

# Identification of Potential Regulatory Elements by Multi-omics Analysis and Haplotype Phasing in Multiple Lung Adenocarcinoma Cell Lines

Graduate School of Frontier Sciences

Department of Computational Biology and Medical Sciences

Laboratory of Systems Genomics (Suzuki-Laboratory)

Sereewattanawoot Sarun ID 47-156432 Master's degree 2017 March

Keywords: Lung Cancer, Adenocarcinoma, Haplotype Phasing, Multi-omics

Academic Adviser: Professor Suzuki Yutaka

## Background

In this study, I intend to elucidate the transcriptional consequences of the somatic mutations (SNVs), which are frequently identified in potential regulatory regions in cancer genomes. While frequently identified, their functional relevance still remains elusive. For this purpose, I first attempted to identify the allelic background of SNVs in potential regulatory regions with regards to SNPs/SNVs of their corresponding transcripts. Then, with the presence of SNVs in the regulatory regions, I selected the genes which showed allelic bias in their transcript levels. In this study, I used a series of cell lines derived from Lung Adenocarcinoma as the model cases. Lung Adenocarcinoma is one of the most prominent and extensively studied cancer both in clinical specimen and cell lines. The most important driver genes, which are now being used or developed into therapeutics targets, such as EGFR, ALK fusion or KRAS, are known. However, even for this intensively studied cancer, genetic background in carcinogenesis of a large population (24.4%) of cases remains unknown. Moreover, the carcinogenic effects of even the most powerful driver genes themselves could not be solely responsible for entire carcinogenesis process. Previous work at our lab (Suzuki et al. 2014 NUCLEIC ACID RES) identified a large number of mutations, both in coding and regulatory regions in 26 Lung Adenocarcinoma-derived cell lines. In addition, the multi-omics data, such as histone modifications and transcriptomes, from the same material have been collected. With these dataset, I attempt to identify SNVs which lead to aberrant transcriptional regulations, thereby contributing to carcinogenesis.

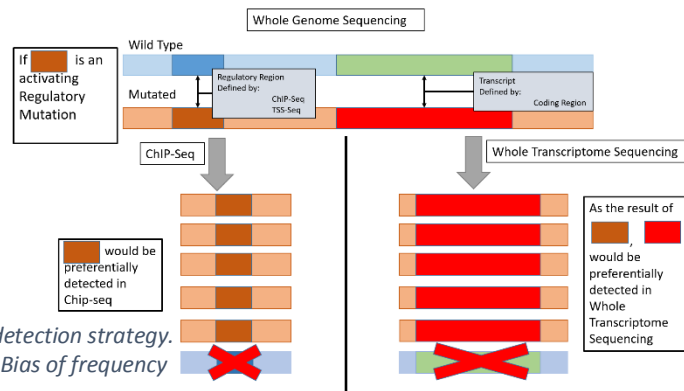
## Material and Methods

The multi-omics dataset I processed and used in this study for each cell line are as follows:

Whole genome sequencing, RNA-seq, TSS-seq and ChIP-seq data sets, ChIP-seq data include Pol-II, H3K4me1, H3K4me3, H3K9me3, H3K27me3, H3K36me3, H3K27Ac and H3K9/K14Ac, all available from previous work (Suzuki et al 2014). Whole genome sequence, RNA-seq and ChIP-seq were re-mapped to UCSC's human reference genome hg38. Regulatory regions were defined by TSS-seq and ChIP-seq. Bias in allele expression were calculated from changes in variant frequency across two-omics domains.

Synthetic long reads from 10x GemCode were available in Whole genome sequence for 2 cell lines and Agilent SureSelect V5 with regulatome regions for 23 others. The 10x GemCode data were handled by 10x LongRanger software for linked-read analysis. Final Phasing was done with my own phasing schemes.

Figure 1 show allele expression/regulatory activity bias detection strategy. Variant frequency were calculated for all omics dataset. Bias of frequency in Chip to WGS and RNA to WGS were noted



## Results



With more cross validation and control data, I believed that major improvement in phasing efficiency and accuracy should be achieved. Indeed, I found that phasing provides many crucial and unique information for SNPs/SNVs phasing and copy number alterations detection. Continuous improvement in this new field would one day turn phasing into new standard in future cancer genome sequencing.