



## Planning Considerations

### Reference Population

A reference population with very dense genotyping can be used as a scaffold to align data from different experiments for imputation. The reference data provides a denser set of markers with minor allele frequency and LD information that can be used to inform imputation across data sets. The reference population should be representative of the experimental sample population<sup>9</sup>. For example, if the experimental data were collected from individuals of Caucasian ancestry, then a Caucasian reference sample (e.g. HapMap CEU samples) should also be used. Likewise, for samples of mixed or alternative ancestry, an appropriate reference sample should be used. Huang et al. present the optimal proportions of CEU/CHB/JPT/YRI HapMap samples for imputing diverse world populations<sup>10</sup>.

### Consistent Strand

When merging data sets, it is essential that genotypes from both data sets are presented consistently from the same strand (e.g., forward or “+” strand). Errors in this consistency will result in an inability to merge data and cryptic strand flips of A/T and C/G SNPs, which can lead to spurious results<sup>7</sup>.

HapMap reference data are provided from a number of sources (see the Reference Data Sets section) on the forward strand. Beginning with the Infinium<sup>®</sup> HD HumanOmni1-Quad BeadChip, Illumina will provide strand annotation files for all its products, which researchers can obtain by contacting Technical Support. These strand annotation files can be used to identify markers assayed on the reverse strand. Researchers can flip reverse strand markers using a program such as PLINK before merging with reference data for imputation.

However, there are suspected strand errors in the HapMap data, so researchers should expect that a few markers (usually not more than a few thousand out of a whole-genome data set) display irreconcilable strand differences. Illumina recommends removing those SNPs from the experimental data set and proceeding with imputation.

### Initial Quality Control

Before genotype imputation, Illumina recommends that researchers carry out basic data quality checks on available genotypes in the experimental data set. This generally includes removal of<sup>7</sup>:

- Markers with low call rate
- Large deviations from Hardy-Weinberg equilibrium
- Large numbers of discrepancies among duplicate samples

- Mendelian inconsistencies
- Markers with very low minor allele frequency (MAF)

Optimal cutoffs for these metrics vary, and researchers should use appropriate scores for their particular study.

### Confidence Threshold for “Hard” Genotype Calls

When using imputation to facilitate the merging of data sets from different sources for a combined analysis, the goal is to fill in missing genotypes across two different data sets. Genotypes can be called at each missing data position for later association analysis.

Imputation provides a probability for each of the three possible genotype classes, and calls are based on the most likely genotype at each position<sup>9</sup>. When a hard genotype call is made, it carries with it a confidence score that corresponds to the likelihood that the called genotype was the correct choice. For example, if the genotype AA had a probability of 95% versus the genotype of AB having a probability of 3%, the confidence score for the choice of AA would reflect the overwhelming likelihood that the true genotype is AA. If the probability of AA was 40% and the probability of AB was 30%, making a hard genotype call is not as clear-cut and would be reflected in a lower confidence score. Imposing a stringent cutoff based on confidence scores will decrease the likelihood of imputation error in downstream association testing<sup>9</sup>. Refer to the accompanying documentation for each type of imputation software for more information on interpreting confidence scores.

### Imputation Accuracy

In addition to imposing a stringent confidence score cutoff, several strategies can be employed to minimize potential errors in imputation. Watching out for these red flags will reduce the chances of inaccurate data interpretation due to imputation error.

Imputation is based on LD, so it will not predict completely independent regions of the genome. Association tests of flanking markers should show similar levels of association compared with an imputed marker. Therefore, an imputed marker with a dramatically different association statistic than the surrounding directly genotyped markers should be treated with caution and investigated carefully. An exception to this could occur when large amounts of data are missing in two or more data sets, which are merged with a reference data set. For example, if 50% of the data was missing at one SNP and 50% of the data was missing at a second neighboring SNP, then after imputation there would be nearly 100% of the genotypes for those markers across all data sets. This imputation scenario would provide such a

Table 1: Commonly Used Imputation Software Packages

Software Name	Institution	URL
<b>Mach</b>	University of Michigan <sup>1,2</sup>	<a href="http://www.sph.umich.edu/csg/abecasis/MaCH/tour/imputation.html">http://www.sph.umich.edu/csg/abecasis/MaCH/tour/imputation.html</a>
<b>Beagle</b>	University of Auckland <sup>3</sup>	<a href="http://faculty.washington.edu/browning/beagle/beagle.html">http://faculty.washington.edu/browning/beagle/beagle.html</a>
<b>Impute</b>	Oxford University <sup>4,5</sup>	<a href="http://mathgen.stats.ox.ac.uk/impute/impute.html">http://mathgen.stats.ox.ac.uk/impute/impute.html</a>
<b>Plink</b>	Massachusetts General Hospital / Broad Institute <sup>6</sup>	<a href="http://pngu.mgh.harvard.edu/~purcell/plink/">http://pngu.mgh.harvard.edu/~purcell/plink/</a>





