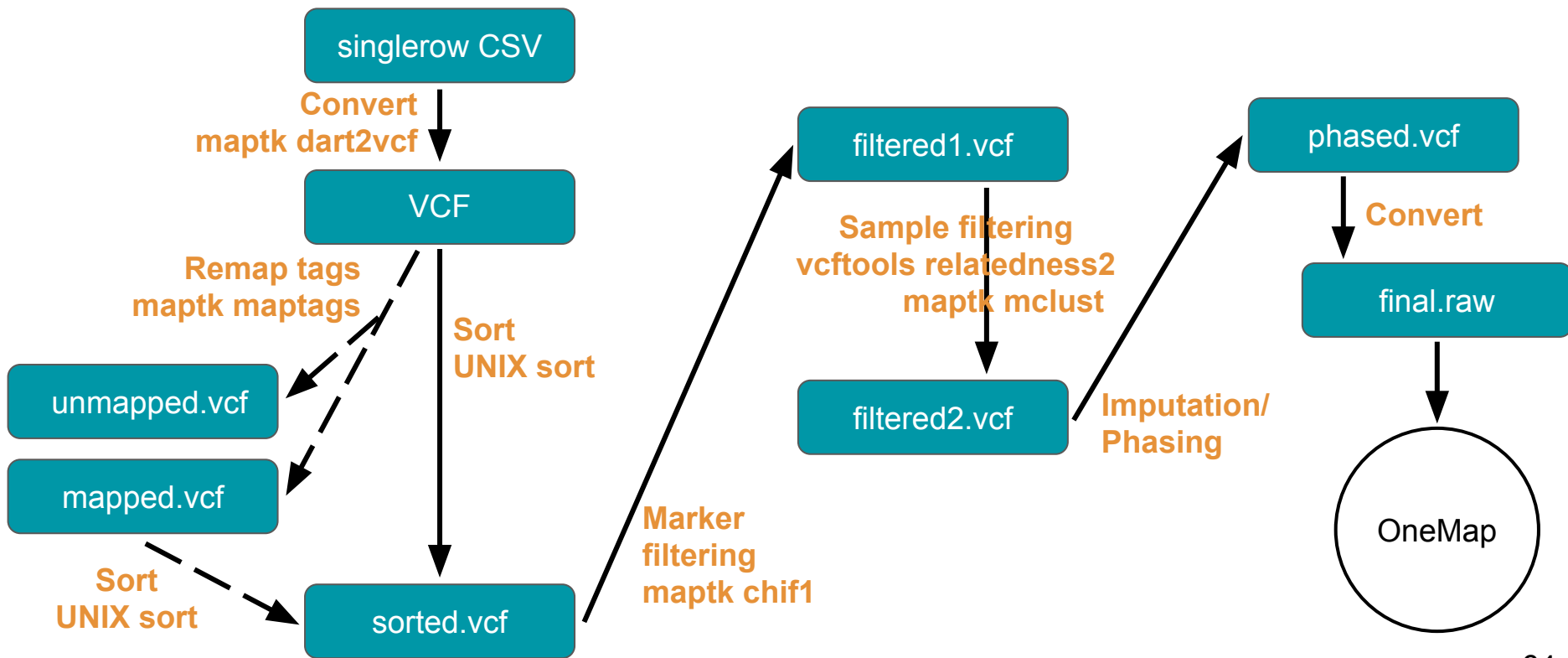
Download the plots to your laptop:

**Mac/Linux**:

```
scp -i ~/fn.pem \
fn@ec2-13-57-194-80.us-west-1.compute.amazonaws.com:/home/fn/mclust.dat.pdf .
```

Where "fn" is your family name

**Windows**:

Download and install WinSCP

# DArTseq to mapping: analysis flow chart



64

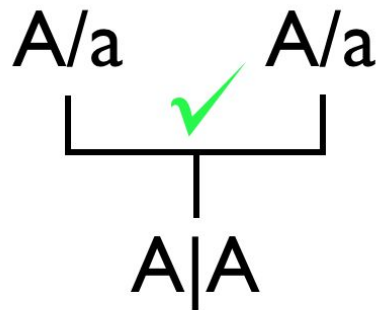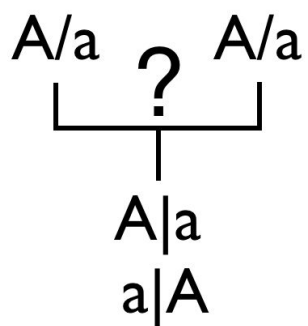# Imputing genotypes
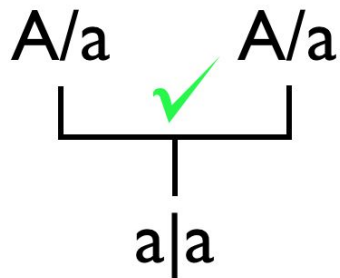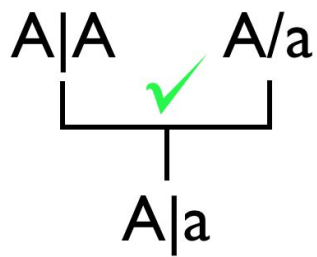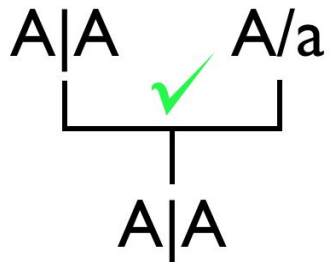
# Imputation and phasing

What is imputation:

➔ Filling in missing values (*i.e.*, genotypes) with data inferred from patterns/correlations in the dataset.

What is phasing?

➔ Assigning alleles to their maternal or paternal chromosome.

phaseF1 script

# Imputation and phasing

# Imputation and phasing

**Mapping algorithms are sensitive to missing data, but *more* sensitive to incorrect data.**
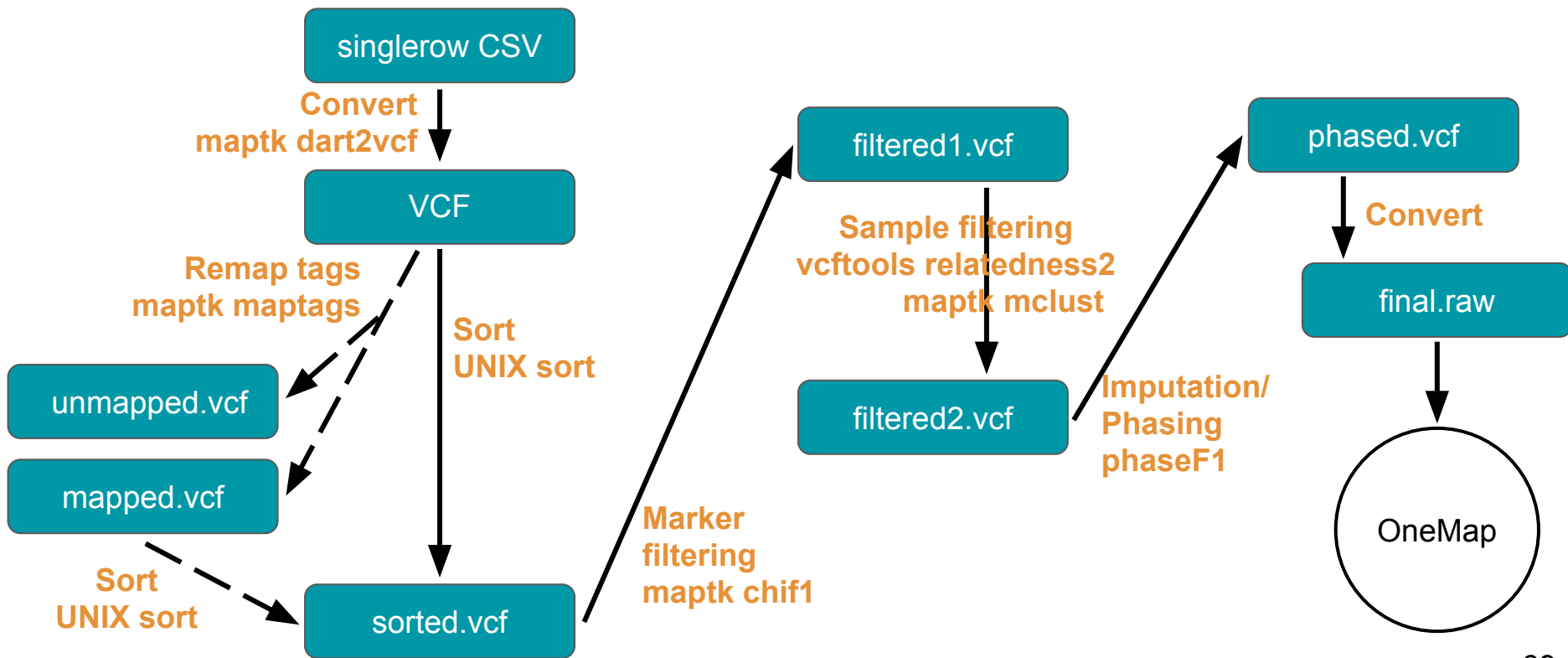
When to impute? When not?

- With OneMap we'll be able to judge correctness of imputation.
- Iteration/parameter sweeps.

Imputation/phasing software:

1. Beagle
2. MACH
3. IMPUTE2

# DArTseq to mapping: analysis flow chart

# DArTseq to mapping: analysis flow chart