





- 5) To determine whether any additional samples should be excluded from the project, reevaluate samples as described in the previous section.
- 6) When the GenomeStudio project contains only final and unique samples, update the SNP statistics.

## Locus Analysis and Reclustering

Standard cluster files provided with Infinium products identify expected intensity levels of genotype classes for each SNP. For human projects, comparing sample intensities to this cluster file is sufficient for generating high-quality data. Standard cluster files for standard human Infinium products are created using a diverse set of over 100 samples from the Caucasian (CEU), Asian (CHB+JPT), and Yoruban (YRI) HapMap populations. Therefore, much of the genetic diversity in these populations is incorporated. Use this standard file when call rates are high and the genetic diversity of a sample set is well represented by the standard cluster file.

However, in some situations, sample intensities might not overlay perfectly onto the standard cluster positions. Reclustering some or all SNPs can optimize the GenomeStudio software's ability to call genotypes and may result in higher overall call rates (Figure 3). In some cases, reclustering is not helpful and the locus must be zeroed (Figure 4). Reclustering using samples from a particular project will redefine cluster positions, making them more representative of the samples in the project.

All or a subset of loci can be reclustered to generate a custom cluster file. To recluster all SNPs in the GenomeStudio project select Analysis | Cluster All SNPs. To recluster a subset of SNPs, highlight relevant rows of the SNP Table, right-click the SNP Table, and select Cluster Selected SNPs. An alternative method is to use filters to identify and select SNPs that should be excluded from reclustering. Consult the GenomeStudio Framework User Guide (part # 11318815) for information about how to create filters for various metrics in the Full Data Table, SNP Table, or Samples Table. Generating a custom cluster file will require a thorough evaluation of newly defined SNP cluster positions to ensure the accuracy of the genotype data. See Evaluating and Editing SNP Cluster Positions later in this document.

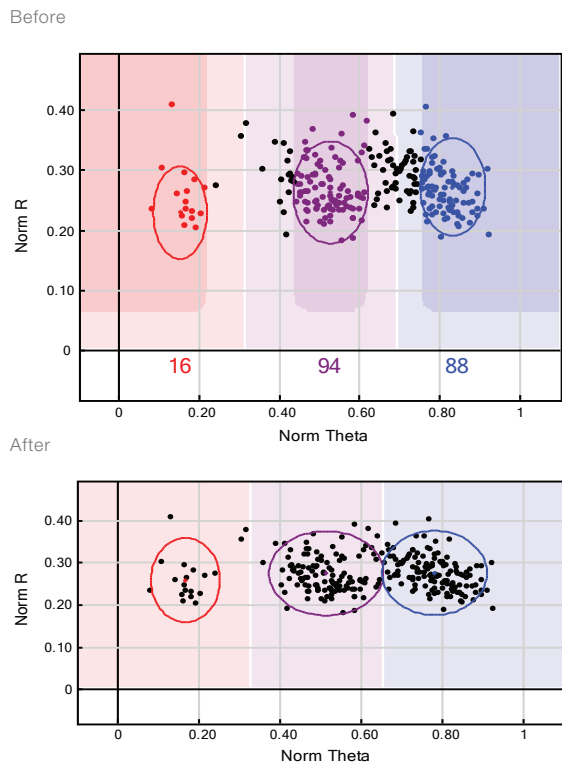
An important consideration in the decision to recluster is that the GenomeStudio clustering algorithms require approximately 100 samples to produce cluster positions that are representative of the overlying population. Therefore, in projects with fewer than 100 different samples, it is best to use the standard cluster file for calling genotypes.

Special consideration must be given to mitochondrial (mtDNA) and Y chromosome (Y-chr) SNPs. The GenomeStudio clustering algorithms do not automatically accommodate loci that lack heterozygous clusters. The mtDNA (Figures 5 and 6) and Y-chr SNPs (Figures 7 and 8) must be manually edited following reclustering. Before reclustering all SNPs, evaluate the mtDNA and Y-chr SNPs to decide whether the predefined Illumina standard cluster positions are appropriate for the current project samples. If appropriate, exclude these loci from reclustering.

## Evaluating and Editing SNP Cluster Positions

To identify loci that should be manually edited or zeroed due to poor performance (< 1%), evaluate newly reclustered SNPs using the metrics listed in the SNP Table. These metrics are based on all samples for each locus, and thus provide overall performance information for each locus. To find loci that may require editing or removal before generating final reports, sequentially sorting the SNP Table by columns and screening loci by each quality metric identified in Table 1 is recommended. To determine hard cutoffs and grey zones, sort data by one column at a time and explore values starting at the extremes of the ranges. A hard cutoff, defined as the level above (or below) which most of the loci are unsuccessful, can be applied and the loci zeroed when necessary. To improve analysis further, the grey zone, defined to contain loci that are 80–90% successful, can be manually adjusted. The upper (or lower) limit of the grey zone is the point at which all loci are successful. SNPs falling in the grey zone should be visually evaluated and either accepted, zeroed, or manually edited by moving cluster positions. Hard cutoffs and grey zones may not transfer between projects because they are highly dependent on

**Figure 4: Example of an Unsuccessfully Reclustered SNP**



At this locus, project samples form diffuse clusters and do not correlate well with the standard cluster positions (before). The genotype calls are unreliable because of poor cluster separation. Reclustering on this locus does not improve data quality. If not zeroed automatically by the algorithm (as shown), this locus should be zeroed based on the low score for cluster separation.

AAAGAATGATAACAGTAACACACTTCTGTTAAACCTTAAGATTACTTGATCCACTGATTCACCGTACCGTAACGAACTATCAAITGAGACTAAATATTAACGTACCATTAAAGAGTACCGTCTTCTGTTAAACCTTAAGATTACTTGATCCACTGATTC  
AATCAACGTACCGTAACGAACTATCAITAAAGATTACTTGATCCACTGATTCACCGTACCGTAACGAACTATCAAITGAGACTAAATATTAACGTACCATTAAAGAGTACCGTCTTCTGTTAAACCTTAAGATTACTTGATCCACTGATTC  
AACGACGAAAGAATGATAACAGTAAACACACTTCTGTTAAACCTTAAGATTACTTGATCCACTGATTCACCGTACCGTAACGAACTATCAAITGAGACTAAATATTAACGTACCATTAAAGAGTACCGTCTTCTGTTAAACCTTAAGATTACTTGATCCACTGATTC  
TTAAAGTACCATTAAGAGCTACCGTGC AACAGTAAACACACTTCTGTTAAACCTTAAGATTACTTGATCCACTGATTCACCGTACCGTAACGAACTATCAAITGAGACTAAATATTAACGTACCATTAAAGAGTACCGTGC AACGACGAAAGAATGAT  
AAAAGAATGATAACAGTAAACACACTTCTGTTAAACCTTAAGATTACTTGATCCACTGATTCACCGTACCGTAACGAACTATCAAITGAGACTAAATATTAACGTACCATTAAAGAGTACCGTCTTCTGTTAAACCTTAAGATTACTTGATCCACTGATTC  
AAGATTACTTGATCCACTGATTCACCGTAAAGATTACTTGATCCACTGATTCACCGTACCGTAACGAACTATCAAITGAGACTAAATATTAACGTACCATTAAAGAGTACCGTCTTCTGTTAAACCTTAAGATTACTTGATCCACTGATTC  
AACGTATCAAITGAGACTAAATATTAACGTACCATTAAAGAGTCTGTTAAACCTTAAGATTACTTGATCCACTGATTCACCGTACCGTAACGAACTATCAAITGAGACTAAATATTAACGTACCATTAAAGAGTACCGTGC AACGAAAGAATGATAAC















