

Genotyping Rare Variants

A simulated analysis achieves high call rates and low error rates from loci containing rare variants.

Introduction

Over the last decade, the International HapMap project has been the major resource of validated SNPs for the development of whole-genome genotyping arrays. SNPs chosen for the HapMap database and used in commercial array products were selected to preferentially include common (minor allele frequency (MAF) > 5%) variants and provide only a small subset of less common (MAF < 5%) variants. As a result, the arrays designed from this data set have primarily targeted common variants.

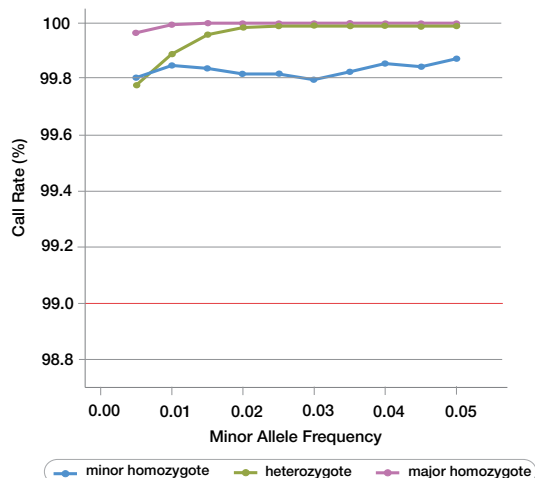
By sequencing thousands of individuals across many populations, the 1000 Genomes Project is seeking to identify SNPs down to 1% MAF, dramatically expanding the catalog of known human variants. This data set is expected to provide a more comprehensive understanding of the true spectrum of variation in human populations and provide a more accurate picture of the linkage disequilibrium (LD) structure of the human genome, enabling researchers to explore the role of rarer variants in traits and diseases. Illumina is leveraging these data to develop the next generation of whole-genome genotyping arrays, improving genomic coverage and fueling new discoveries.

As the target MAF of the next-generation microarrays decreases, it is essential that the quality of these rarer variant genotyping calls

remain high. Genotype calls are normally made using a training data set, where the assay signal is referenced against the expected cluster positions of the three genotype classes: homozygous major, homozygous minor, and heterozygous. However, as the MAF for a given SNP becomes increasingly rare, the training data sets will contain fewer examples of the homozygous rare and heterozygous genotypes. This can impact the reliability of genotype calls made using this limited data set. For common variants, it is well understood that re-clustering the experimental data set using GenomeStudio® software can significantly improve both call rates and reproducibility for a given study, but there is limited data showing how this re-clustering approach will perform for rare variants.

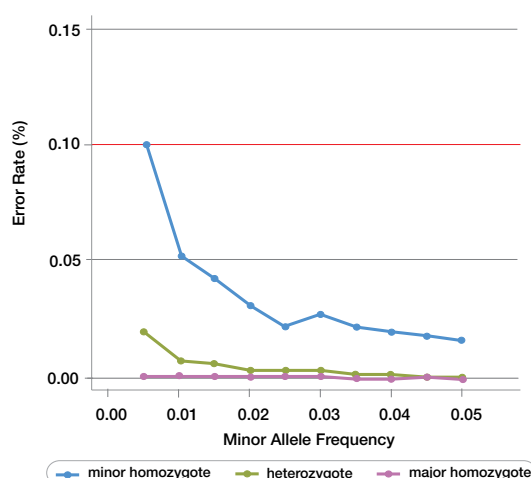
One challenge with determining the quality of rare variant genotyping calls is the lack of a “gold standard” reference source for the true genotypes of a given SNP in a given sample. One approach is to calculate transmission rates and Mendelian inheritance errors from a large data set of related individuals; however, data sets of this type are not widely available. An alternative approach is to simulate rare variants by down sampling from common, high-quality, and well-accepted SNP genotype data and comparing the new call rate and genotype information with the original genotype calls. This technical note presents the results of such a rare variant genotyping simulation and demonstrates that these variants can be called with high accuracy.

Figure 1: Call Rates for Simulated Genotype Calling



A simulated analysis shows that high call rates (> 99%) can be achieved for alleles that exhibit low minor allele frequencies. The red line indicates the Illumina product specification, 99%.

Figure 2: Error Rates for Simulated Genotype Calling



Low error rates (< 0.1%) can be achieved for all three genotype classes for alleles that exhibit very low minor allele frequencies. The red line indicates the Illumina product specification, 0.10%.

CGAATGATAACAGTAACACACTTCTGTTAACTTAAAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTACCGTCTTTCTGTTAACTTAAAGATTACTTGATCCACTGATTCAACG
 TCAACGTACCCTAACGAACGTATCAATTAAAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTACCGTCTTTCTGTTAACTTAAAGATTACTTGATCCACTGATTCAACG
 CGACGAAAGAAATGATAACAGTAACACACTTCTGTTAACTTAAAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTACCGTCTTTCTGTTAACTTAAAGATTACTTGATCCACTGATTCAACG
 AACGTACCATTAAAGAGCTACCGTCTTTCTGTTAACTTAAAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTACCGTCTTTCTGTTAACTTAAAGATTACTTGATCCACTGATTCAACG
 AGAATGATAACAGTAACACACTTCTGTTAACTTAAAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTACCGTCTTTCTGTTAACTTAAAGATTACTTGATCCACTGATTCAACG
 GATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTACCGTCTTTCTGTTAACTTAAAGATTACTTGATCCACTGATTCAACG
 CGTATCAATTGAGACTAAATATAACGTACCATTAAAGATTACTTGATCCACTGATTCAACGTACCCTAACGAACGTATCAATTGAGACTAAATATAACGTACCATTAAAGAGCTACCGTCTTTCTGTTAACTTAAAGATTACTTGATCCACTGATTCAACG

