

Informatic Updates for Whole-Genome Microarrays

Illumina has released new product files (.bpm, .egt, and .xml) for all currently available human whole-genome genotyping microarrays. These new files provide a number of improvements requested by researchers, including:

- **Marker positions in human genome build 37/hg19**
- **Genomic strand information (+/-)**
- **Removal of indeterminate loci**

What is a new genome build?

While the human genome reference sequence has been available to researchers for many years, groups such as the Genome Reference Consortium are constantly striving to improve the sequence accuracy. Periodically (every few years), they release a new reference genome that incorporates the latest data, fills sequence gaps, corrects any known errors, and maps any newly discovered markers or features across the genome.

What is the most recent genome build?

The most recent human genome build is commonly referred to as hg19 or build 37, released in February 2009. The new genome build can be visualized using publically available tools like the UCSC Genome Browser (<http://genome.ucsc.edu/>).

What is the difference between build 36 and build 37?

Build 37 offers a large number of small improvements to the accuracy of the reference sequence. For complete release notes, please review the information provided by the Genome Reference Consortium, which lists the changes made to each chromosome in detail. The release notes are available from the following web address: <http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/human/index.shtml>.

How does using build 37 information help my research?

Build 37 represents the most up-to-date information from the human genetics community regarding the location of features (i.e. genes, exons, SNPs, repeat regions etc.) across the human genome. It is important that research findings be placed in the correct context across the genome for data interpretation.

What is the benefit of strand information?

SNPs can be assayed against either strand of DNA. Knowing this information allows researchers to convert genotypes to a consistent strand, which enables data sets to be merged across different microarrays for imputation analysis and facilitates comparison with publicly available data such as that from the HapMap and 1000 Genomes Projects.

Will I need a new version of GenomeStudio® software to read these new manifest files?

GenomeStudio v1.6 (version 2010.1) or higher is needed to read the new manifest files that include +/--strand information.

What do you mean by “indeterminate loci”?

As researchers begin to incorporate more rare variants into their studies, having very high genotype accuracy is increasingly important. To address this, Illumina scientists are incorporating additional tests into the quality control (QC) pipeline for whole-genome genotyping arrays. These tests include more diverse DNA samples during the cluster training process, increased standards for reproducibility of genotypes (within and across products), as well as concordance with high-quality sequence data. Moving forward, these QC tests will be a standard part of the GenTrain process (cluster training process) for new products. These stringent standards have been applied to all currently available whole-genome genotyping arrays. Markers that do not pass the increased stringency of these new tests are referred to as “indeterminate loci.”

How should I use these files if I am just starting a new project?

Researchers should process samples according to standard protocol using the new manifest, cluster, and product descriptor files.

