

Whole exome sequencing analysis pipeline for the discovery of mutations causative of human rare diseases

Francisco J. López Domingo¹, Antonio Rueda¹, Javier Pérez Florido¹, Alicia Vela Boza¹, Pablo Arce¹, Luis Miguel Cruz¹, Javier Escalante¹, Ana Isabel López¹, Federica Trombetta¹, Guillermo Antiñolo^{1,2}, and Javier Santoyo¹

¹ Andalusian Human Genome Sequencing Centre (CASEGH),
Medical Genome Project (MGP),

INSUR building, Albert Einstein st, Cartuja 93 Scientific and Technology Park,
41092 Seville, Spain

² Unidad de gestión clínica de genética, reproducción y medicina fetal. Instituto de Biomedicina de Sevilla (IBIS), Hospital Universitario Virgen del Rocío-CSIC-University of Seville,
Manuel Siurot Av., 41013, Seville, Spain
javier.santoyo@juntadeandalucia.es

Abstract

Recent advances in high-throughput sequencing technologies have made exome sequencing to be an outstanding tool for finding disease associated mutations at a relatively low cost. However, it is a non-trivial task to transform the vast amount of sequence data into meaningful variants to improve disease understanding. Several challenges arise when dealing with this approach, being critical checkpoints the raw read preprocessing, mapping procedure, variant calling and posterior variant selection. A number of computational algorithms and pipelines have been reported for variant analysis, none of them providing a complete strategy from raw data to mendelian analysis results. Here, we present a methodology that spans from SOLiD raw reads processing to mendelian analysis and variant selection, and its application over a set of samples from The Medical Genome Project, which proves the good performance of the applied methodology.

As stated above, the input of the pipeline is an xsq file generated by Applied Biosystem SOLiD 5500 XL sequencers, while the output is the result of variant annotation and mendelian analysis, assuming samples to be derived from a group or a family. A brief description of the steps is provided below:

1. Fasta and qual files generation from xsq files.
2. Duplicated reads removal.
3. BLAT-like Fast Accurate Search Tool v0.7.0a (BFAST) [1] for read mapping.
4. BAM cleaning: duplicated alignments and mismatched reads removal.
5. BAM realignment and SNV calling using the Genome Analysis Toolkit v1.4.14 (GATK) [2].
6. Variant quality filter based on GATK Best Practices V3 and depth filter.

7. Annotate Variation package (ANNOVAR) for variant annotation [3]; SIFT [4] and Polyphen [5] for variant function impact prediction; 1000 genomes [6] and dbSNP [7] for assessment of variant frequency.
8. Mendelian filter of deleterious variants.

The Medical Genome Project (MGP) aims to characterize a large number of rare genetically-based diseases. As a proof of concept, we selected from the MGP a set of affected individuals by several hereditary rare diseases, their healthy relatives and a set of 50 control healthy individuals from Spanish population. The full methodology was run and the results reveal a number of deleterious haplotypes in several genes which could be directly associated with the diseases.

The validation of some of the predicted variants by the pipeline demonstrates the good performance of our methodology analysis. Critical aspects to achieve such good performance are (i) BAM filtering, since an excessive number of mismatches are allowed by BFAST for short reads; (ii) the selection of variant filters and quality thresholds as recommended by GATK Best Practices V3 in combination with a depth threshold allowing high quality calls and (iii) the inclusion of control individuals in the analysis, which is essential since they remove population variants which can disturb the interpretation of the final variant set.

Keywords: NGS analysis pipeline, whole exome sequencing, variant analysis, rare hereditary diseases

References

1. Homer, N., Merriman, B., Nelson, SF.: BFAST: An Alignment Tool for Large Scale Genome Resequencing. PLoS ONE 4(11) (2009)
2. DePristo, M., Banks, E., Poplin, R., et al.: A framework for variation discovery and genotyping using next-generation DNA sequencing data. Nature Genetics. 43,491-498 (2011)
3. Wang, K., Li, M., Hakonarson, H.: ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. Nucl. Acids Res. 38(16) (2010)
4. Kumar, P., Henikoff, S., Ng, P.C.: Predicting the effects of coding non-synonymous variants on protein function using the SIFT algorithm. Nature protocols 4, 1073-1081 (2009)
5. Adzhubei, IA., Schmidt, S., Peshkin, L., et al.: A method and server for predicting damaging missense mutations. Nature methods 7(4), 248-249 (2010)
6. The 1000 genomes project consortium: A map of human genome variation from population-scale sequencing. Nature 467, 1061-1051 (2010)
7. Sherry, S.T., Ward, M.H., Kholodov, M. et al.: dbSNP: the NCBI database of genetic variation. Nucleic Acid Research 29, 308-311 (2001)