

Lung Precancer Analysis

Jaewoong Lee Semin Lee

Department of Biomedical Engineering
Ulsan National Institute of Science and Technology

jwlee230@unist.ac.kr

2021-04-13

Overview

1 Introduction

2 Materials

3 Methods

4 Results

Introduction

Lung Cancer

- Squamous cell carcinoma
- Adenocarcinoma

Precancer

Introduction

Study Objectives

Study Objectives

- Find different mutations
 - between WES
 - between WTS
 - from cancer
 - from precancer
- Pathway examine from the mutations
 - of WES
 - of RNA-seq
- Ultra-deep sequencing to find an *infinitesimal* quantity of Non-Circulating Tumor DNA
 - from blood
 - from urine
 - from bronchus
- Diagnostic performance

Materials

Lung Cancer Data

- WES + WTS
- Normal + {Primary, CIS + AIS, AAH, Dysplasia, MIA}
- Total 112 samples

Materials

Cancer Types

CIS + AIS

- Carcinoma *in situ* + Adenocarcinoma *in situ*

- Atypical adenomatous hyperplasia

Dysplasia

- Minimally invasive adenocarcinoma

Materials

Sample Count

Sample Count in WES

Sample Count in Transcriptome

Methods

Methods

Workflows

Data pre-processing for variant discovery

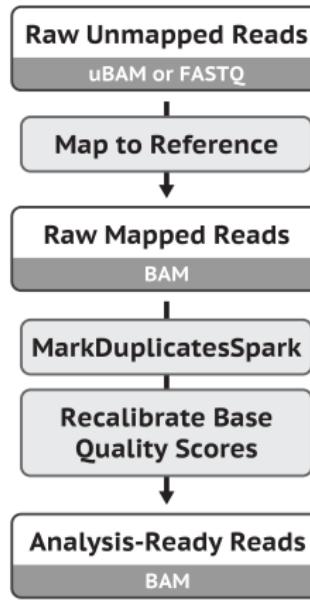


Figure: Data pre-processing for variant discovery (Van der Auwera et al., 2013; DePristo et al., 2011)

Somatic short variant discovery

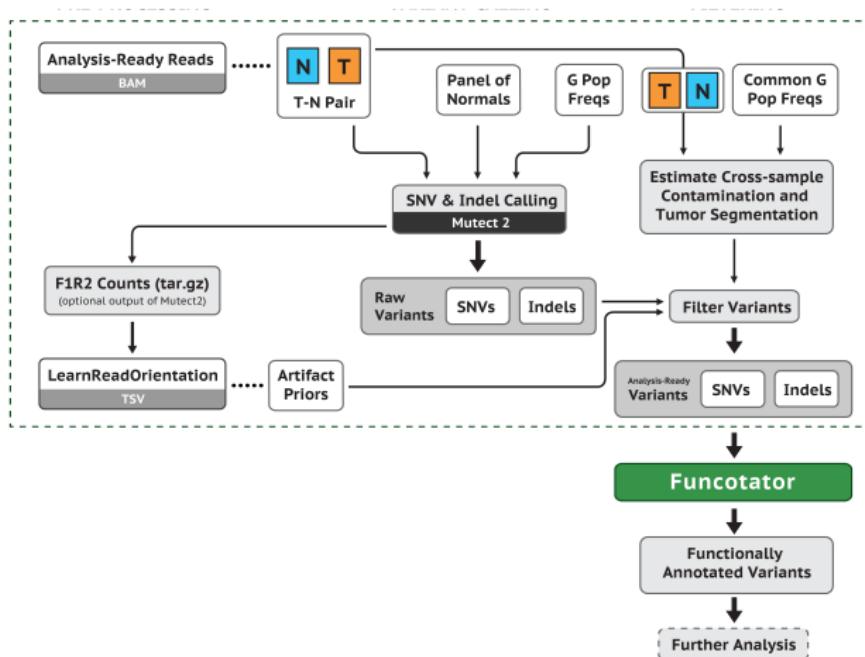


Figure: Somatic short variant (SNVs + Indels) discovery workflow (Van der Auwera et al., 2013; DePristo et al., 2011)

Germline short variant discovery

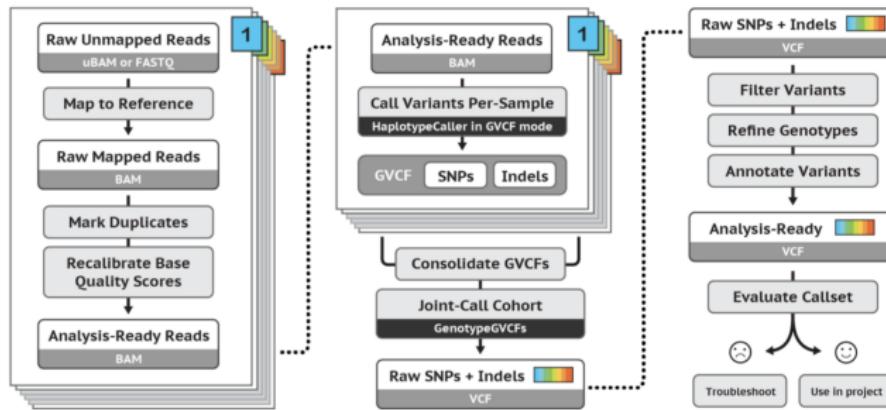


Figure: Germline short variant (SNVs + Indels) discovery workflow (Van der Auwera et al., 2013; DePristo et al., 2011)

RNA-seq short variant discovery

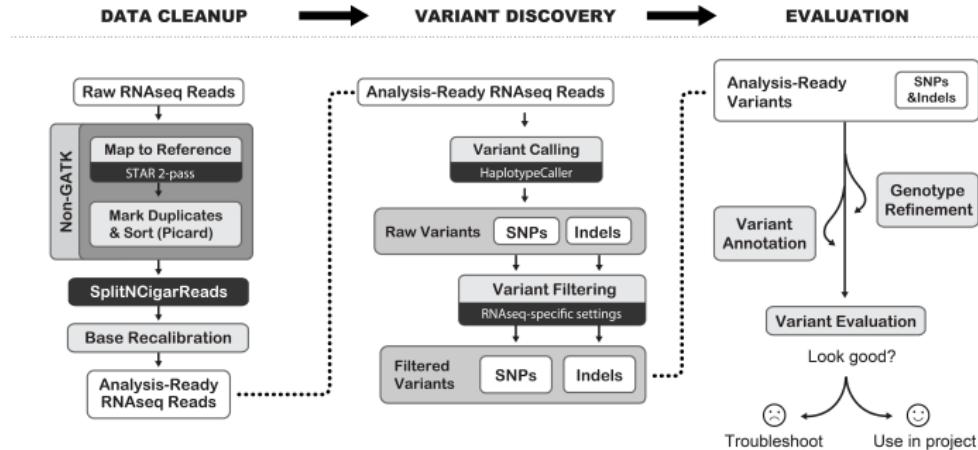


Figure: RNA-seq short variant (SNVs + Indels) discovery workflow (Van der Auwera et al., 2013; DePristo et al., 2011)

Methods

Miscellaneous

Used Bioinformatics Tools

- FastQC (Andrews et al., 2012)
- Sequenza (Favero et al., 2015)
- BWA (Li & Durbin, 2009; Li, 2013)
- STAR (Dobin et al., 2013)
- Bowtie2 (Langmead & Salzberg, 2012)
- Samtools (Li et al., 2009)
- GATK (Van der Auwera et al., 2013; DePristo et al., 2011)
- Picard (*Picard toolkit*, 2019)
- VCF2MAF (Kandoth et al., 2018)
- VEP (McLaren et al., 2016)

Python Packages

- Pandas (pandas development team, 2020; Wes McKinney, 2010)
- Sequenza-utils (Favero et al., 2015)
- Matplotlib (Hunter, 2007)
- Seaborn (Waskom & the seaborn development team, 2020)
- CoMut (Crowdis, He, Reardon, & Van Allen, 2020)

Results

Results

Quality Checks with FastQC

FastQC?

FastQC on WES

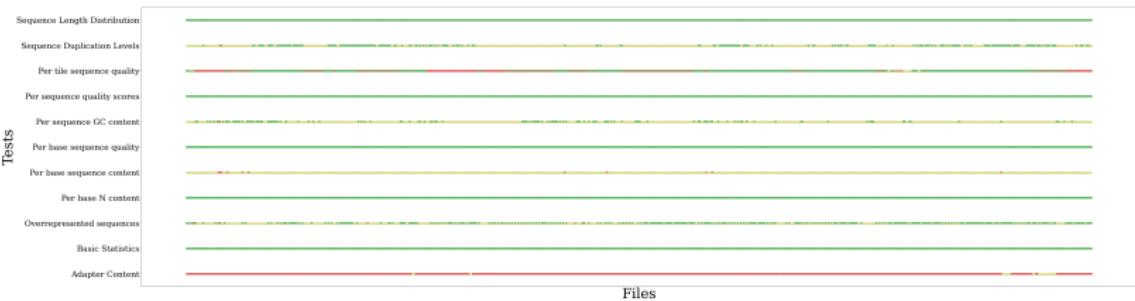


Figure: FastQC with WES data

∴ Only 33P1 has more than 3 failures: 6 FAILs.
∴ 33P1 is excluded at further analysis.

FastQC on WTS

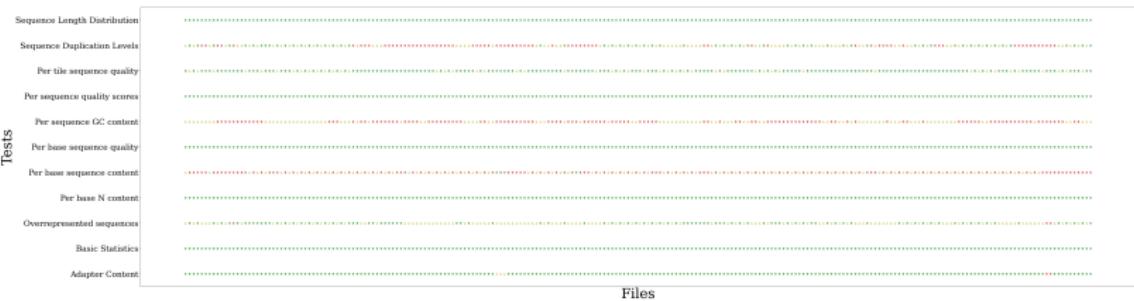


Figure: FastQC with WTS data

∴ No sample has more than 5 failures.
∴ All sample are good to analysis.

Results

Quality Checks with Sequenza

Sequenza?

Cellularity & Ploidy on WES

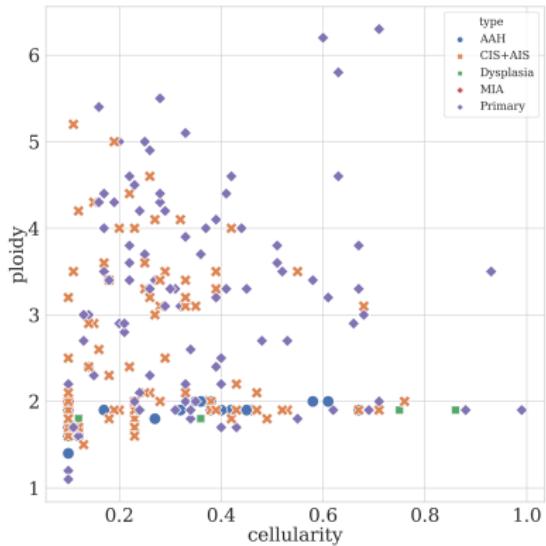


Figure: Cellularity and Ploidy from Sequenza

Copy Number Variation on WES

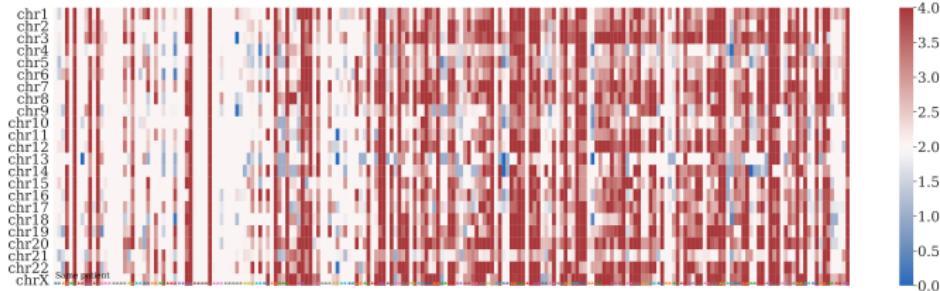


Figure: Copy Number Variation by patients and chromosomes

Results

Mutect2

Mutect2?

References I

- Andrews, S., Krueger, F., Segonds-Pichon, A., Biggins, L., Krueger, C., & Wingett, S. (2012, January). *FastQC*. Babraham Institute. Babraham, UK.
- Crowdis, J., He, M. X., Reardon, B., & Van Allen, E. M. (2020). Comut: visualizing integrated molecular information with comutation plots. *Bioinformatics*, 36(15), 4348–4349.
- DePristo, M. A., Banks, E., Poplin, R., Garimella, K. V., Maguire, J. R., Hartl, C., ... others (2011). A framework for variation discovery and genotyping using next-generation dna sequencing data. *Nature genetics*, 43(5), 491.
- Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., ... Gingeras, T. R. (2013). Star: ultrafast universal rna-seq aligner. *Bioinformatics*, 29(1), 15–21.

References II

- Favero, F., Joshi, T., Marquard, A. M., Birkbak, N. J., Krzystanek, M., Li, Q., ... Eklund, A. C. (2015). Sequenza: allele-specific copy number and mutation profiles from tumor sequencing data. *Annals of Oncology*, 26(1), 64–70.
- Hunter, J. D. (2007). Matplotlib: A 2d graphics environment. *IEEE Annals of the History of Computing*, 9(03), 90–95.
- Kandoth, C., Gao, J., qwangmsk, Mattioni, M., Struck, A., Boursin, Y., ... Chavan, S. (2018, February). *mskcc/vcf2maf: vcf2maf v1.6.16*. Zenodo. Retrieved from
<https://doi.org/10.5281/zenodo.1185418> doi:
10.5281/zenodo.1185418
- Langmead, B., & Salzberg, S. L. (2012). Fast gapped-read alignment with bowtie 2. *Nature methods*, 9(4), 357.
- Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with bwa-mem. *arXiv preprint arXiv:1303.3997*.

References III

- Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with burrows-wheeler transform. *bioinformatics*, 25(14), 1754–1760.
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., ... Durbin, R. (2009). The sequence alignment/map format and samtools. *Bioinformatics*, 25(16), 2078–2079.
- McLaren, W., Gil, L., Hunt, S. E., Riat, H. S., Ritchie, G. R., Thormann, A., ... Cunningham, F. (2016). The ensembl variant effect predictor. *Genome biology*, 17(1), 1–14.
- pandas development team, T. (2020, February). *pandas-dev/pandas: Pandas*. Zenodo. Retrieved from <https://doi.org/10.5281/zenodo.3509134> doi: 10.5281/zenodo.3509134
- Picard toolkit. (2019). <http://broadinstitute.github.io/picard/>. Broad Institute.

References IV

- Van der Auwera, G. A., Carneiro, M. O., Hartl, C., Poplin, R., Del Angel, G., Levy-Moonshine, A., ... others (2013). From fastq data to high-confidence variant calls: the genome analysis toolkit best practices pipeline. *Current protocols in bioinformatics*, 43(1), 11–10.
- Waskom, M., & the seaborn development team. (2020, September). *mwaskom/seaborn*. Zenodo. Retrieved from <https://doi.org/10.5281/zenodo.592845> doi: 10.5281/zenodo.592845
- Wes McKinney. (2010). Data Structures for Statistical Computing in Python. In Stéfan van der Walt & Jarrod Millman (Eds.), *Proceedings of the 9th Python in Science Conference* (p. 56 - 61). doi: 10.25080/Majora-92bf1922-00a