# MAKING AN I.M.P.A.C.T
## ADVANCING THE COMPUTATION OF NEXT-GENERATION SEQUENCING DATA

Krishan Kania

## Part I: My Story

After competing in the Intel Science Talent Search, I am sometimes asked about what I have gained or learned, aside from specific academic knowledge. Well, in the process, I believe I have gained clarity pertaining to science in general.

Clarity: I've always thought of biology as the qualitative science-- that biology is applied chemistry, chemistry is applied physics, and physics is applied mathematics. And yes, from my research experience this past summer, I've experienced first hand that the tertiary and quaternary principles that explain many of our observations in biology do rest in these quantitative disciplines. However, I've also realized that these "degrees of separation" between biology and math are meaningless.
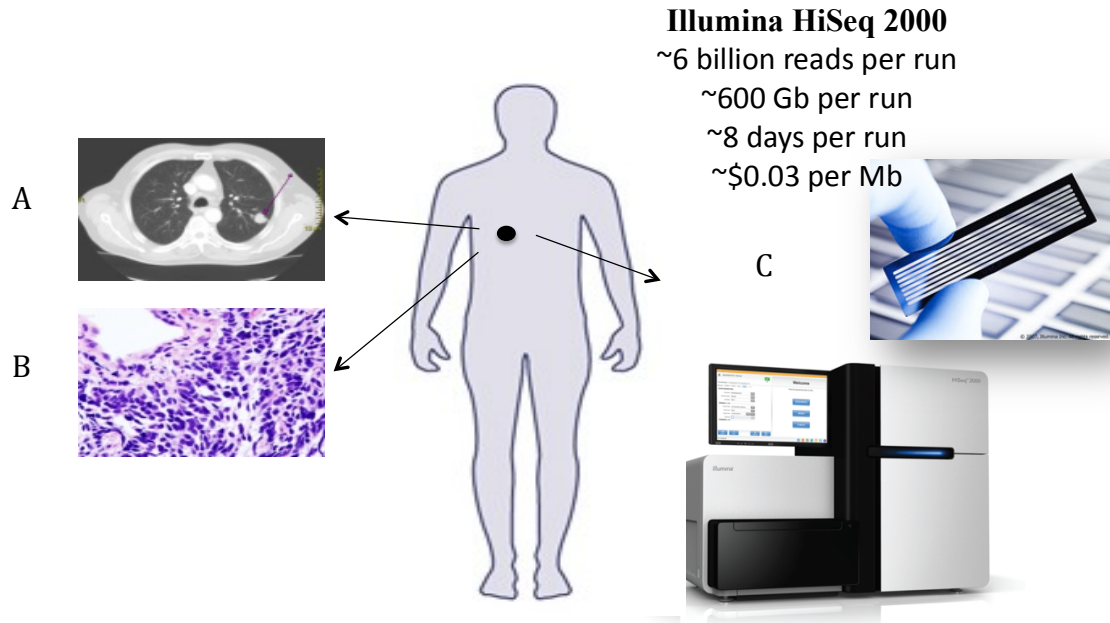
In order to accomplish my project goals, I needed to combine number theory, statistics, genetics, and computer science. My project was never to just observe tissue samples, nor was it to crunch numbers and perform calculations. It was to do both in a way that allows me to approach a question in a more comprehensive way than I have ever been exposed to in any of my previous academic experiences. My research experience has clarified the value (and necessity) of uniting all scientific and quantitative disciplines when answering a question, responding to a worthy cause.

## Part II: The Science

### INTRODUCTION:

Next-generation sequencing (NGS) has allowed substantial advances in cancer genomics. In fact, large-scale discovery efforts have propelled the identification of hundreds of cancer-related genes in recent years. To be truly transforming, however, key cancer-associated mutations must be profiled systematically in the clinical and translational arena to guide rational cancer therapeutics. This aim has yet to be achieved on a large scale, mainly because many methodologies cannot be applied efficiently and reliably on formalin-fixed paraffin embedded (FFPE) tumor samples that are routinely encountered in the clinic and in archived tumor banks. This project is a part of the computational effort to develop and apply a robust and cost-effective methodology, empowered by solution-phase exon capture and massively parallel next-generation sequencing, by which any FFPE tumor may be characterized for somatic base mutations and copy number changes in all known cancer genes. With the programming language "R," the computational analysis of NGS data for assays running clinical samples has been redeveloped, automated, and graphically represented. Moreover, such analysis, such as copy-number graphs or QC metrics, can be computed at a speed that is 568 times as fast as the traditional, and manual, computational techniques of alternative

methodologies. Furthermore, the program is built with careful considerations to make an even
mo[...]
con[...]



**Illumina HiSeq 2000**
~6 billion reads per run
~600 Gb per run
~8 days per run
~$0.03 per Mb

A

B

C

**Figure 1:** The paradigm of IMPACT Personalized Medicine[2]
(A) Radiology, (B) Pathology, (C) IMPACT

## Integrated Mutation Profiling of Actionable Cancer Targets (IMPACT)

IMPACT refers to the assay performed in this study. While NGS assays are generally similar,
there are important differences between IMPACT and competing assays of other labs. IMPACT
takes advantage of capture-based sequencing that targets a subset of the genome using "baits"
that select for specific DNA sequences. Scientists can lightly sequence the entire genome or
exome; this is useful in discovery projects. IMPACT, however, uses baits (Roche NimbleGen) to
only "pull down" the 275 proto-oncogenes, tumor-suppressor genes, and any other genes
involved in cell growth or division (Figure 2). This allows for deep coverage of the genes that
other labs have defined as most relevant to cancer20. By sequencing only targeted regions of
the genome, this technique not only allows for deep coverage of key genes, but also, can detect
low frequency mutations that occur in heterogeneous tumors or impure samples21-24.
Furthermore, capture-based sequencing data can be used to identify structural rearrangements
when at least one of the breakpoints is located in a targeted region25. On average, IMPACT
provides 700-1200 reads of coverage.

| ABL1 | CBLC | DNMT1 | FGFR1 | IGF1R | MDM2 | NOTCH2 | PNRC1 | SPOP |
| ABL2 | CCND1 | DNMT3A | FGFR2 | IGFBP7 | MDM4 | NOTCH3 | PREX2 | SRC |
| AKT1 | CCNE1 | DNMT3B | FGFR3 | IKBKE | MEN1 | NOTCH4 | PRKAR1A | STK11 |
| AKT2 | CD79B | EGFR | FGFR4 | IKZF1 | MET | NPM1 | PRKCI | SUFU |
| AKT3 | CDC42EP2 | EIF4EBP1 | FH | INSR | MITF | NRAS | PTCH1 | TBK1 |
| ALK | CDC73 | EP300 | FLCN | IRS1 | MLH1 | NTRK1 | PTEN | TEK |
| ALOX12B | CDH1 | EPHA3 | FLT1 | IRS2 | MLL | NTRK2 | PTPN11 | TERT |
| APC | CDK4 | EPHA5 | FLT3 | JAK1 | MLL2 | NTRK3 | PTPRD | TET1 |
| AR | CDK6 | EPHA6 | FOXL2 | JAK2 | MLL3 | PAK7 | PTPRS | TET2 |
| ARAF | CDK8 | EPHA7 | GATA1 | JAK3 | MLST8 | PARK2 | RAF1 | TGFBR2 |
| ARHGAP26 | CDKN2A | EPHA8 | GATA2 | JUN | MPL | PARP1 | RARA | TMPRSS2 |
| ARID1A | CDKN2B | EPHB1 | GATA3 | KDM5C | MSH2 | PAX5 | RB1 | TNFAIP3 |
| ASXL1 | CDKN2C | EPHB4 | GNA11 | KDM6A | MSH6 | PBRM1 | REL | TOP1 |
| ATM | CEBPA | EPHB6 | GNAQ | KDR | MTOR | PDGFRA | RET | TP53 |
| ATRX | CHEK1 | ERBB2 | GNAS | KEAP1 | MYB | PDGFRB | RICTOR | TP63 |
| AURKA | CHEK2 | ERBB3 | GOLPH3 | KIT | MYC | PHOX2B | RPTOR | TSC1 |
| BAP1 | CREBBP | ERBB4 | GRIN2A | KLF6 | MYCL1 | PIK3C2G | RUNX1 | TSC2 |
| BCL2L1 | CRKL | ERG | GSK3B | KRAS | MYCN | PIK3CA | SDHB | TSHR |
| BCL6 | CRLF2 | ESR1 | HDAC2 | LDHA | NCOA2 | PIK3CB | SETD2 | VHL |
| BIRC2 | CSF1R | ETV1 | HIF1A | LGR6 | NF1 | PIK3CD | SHQ1 | WT1 |
| BRAF | CTNNB1 | ETV6 | HMGA2 | MAGI2 | NF2 | PIK3CG | SMAD4 | YAP1 |
| BRCA1 | CYLD | EZH2 | HNF1A | MAP2K1 | NFE2L2 | PIK3R1 | SMARCA4 | YES1 |
| BRCA2 | DAXX | FAM123B | HRAS | MAP2K2 | NFKB1 | PIK3R2 | SMARCB1 | |
| CARD11 | DDR2 | FAM46C | HSP90AA1 | MAP2K4 | NFKB2 | PIK3R3 | SMO | |
| CBL | DICER1 | FAS | IDH1 | MAP3K8 | NKX2-1 | PKM2 | SOCS1 | |
| CBLB | DIS3 | FBXW7 | IDH2 | MCL1 | NOTCH1 | PLK2 | SOX2 | |

**Figure 2:** IMPACT panel of captured genes.

<u>MATERIALS & METHODS:</u>

## 2.1 Targeted Sequencing Methodology in IMPACT

Once Formalin-fixed paraffin embedded (FFPE) tumor tissue is obtained, the genomic DNA is extracted and sheared to a mean fragment length of 200-300 base pairs. Adaptors containing sequencing primer sites and a unique barcode are ligated to the ends of DNA fragments to create a sequencing library (approximately 24 barcoded libraries are combined in an equimolar pool). These libraries are hybridized in solution to biotinylated capture oligonucleotides (baits) complementary to the exons of 275 cancer genes. Captured DNA is enriched via streptavidin-coated magnetic beads and eluted. The DNA is then sequenced on one lane of an Illumina HiSeq 2000. After QC metrics and other metric analysis of Illumina's fastQ file, the sequence reads are aligned to the reference human genome, and target genes are examined for mutations, InDels, copy number alterations, and rearrangements (figure 3).



**Figure 3:** Overview of IMPACT
(A) IMPACT uses a hybrid capture method to sequence multiplexed libraries of 12-24 samples. (B) Reads are aligned to the reference genome and can be visualized by the Integrated Genomics Viewer. (C) This method allows
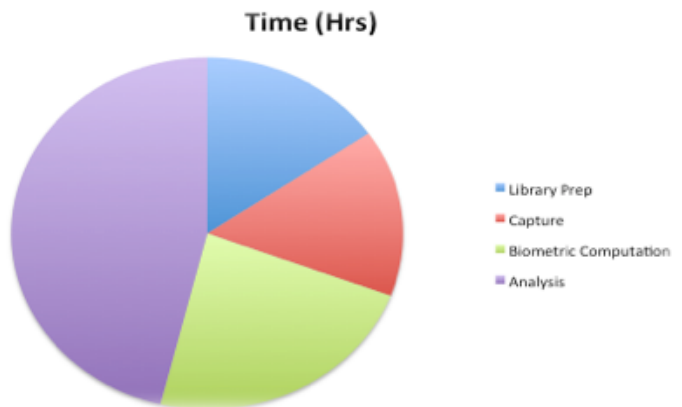
## 2.2 Metrics & Quality Control (QC):

To become familiar with the computational side of IMPACT, with the instruction of my mentor, I performed the following calculations from raw data for ovarian, colorectal, and melanoma projects of collaborating labs. These experiences would later expand into this project, where I will go on to develop these metrics into an R-scripted program that relates each calculation with graphical support:

**Metrics:**

*Cluster Density & Alignment Rate*
*Base Quality Scores*
*Insert Size Distributions*
*Fingerprints*
*Contamination*
*Capture Specificity*
*Library Complexity*
*Mean Target Coverage*
*GC Loess Normalization*

**Figure 5:**



Time (Hrs)

- Library Prep
- Capture
- Biometric Computation
- Analysis

**Context:**

All metrics are derived from fundamental concepts in biology, chemistry, physics, and statistics. For example, "Fingerprints" is a metric that checks the alleles of a tumor/normal pair at 42 sites of single-nucleotide polymorphism (SNP). SNP is a single nucleotide base-pair site where variation is found in at least 1% of the population[26]. Specifically, the SNP's we capture are in tiling regions, regions of the genome very close to the 275 cancer genes. When 38 out of 42 SNP sites (arbitrary threshold) match between tumor and normal tissue, there is confidence that the tumor is paired with its matched normal. "Fingerprints" allows an investigator to identify contamination, sample swamps, and even loss of heterozygosity. GC content is the percentage of nitrogenous bases on a DNA molecule in a particular region that are either guanine or cytosine.

$$GC\ Content = \frac{G + C}{A + T + G + C} \times 100$$

A GC pair is bound by three hydrogen bonds, compared to AT pairs, which are bound by only two hydrogen bonds. A high GC-content implies a higher annealing temperature and higher melting temperature in PCR experiments. Moreover, high GC-content often implies major technology-related artifacts and biases due to the weakness of sequencing technology. Since levels GC bias is varied across the genome, the GC effect can be hard to tell apart from the true signal[27]. Even more challenging, the effect is not consistent between repeated experiments, or even libraries within the same experiment. Unsurprisingly, estimating and directly correcting for this effect has become a well-established step in protocol design. Normalization is therefore

essential to ensure accuracy, particularly in GC rich regions; the statistical method of choice is the Loess model. GC Loess Normalization will be incorporated into the R-Scripted program.

Traditionally, the two biggest limitations of metric computation in IMPACT and similar assays are speed and confidence. In a typical project, these metrics require hours of valuable technician time to do rudimentary calculations or manipulations. As a corollary, this process is not as comprehensive as it can be, resulting in more scenarios where mutation calls are questionable and thus subject to primer evaluation. Speed and confidence, arguably the two most important concerns of a clinical setting, will be better addressed with the R-scripted program.

## 2.3 Mutation Calling & Primer Evaluation

Following bench work, sequencing, and metric computation, scientists are now ready to confidently discuss genomic mutations:

### *Copy-Number Alterations:*

IMPACT can determine the pattern of copy-number alterations: the gain/loss of chromosome arms or focal amplifications and deletions that might range from tens of kilobases to tens of megabases in size. In a normal genome, across the autosomes, there should be two copies of every gene, one maternal and one paternal. However, in the cancerous genome there is usually a gain or loss of one or several exons. Of course, these events have important therapeutic implications. Copy-number was traditionally visualized in Microsoft Excel in manually generated graphs depicting the tumor/normal ratio of exon reads between a given tumor and its matched normal[28]. The limitations of this strategy will be revisited in results.

### *Somatic nucleotide substitutions and small insertion and deletion mutations:*

Nucleotide substitution mutations are the most frequent somatic genomic alteration in cancer, occurring at the rate of about one somatic nucleotide substitution per million nucleotides; insertion and deletion events are approximately tenfold less common[29]. The detection of somatic mutations in cancer requires mutation calling in both the tumor DNA and the matched normal DNA, coupled with comparison to a reference genome and an assessment of the statistical significance of the number of counts of the mutation in the cancer sequence and its absence in the matched normal sequence[30].

*Primer Evaluation:*

The validation [text obscured] PCR amplification of [text obscured] surrounding the [text obscured] followed by Sanger sequencing.



**Figure 6:** Suspected substitution event in PIK3CD shown in Integrated Genomic Viewer (IGV) is confirmed with primer evaluation.

*Primer Development:*

A length of 18-25 bases

%GC content between 40 and 60%

Tm (melting temperature) between 50 and 60°C

No secondary priming sites (BLAT)

No dimerization capability

No significant hairpins (> 3 bp)



TCGA-13-1481   PIK3CD   p.R184W

Forward Primer

Hypothetical Substitution in PIK3CD

Reverse Primer

Melting point forward: 59.15°C
Melting point reverse: 60.45°C
Amplicon length: 307
BLAT for repetitive regions: Pass
No SNP's

**Figure 7:** Choosing the best primer to evaluate the substitution in PIK3CD

## 2.4 "R" Package & R-Scripted Programming

R is a programming language widely used among biostaticians and is highly relevant in bioinformatics (gene expression data, serial analysis of gene expression, etc.). R is a useful tool for plotting graphics, analyzing data, and fitting data to statistical models. It is open source, free, and maintained by a team of developers around the world.      Traditionally, the metrics, QC metrics, and copy-number analysis of IMPACT were done on Microsoft Excel; however, they are now all redeveloped, improved, and automated in R-Studio. *R Cookbook*, *R Graphics Cookbook,* and *R in a Nutshell* have all been very helpful in writing code and integrating it into Next-Generation Sequencing assays, including IMPACT[31,32,33].

## RESULTS:

To illustrate the performance of the R-scripted program to compute metrics and copy number alterations, results will be presented for a project to identify biomarkers of metastasis in colorectal cancer. DNA was sequenced from primary tumors, metastatic tumors, and matched normal blood from 36 patients. Nevertheless, this program can and should be applied to enable discovery across virtually all cancers studied with NGS.

### 3.1 Visualization, Automation, and Development of Metric Calculations

The R-scripted program relates each metric computation to a comprehensible graph that will quickly allow investigators to match tissues, identify potential contamination, evaluate the performance of the sequencing technology, and document quality control (figure 8).

**Figure 8:** An Automated and more developed QC Metric Computation on colorectal cancer NGS data

(A) Cluster Density & Alignment Score, (B) Capture Specificity, (C) Base Quality Score, (D) Insert Size Distribution, (E) Fingerprint Matching, (F) Contamination, (G) Estimated Duplication Rate, (H) Estimated Library Size. These calculations (and graphs) are now automatically computed from raw data with the R-Scripted program (I). This program also offers previously uncalculated metrics such as mean exon coverage (in addition to mean target coverage). Most importantly, the program is universally functional across any permutation or combination of up to 24 tumors/normals/etc. prioritizing speed and automation across virtually all laboratory scenarios encountered in IMPACT or the other targeted cancer assays in published literature.

## 3.2 Automated Copy-Number Plotting with "Normal to Normal", "Tumor to all Normals", and "Tumor to Tumor-Median" Capability

IMPACT is run on a diverse set of projects, which include the conventional tumor to normal comparisons, but also, primary tumor to metastatic tumor to normal comparisons (for discovery projects on tumor evolution), tumor to all tumor comparisons, and other variations, which have severely limited automation programs in the past. However, this project, with the tools of R, is built with the flexibility to approach virtually all variations of NGS on cancerous tissue.

Furthermore, with the computational power of R, copy-number graphs exist for each normal compared to all normals, and each tumor compared to all normals. On Microsoft Excel, it would be possible to make these comparisons, however, this would make an already long process even longer.

Manual output from
Microsoft Excel:



Automated and more
informative output
from "R":



**Figure 9:** Advantage of Copy-Number Plotting with R compared to competing methodologies

results in a "messy" copy-number graph, where amplifications or deletions are not clear. (This usually occurs due to machine artifacts during sequencing, differences in DNA quality, or differences in input size during library preparation). With the R-scripted program, tumor tissue is automatically compared to not only the matched normal, but also, all other normals, often times resulting in a cleaner tumor/normal copy-number graph. It is scientifically sound to compare a given tumor to an unmatched normal because in theory, all normal samples should have two copies of every exon in the twenty-two autosomes sequenced in IMPACT. If there is suspicion that there is a germ-line amplification or any contamination in the normal compared to the given tumor, the investigator can always consult the normal to all normal comparisons.

In figure 9, a traditional pipeline using Microsoft Excel was too messy to provide any meaningful conclusions for this patient with colorectal cancer. However, with R, it is clear that there is a EGFR amplification, an actionable cancer target, on the 7th chromosome. The copy-number is cleaner on R because the program is actually comparing the tumor to all normals in the pool, rather than just the matched normal, allowing the investigator to choose the most informative graph. After comparing this "miracle" normal to other normals in the pool, the normal-to-normal comparison indicates that even though the normal does not match the tumor, it does not suffer from any autosomal bias, and is perfectly legitimate (Figure 10). In a prospective setting, detecting the EGFR amplification, with the tumor to all normals comparison and the added confidence of the normal to all normals comparison, would open important doors to

**9**

targeted therapies such as cetuximab, gefitinib, erlotinib, and panitumumab. In addition, in the event that there are no normals in the pool, the program is prepared to compare the given tumor to the median exon values across all tumors.



**Figure 10:** Relatively quiet copy-number comparison between the normal used in Figure 7 and another normal in the given pool.

## 3.3 Comparison to The Cancer Genome Atlas (TCGA)

Interestingly, these computational efforts, compounded with the existing experimental design of IMPACT, has enabled mutation calling that does not only match, but occasionally, even surpasses the NGS assay of The Cancer Genome Atlas (TCGA). When frozen and FFPE tissue first screened by the TCGA[34] was run on IMPACT, the investigators using IMPACT called all 17 mutations found by TCGA, and 8 additional mutations not found by TCGA. These mutations were all confirmed with sanger sequencing (primer evaluation).



**Figure 11: Comparison with TCGA**
A comparison of 6 frozen ovarian tumors sequenced by IMPACT and TCGA revealed that (A) all 17 mutations found by TCGA in IMPACT genes were detected by IMPACT. Additionally, 8 mutations not found by TCGA were detected by IMPACT, as seen in (B) IGV screenshots. (C) These mutations were at low frequency in both tumors but (D) detected due to higher coverage.

**10**

## DISCUSSION:

Automating the metric computation of NGS has provided an easier transition from the raw data output of sequencing (Illumina HiSeq 2000 FastQ file) to mutation calling. The programing language R has served to bridge the gap in this transition, replacing the traditional, and less robust, approach of Microsoft Excel and the other competing programs. While many labs performing NGS are not core-facilities, they can still benefit from the computational and statistical power of a well-scripted program.

**Speed:** To begin with, for IMPACT, as well as other assays, this project guarantees speed. On average, running 12 patients (24 tissue samples) on IMPACT implies at least 3 hours behind the computer or calculator performing metric computation or plotting copy-number graphs. This process, on average, now takes just 19 seconds with the R-Scripted program. Thus, an investigator can shift laboratory resources from computation to interpretation, while performing critical quality controls checks to identify artifacts that could lead to false positive mutation calls or spurious conclusions.



**Figure 12**: Average time spent on metric computation with R $\approx$ **19 seconds** ($\mathbf{CPU = 2.2\ GHz, RAM = 4\ Gb}$)

**Confidence:** The R-Scripted program has not only made the assay faster, but also, more informative. In addition to the traditional metric calculations of the existing pipeline, the R-Script provides 100% bar plots, average exon coverage, normal to all normal copy number, and tumor to all normal copy number data. With these additional resources, labs can be more confident in their calls as they move into Integrated Genomic Viewer (IGV) and rate mutations. This will also allow for more informed choices when contemplating primer evaluation on a questionable mutation.

## FUTURE WORK:

One important limitation of the R-Scripted program is that the investigator must visually inspect the copy number graphs and then manually choose the most informative one. The program produces graphs that compare one tumor to every normal in a given pool. Currently, all graphs are displayed in the PDF, including those that are not particularly informative. For example, in a pool of 12 patient samples, the investigator would be looking for the 12 cleanest tumors to normal graphs from a PDF of 78 graphs. While it is important to preserve this element of human intuition, I am experimenting with computational methods that can output graphs that would match visual inspection.

One method that is currently being tested is a square-adjacent exon calculation: the addition of the square of the difference in tumor/normal ratio between adjacent exons across all exons, where the smallest sum relates to the most informative graph.

$$\sum_{n=1}^{4655} (exon_{n+1} - exon_n)^2$$

*There are 4656 exons captured with the current version of baits in IMPACT

Taking the square, will exaggerate the sum of graphs with particularly messy copy numbers more than simply taking the sum of an absolute value of the difference between adjacent exons. The "R" script could then easily select to output (PDF) only the graph with the smallest sum in a for-loop for each tumor.

Another option is to segment the copy number data using circular binary segmentation via DNAcopy, an "R" package developed by BioConductor that has resolved similar issues.

## CONCLUSION:

The establishment of the experimental and computational efforts of IMPACT will have immediate, far-reaching benefits for translational and clinical research and will provide the foundation for personalized cancer medicine. Systematic profiling of every cancer gene in tumor DNA from every cancer patient would improve diagnosis and reveal the spectrum of alterations across tumor types, the presence of mutations with potential therapeutic implications in unexpected contexts, and their patterns of co-occurrence that might direct treatment choice. Profiling these same genes retrospectively across a vast collection of clinically annotated FFPE tumors would enable the discovery of significant oncogenic mutations in rare or understudied tumor types and the identification of genomic biomarkers exhibiting correlations with clinical outcomes or phenotypes in every cancer. These efforts collectively embody the goal to produce better outcomes in cancer patients and make cancer a more manageable disease.

## ACKNOWLEDGEMENTS:

# REFERENCES

[1] MacConaill, L. E., & Garraway, L. A. (2010). Clinical implications of the cancer genome. *Journal of Clinical Oncology*, 5219-5228.

[2] Wagle, N., & Berger, M. F. (2012). High-throughput detection of actionable genomic alterations in clinical tumor samples by targeted, massively parallel sequencing. *Cancer Discovery*, 83-93.

[3] Bentley, D., Balasubramanian, S., & Swerdlow, H. (2008). Accurate whole human genome sequencing using reversible terminator chemistry. *Nature*, 53-59.

[4] Drmanac, R. *et al.* (2010). Human genome sequencing using unchained base reads on self-assembling DNA nanoarrays. *Science*, 78-81.

[5] Marguilies, M. *et al.* (2005). Genome sequencing in microfabricated high-density picolitre reactors. *Nature*, 376-380.

[6] Shendure, J. *et al.* (2005). Accurate multiplex polony sequencing of an evolved bacterial genome. *Science*, 1728-1732.

[7] Wheeler, D. *et al.* (2008). The complete genome of an individual by massively parallel DNA sequencing. *Nature*, 872-876.

[8] Ley, T. J. *et al.* (2008). DNA sequencing of a cytogenetically normal acute myeloid leukemia genome. *Nature*, 66-72.

[9] Mardis, E. R. *et al.* (2009). Recurring mutations found by sequencing an acute myeloid leukemia genome. *New England Journal of Medicine*, 1058-1066.

[10] Pleasance, E. D. *et al.* (2010). A comprehensive catalogue of somatic mutations from a human cancer genome. *Nature*, 191-196.

[11] Weir, B., Zhao, X., & Meyerson, M. (2004). Somatic alterations in the human cancer genome. *Cancer Cell*, 433-438.

[12] Mitsudomi, T. *et al.* (2009). Gefitinib versus cisplatin plus docetaxel in patients with non-small-cell lung cancer harboring mutations of the epidermal growth factor receptor: an open label, randomized phase 3 trial. *Lancet Oncology*, 121-128.

[13] Rosell, R. *et al.* (2009). Screening for epidermal growth factor receptor mutations in lung cancer. *New England Journal of Medicine*, 958-967.

[14] Campbell, P. J. *et al.* (2008). Identification of somatically acquired rearrangements in cancer using genome-wide massively parallel paired-end sequencing. *Nature Genet.*, 722-729.

[15] The Cancer Genome Atlas Research Network (2008). Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature*, 1061-1068.

[16] Wood, H.M. *et al.* (2010). Using next-generation sequencing for high-resolution multiplex analysis of copy number variation from nanograms quantities of DNA from formalin-fixed paraffin-embedded specimens. *Nucleic Acids*, 38.

[17] Gallegos, R.M. *et al.* (2007). EGFR and K-Ras mutation analysis in non-small cell lung cancer: comparison of paraffin embedded versus frozen specimens. *Cell Oncology*, 257-264.

[18] Meyerson, M., Gabriel, S., & Getz, G. (2010). Advances in understanding cancer genomes through second-generation sequencing. *Nature Reviews: Genetics*, 685-691.

[19] Carter, H. *et al.* (2009). Cancer-specific high-throughput annotation of somatic mutations: computational prediction of driver missense mutations. *Cancer Res.*, 6660-6667.

[20] Kerick, M., Isau, M., & Timmermann, B. (2011). Targeted high throughput sequencing in clinical cancer settings: formaldehyde fixed-paraffin embedded (FFPE) tumor tissues, input amount and tumor heterogeneity. *Medical Genomics*, 68-81.

[21] Blumenstiel, B., Kristian, C., & Fisher, S. (2010). Targeted exon sequencing by in-solution hybrid selection. *Current Protocols in Human Genetics*, 66-84.

[22] Fisher, S., Barry, A., & Abreu, J. (2011). A scalable, fully automated process for construction of sequence-ready human exome targeted capture libraries. *Genome Biology*, 12-27.

[23] Meyer, M., & Kircher, M. (2010). Illumina sequencing library preparation for highly multiplexed target capture and sequencing. *CSH Protocols*, 1-9.

[24] Parkinson, N., Maslau, S., & Ferneyhough, B. (2012). Preparation of high-quality next-generation sequencing libraries from picogram quantities of target DNA. *Genome Research*, 125-133.

[25] Campbell, P. J. *et al.* (2008). Identification of somatically acquired rearrangements in cancer using genome-wide massively parallel paired-end sequencing. *Nature Genet.*, 722-729.

[26] Ramensky, V., & Bork, V. (2002). Human non-synonymous SNPs: server and survey. *Nucleic Acids*, 3894-3900.

[27] Aird, D., Ross, M., & Chen, W. (2011). Analyzing and minimizing PCR amplification bias in illumina sequencing libraries. *Genome Biology*, 1-14.

[28] Beroukhim, R. *et al.* (2010). The landscape of somatic copy-number alteration across human cancers. *Nature*, 899-905.

[29] Stephens, P.J. *et al.* (2009). Complex landscapes of somatic rearrangement in human breast cancer genomes. *Nature*, 1005-1010.

[30] Thomas, R.K. *et al.* (2007). High-throughput oncogene mutation profiling in human cancer. *Nature Genetics*, 347-348.

[31] Teetor, P. (2011). *R cookbook*. (1st ed.). San Francisco, CA: O'Reily.

[32] Adler, J. (2009). *R in a nutshell*. (1st ed.). San Francisco, CA: O'Reily.

[33] Mittal, H. (2011). *R graphs cookbook*. (1st ed.). Birmingham, UK: Packt Publishing

[34] Cancer Genome Atlas Research Network. (2011). Integrated genomic analysis of ovarian carcinoma. *Nature*, 609-615.