

Processing Ion AmpliSeq™ Cancer Panel Data using NextGENe® Software

November 2011

John McGuigan, Megan Manion, Kevin LeVan, CS Jonathan Liu

Introduction

The Ion AmpliSeq™ Cancer Panel uses highly multiplexed PCR in order to generate amplicons from 46 different cancer genes. These amplicons are sequenced on the Ion PGM™, allowing for rapid turnaround and low cost. NextGENe® software allows for a customized analysis of the results, including detection of novel variants or alleles, alleles found at lower frequencies (less than 5%), or alleles in regions with less than 500x coverage. NextGENe includes many useful features, such as quality control reports, functional prediction scoring, and advanced project comparison. When using the Ion Torrent paired-end protocol, the overlapping reads can be merged into high-quality single reads in order to improve sequence accuracy. This can reduce the amount of filtering that is needed and increase the percentage of reads that successfully align. The analysis of this type of data is covered in the “Processing Paired End Data from the Ion PGM™ Sequencer Using the Floton Paired-End Merger in NextGENe® Software” app note.

Procedure

A given dataset can be processed in less than 15 minutes on a desktop computer running a 64-bit Windows operating system. All steps are performed with an easy-to-use point-and-click interface with no scripting required. Three datasets (one 316 chip and two 314 chips) were processed in this analysis.

Format Conversion

- As with all NextGENe projects, the raw data (FASTQ or SFF in this case) is first converted to FASTA format
- Basecall quality scores are used to filter and trim the raw data

Alignment to the human genome

- Alignment is performed against a pre-index human genome reference in order to avoid false positives caused by nonspecific amplification of untargeted regions
- Up to one mismatch is allowed in a read before the reads are broken into seeds.
- Mismatches are called as variants in the mutation report if they pass the mutation filter. Suggested settings require more than 100x coverage, 1% frequency, and 20 total reads with the mutation (Figure 1).
- Variants with F/R read balance < 0.2 and small homopolymer indels with F/R balance < 0.25 are removed from the final report in order to reduce the number of false positives.

The screenshot shows the 'Mutation Filter' section of the NextGENe software interface. It includes the following settings:

- Sample Trim:**
 - Select Sequence Range (From 1 Bases To 30 Bases)
 - Hide Unmatched Ends
- Mutation Filter:**
 - Mutation Percentage <= 1
 - SNP Allele <= 20 Counts
 - Total Coverage <= 100
 - Except for Homozygous
 - Use Original
 - Allow Software to Delete Mutations
 - Forward and Reverse Balance <= 0.2
 - Delete Small Homopolymer Indels if F/R <= 0.25

Figure 1 – Suggested mutation filter settings

Project review

- The mutation report, expression report, and coverage curve reports are filtered by loading a BED file specifying amplicon or targeted loci regions. The latter two reports provide information about coverage in the regions of interest. The mutation report has many additional filtering options.

Results

Table 1 lists the results of format conversion. Roughly 90% of the original reads were kept for each sample. Table 2 lists alignment results- around 95% of reads were successfully aligned, and about 95% of those reads were aligned in the amplicon regions. Figure 2 shows several mutations detected in a KRAS amplicon. One mutation is known (listed in the dbSNP database) and two others are complex, with multiple mutant alleles. Figure 3 lists the number of mutations found in each project. All 37 mutations called by the Ion AmpliSeq™ Cancer Variant Caller plugin were also found by NextGENe.

	FOZ-214 (316)	SUR-173 (314)	KER-780 (314)
Total Reads	1,167,000	578,080	569,813
Kept Reads	1,037,550	537,355	539,215
% Kept	88.91%	92.96%	94.63%

Table 1 – Format Conversion Results

	FOZ-214	SUR-173	KER-780
Aligned Reads	976,766	510,793	520,010
% Aligned	94.14%	95.06%	96.44%
Reads In Amplicon Regions	943,691	484,538	494,918
% in Regions	96.61%	94.86%	95.17%
Average Coverage in Amplicons	2,667x	1,621x	1,685x

Table 2 – Alignment Results

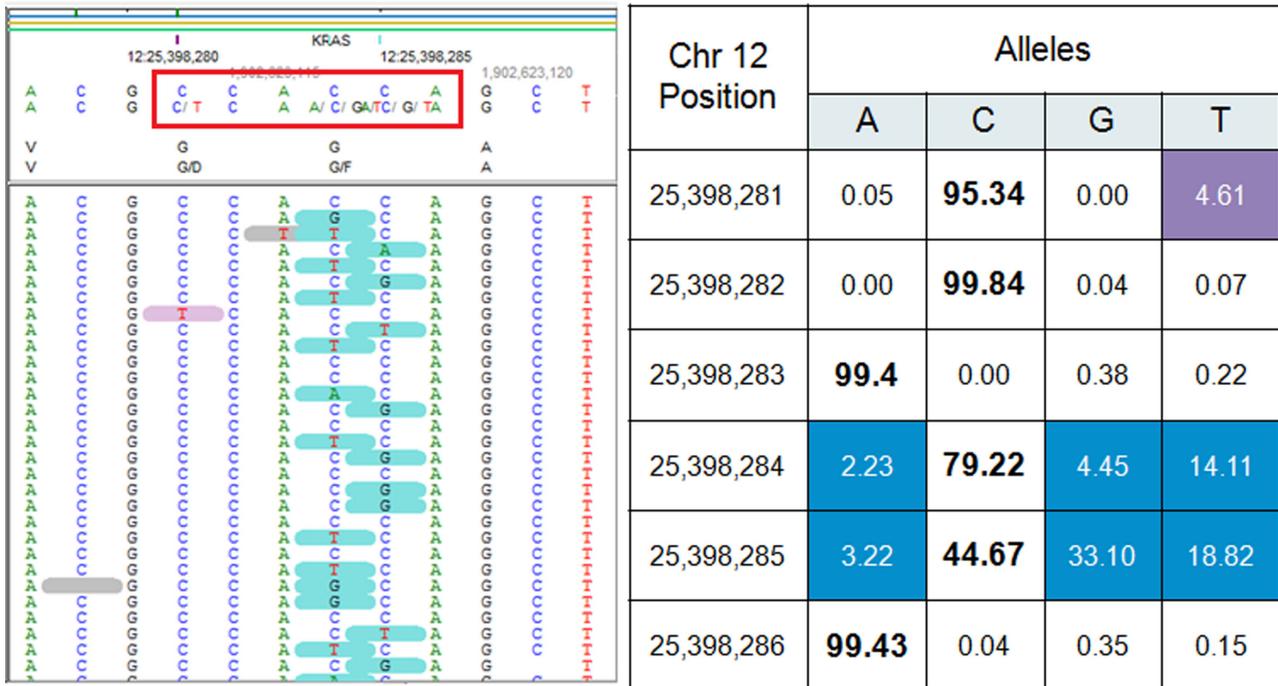


Figure 2 - A series of mutations found in KRAS. The first (purple) is a known mutation found in dbSNP. The other two have multiple mutant alleles at various frequencies. Reference allele frequencies are listed in bold.

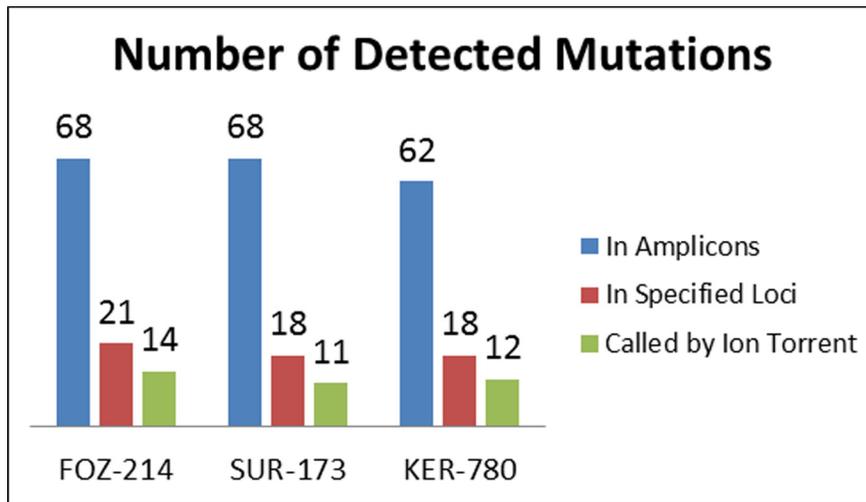


Figure 3 – Number of detected mutations in each project. NextGENe was able to detect more mutations than the built-in plugin because it can examine a larger portion of the amplicons and detect lower frequency alleles.

Discussion

When choosing alignment settings, it is important to consider the expected results. Increasing the minimum depth of coverage will reduce the number of false positives (even at lower mutation frequencies), but it may also decrease sensitivity. Setting a minimum number of mutant allele reads will allow for detection of low frequency variants in high coverage regions without allowing low frequency false positives to be called in low coverage regions. The coverage curve report is very useful for measuring potential loss of sensitivity (figure 4). The expression report can also be used to measure coverage (figure 5). This can be done for entire amplicons or for specific loci.

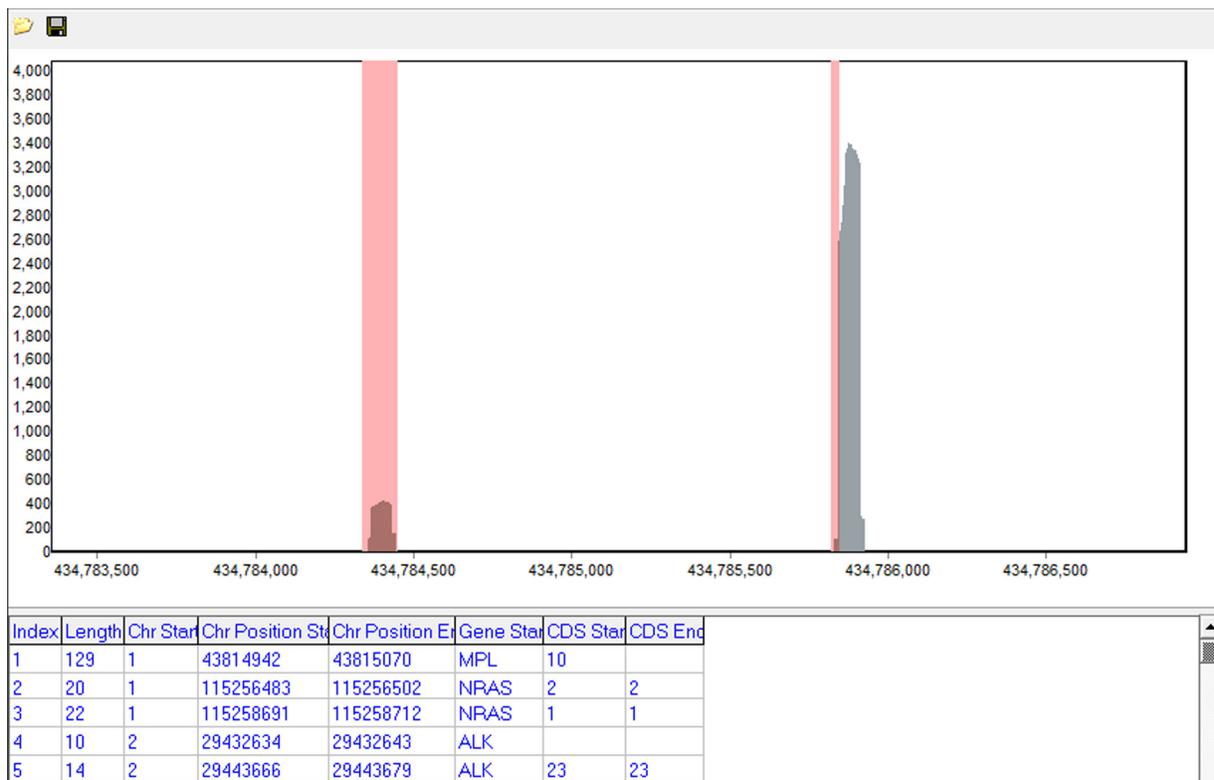


Figure 4 - The coverage curve report can be used to find positions (pink) lacking the desired level of coverage (gray)

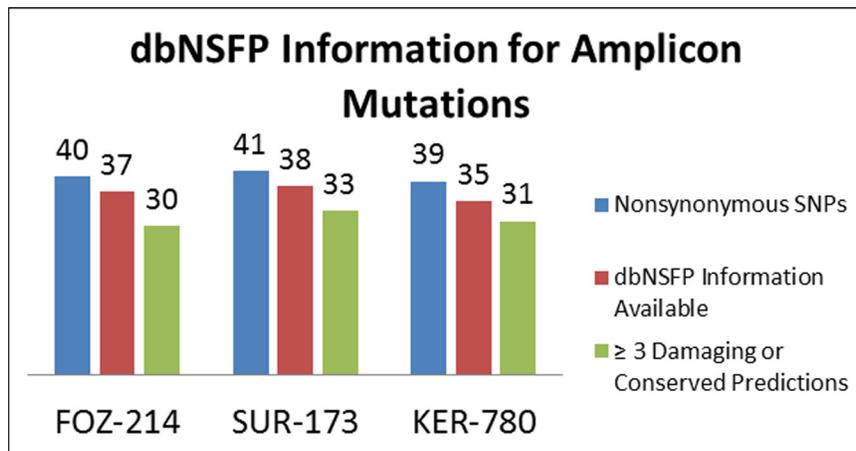


Figure 7 – Summary of dbNSFP Information. Most non-synonymous SNPs had dbNSFP information available, and most of these positions had several scores predicting damage/conservation.

Acknowledgements

We would like to thank Life Technologies for supplying the AmpliSeq data used in this analysis.

References

- [1] Liu, X., Jian, X. & Boerwinkle, E. dbNSFP: A lightweight database of human nonsynonymous SNPs and their functional predictions. *Human Mutation* 32, 894-899 (2011).
- [2] Forbes, S.A. et al. COSMIC: mining complete cancer genomes in the Catalogue of Somatic Mutations in Cancer. *Nucleic Acids Research* 39, D945-D950 (2010).