



**Table 1: Q-scores for Variant Calls Given Various Coverages**

	Reference Bases Called	Variant Base Calls	Expected Miscalls (at 1% error rate <sup>1</sup> )	Depth of Coverage	P Value <sup>b</sup>	Q-Score <sup>c</sup>
No Variant Present	100	0	1	100	1	0
5% Frequency SNV	285	15	3	300	6.7 x 10 <sup>-7</sup>	62
5% Frequency SNV	475	25	5	500	1.6 x 10 <sup>-10</sup>	98
5% Frequency SNV	950	50	10	1000	0	100 <sup>d</sup>
5% Frequency SNV	4750	250	50	5000	0	100 <sup>d</sup>
5% Frequency SNV	9500	500	100	10,000	0	100 <sup>d</sup>

a. <sup>1</sup>Conservative error rate used by somatic variant caller to provide high confidence calls.  
 b. The P value indicates the probability that the SNP is a false positive.  
 c. The Q-score indicates the degree of confidence.  
 d. <sup>d</sup>Q-scores above 100 are not reported.  
 e. The Amplicon - DS somatic variant caller computes the Q-score for an SNV based on a Poisson model. For example, the 5% frequency SNV in the second row identified at a read depth of 300. Assuming a conservative miscall error rate of 1%, this call has a Q-score of 62 and a P value of 6.7 x 10<sup>-7</sup>.

variants called in only one pool are flagged and filtered as probe pool biased (Figure 1).

### Superior Signal Resolution

The Amplicon - DS somatic variant caller identifies SNPs and short insertions and deletions (indels) present at low frequency in a DNA sample while minimizing false positives. To identify a variant, the Amplicon - DS somatic variant caller considers each position in the reference genome separately, counts the overlapping aligned reads at a given position, and computes a variant score that measures the quality of the call. Variant scores are computed using a Poisson model that excludes calls with scores below Q20 (> 99% predicted accuracy). Also, the algorithm only calls variants for bases that are covered at a depth of 300x or greater for a single amplicon (Table 1). Table 1 shows the likelihood of a false variant being called from a single strand given various depths of coverage. Variants are first called separately for each pool and are then compared and combined into a single output file. If a variant is present in both pools, yields at least 1,000x cumulative coverage between the 2 strands, and has a frequency of 3% or greater, it is marked as passing.

For simplicity, and to be conservative, the Amplicon - DS somatic variant caller only considers 2 alleles: reference and variant. Any other calls are considered reference calls. For example, A/C/G/T counts of 100/10/1/1 for reference A are considered to be K = 10 variants out of N = 112 coverage. Under the null hypothesis, it is assumed that no variant is present and that any nonreference calls are due to noise. Given a Q20 base filter, the acceptable noise level is 1%. For simplicity, it is assumed that the expected number of nonreference calls due to noise should follow a Poisson distribution with a mean of  $\lambda = 0.01 * N$ . The equation  $P = 1 - \text{CDF}(K - 1, \lambda)$  represents the probability (P) of having K or more variant calls, where CDF is the cumulative distribution function of the Poisson distribution. P is the probability that no variant is present, given K or more observations. In this way, P is the theoretical false-positive rate, and this probability is converted to a Q-score with the maximum Q-score set to 100.

The Amplicon - DS somatic variant caller addresses indels similarly. It analyzes aligned reads covering a given position that includes a particular indel to determine the variant count versus the overall coverage at that position. Due to the amplicon-based generation of the sequencing libraries, the variant caller does not perform some of the indel realignment steps included in other variant callers, such

as Genome Analysis Toolkit (GATK). Therefore, small indels (up to 25 bases in length) are reliably detected with this aligner. Illumina recommends that users review and confirm indels using visualization software such as the Integrative Genomics Viewer (IGV)<sup>1,2</sup>. All identified variants are provided in a VCF file, along with the frequencies at which they are detected in either the single or combined pool data.

### High Reproducibility

The Amplicon - DS somatic variant caller delivers a high level of reproducibility in detecting low-frequency mutations associated with cancer. In a survey of reference standards with quantified mutations in *BRAF*, *EGFR*, and *KRAS* genes, the Amplicon - DS somatic variant caller results were within ~1% of those reported by the manufacturer (Table 2). These results highlight the reproducibility and reliability of TruSight Tumor 26 and the MiSeq system when detecting somatic mutations.

### Data Analysis

The Amplicon - DS somatic variant caller is compatible with MiSeq Reporter software v2.2, or later. When used with the Amplicon - DS workflow, the algorithm exports a VCF file that contains the identified variants and their associated Q-scores and a genome VCF (gVCF) file. It is important to note that variant calls will only be made for those regions that meet the minimum coverage requirements mentioned previously. Illumina recommends that users review the Amplicon Coverage report generated by MiSeq Reporter software to identify any regions where coverage may be below the minimum requirement. In addition to assay variation, large deletions (greater than 25 bases) or complex substitutions can lead to a loss of coverage in the regions where they occur. This loss of coverage, unlike that caused by assay variation, will be evident in the data from both the FPA and FPB assay pools for such variants. These variants will not be captured in the VCF file and will therefore require manual annotation.

Q-scores calculated with the Amplicon - DS somatic variant caller are not equivalent to those generated using the standard TruSeq<sup>®</sup> Custom Amplicon analysis workflow. As the statistical model and details of

