

Master's Thesis

Implementation and Comparative Assessment of Diagnostic Cancer Gene Panels in the Molecular Pathology Laboratory

University of Luxembourg

Faculty of Science, Communication and Technology

Master in Integrated Systems Biology

by

Ben Flies

(010081174D)

Abstract

Contents

List of Abbreviations	i
List of Figures	ii
List of Tables	iii
1 Introduction	1
1.1 Targeting Cancer	1
1.2 Targeting the EGFR Pathway in Solid Tumors	3
1.3 Targeted Sequencing	3
1.3.1 Target Enrichment Methods	3
1.3.2 Illumina Sequencing Chemistry	3
1.4 NGS Data Analysis	3
1.4.1 GATK Best Practices	3
1.5 Practical Implications in the Laboratory	3
1.6 Aims of the Thesis	3
2 Material and Methods	3
2.1 Library Preparation	4
2.1.1 Patients	4
2.1.2 DNA Extraction, Quantification and Quality Control	4
2.1.3 Agilent Haloplex ClearSeq Cancer	4
2.1.4 Illumina TruSight Tumor 15	4
2.2 Bioinformatic Analysis	4
2.2.1 Agilent SureCall	4
2.2.2 Illumina BaseSpace TruSight Tumor 15 App	4

2.2.3	Custom In-House Pipeline (Velona)	5
2.2.4	Variant Calling Algorithms	5
3	Results	6
3.1	Sample Preparation	6
3.2	NGS Data Quality	6
3.3	Coverage Analysis	9
3.4	Variant Calling Algorithm Comparison	11
3.4.1	Detection of Known Single Nucleotide Variants and Deletions	11
3.4.2	Sensitivity Analysis	11
4	Conclusions	13
	References	14

List of Abbreviations

NGS Next Generation Sequencing
LNS Laboratoire National de Sante
SGMB Service of Genetics and Molecular Biol-
ogy TST15

Illumina

TruSight

Tumor

15

List of Figures

1	Electropherograms of representative sequencing libraries prepared by Agilent Halo-plex ClearSeq Cancer and Illumina TruSight Tumor 15. (*) represents the lower marker, (**) represents the upper marker	6
2	Scatter plot of the corrected peak area (X axis) of the regions corresponding to the sequencing libraries defined in the blabla software and the dCt (Y axis). Agilent Halo-plex ClearSeq Cancer data are represented as blue dots, Illumina TruSight Tumor 15 data are represented as red dots.	7
3	Comparison of coverage distributions per amplicon as reported by FastQC	8
4	Comparison of Coverage Distributions per Amplicon	9
5	Blablabla	12

List of Tables

1	ISV	7
2	samtools _{flagstat}	8
3	failed _{halo}	9
4	failed _{halo}	10
5	failed _{halo}	12
6	sensitivity _{analysis}	13

1 Introduction

Cancer represents a huge burden for health care systems worldwide and one of the leading death causes. In 2012, there were an estimated 14.1 million new cancer cases with estimated 8.2 million cancer deaths [3]. Lung cancer is the most common cancer, both in terms of new cases (1.8 million) and deaths (1.6 million). Breast cancer is the second most common cancer (1.7 million cases) but only ranks 5th as cause of death (522,000 deaths). Colorectal cancer (1.4 million cases; 694,000 deaths), prostate cancer (1.1 million cases; 307,000 deaths), stomach cancer (951,000 cases; 723,000 deaths) and liver cancer (782,000 cases; 723,000 deaths) are following.

Scientific discoveries in the last decade have had an enormous impact on our understanding of the underlying causes of cancer. The development of omics techniques, in combination with enhanced computational power, has lead to an explosion of biological data. It has become clear that cancer is an incredibly complex malignancy. The research community is trying to interpret this vast amount of data with the goal to get a deeper understanding of cancer and to cure it eventually. In recent years, several drugs have been approved that target proteins needed for cancer development, proliferation or metastasis. Molecular testing is employed to check whether these targeted drugs would be of benefit. In that regard, Next-Generation Sequencing (NGS) is an interesting method to gain deep insights into the genetic information of a tumor and to guide personalized therapy.

1.1 Targeting Cancer

Cancers are characterized by several properties known as cancer hallmarks, which comprise blablabla

In the last seventy years, five main models have dominated cancer research:

- The discovery of the carcinogenic potential of tobacco smoke lead to the the first and oldest model ('mutational'): long-term exposure to chemical carcinogens such as PAH or tobacco smoke induces cancerogenesis. The importance of viral DNA sequences and mutations induced by bacteria is also covered by this first model. In the induced mutational model, cancer is caused by external factors.
- The second model ('genome instability') puts emphasis on genome integrity, tumor-suppressor genes and oncogenes and DNA mismatch repair. Genomic instability can affect the genetic information at the level of nucleotides, microsatellites, whole genes or chromosomes. For instance, chromosome rearrangements, chromosome number alterations, loss of heterozy-

gosity (LOH) or gene amplification contribute to chromosomal instability (CIN), which occurs in 50–85% of colorectal cancers (CRCs). Even though numerous somatic mutations occur in tumors, only a small subset contributes to tumor progression. These 'driver' mutations

- Third model ('non-genotoxic') emphasizes on several important factors of cancer risk (obesity, diet, activity level, hormones). These factors do not induce structural changes, but rather act by inducing functional changes through epigenetic events such as histone methylation or DNA acetylation. Cancer then develops through selective selection of cells that have acquired a proliferative advantage through these processes.
- The fourth model ('Darwinian') is also based on clonal expansion, but puts emphasis on the macro- and micro-environment.
- Fifth model ('tissue organization'):

Cancer has been traditionally typified by a stepwise accumulation of mutations in key oncogenes and tumor suppressors. For decades, accumulation of these traits in somatic cells has been considered as the foundation of a developmental model of tumor progression where cells transition from a normal, healthy state to pre-malignant, malignant, and migratory phenotypes.

Meanwhile, tumors are often described as heterogeneous, owing to the intricate genetic diversity and assorted morphological phenotypes they embody [2]. Intratumor heterogeneity specifically refers to heterogeneity within a tumor, while intertumor heterogeneity refers to heterogeneity across several different tumors [3]. The current view of tumor heterogeneity recognizes basic principles of Darwinian evolution at the core of neoplastic development and outgrowth: a single somatic cell with a heritable fitness-promoting mutation proliferates, conferring a survival advantage that allows cells to outlast the less 'fit' cells [3, 4]. Natural selection leads to sequential waves of clonal expansion, resulting in various subclones with differing capacities for proliferation, migration, and invasion [5].

Advances in next-generation sequencing techniques and the inception of The Cancer Genome Atlas (TCGA) have revealed extensive heterogeneity at the molecular level [8]. Genetic heterogeneity of tumors is rooted in one of the key hallmarks of cancer: genetic instability [2]. Several mechanisms are in place in normal cells that protect against chromosome and nucleotide damage by preventing DNA replication until damage is repaired; however, genes controlling these critical checkpoints (e.g. p53) are often perturbed in cancer cells [16]. Genetic instability in cancer has been demonstrated at both the nucleotide level in point mutations and chromosome level in translocations, deletions, amplifications, and complete chromosome aneuploidy [17].

Tumor cells undergo a series of genetic events that contribute to genomic instability throughout tumor progression (Figure 2A). However, the specific mechanisms and precise order in which they occur have yet to be elucidated [21]. Studies have pursued these mechanisms and found that the rate at which mutations occur in somatic cells is insufficient to cause the striking number of mutations present in cancer genomes. Over the past few decades, a 'mutator' hypothesis tumor evolution has emerged, speculating that a mutator phenotype characterized by genomic instability drives multi-step carcinogenesis and explaining the mutation rate discrepancy observed in normal and malignant cells [22]. The current mutator hypothesis speculates that a small number of 'driver' alterations exist and, once acquired by somatic mutation, confer the cancer phenotype; however, seemingly insignificant 'passenger' mutations result via mechanisms yet to be elucidated [26]. McFarland et al. challenged this with stochastic simulation of tumor evolution and reasoned that, though individually weak, the cooperative burden of small-scale accumulated passenger mutations has a present role in tumor progression, and may be the cause for complex oncological events that remain unanswered by the driver-centric model [27].

Virtually all cancers tend to accumulate mutations during their progression. Studies have demonstrated that a typical cancer genome comprises about 40–80 amino acid changing mutations. Some of these mutations increase the cancer's 'fitness' over that of surrounding cells. The term 'fitness' here is defined by the difference between cell proliferation and cell death (net replication rate). Mutations that enhance the selective proliferative advantage of the cancer cell are called 'driver' mutations. It is often difficult to separate 'driver' from 'passenger' mutations, especially for low frequency variants. Passenger mutations occur in cancer cells subsequently or coincidentally to driver mutations and do not alter the cell's fitness. Typical solid tumors may contain 40–80 amino acid changing mutations, but only 5–15 of them are driver mutations.

Which genetic aberrations? Mutational inactivation of tumor-suppressor genes, activation of oncogene pathways, and then you do see solid tumors, and then you see the EGFR pathway, and then you can explain the targeted drugs.

Problem: heterogeneity resistances

1.2 Targeting the EGFR Pathway in Solid Tumors

EGFR Pathway

Targeted Drugs

1.3 Targeted Sequencing

Next generation sequencing can facilitate personalized cancer therapy approaches by identifying actionable somatic events in tumor samples (1). Furthermore, high-quality sequencing data can reveal associations with sensitivity or resistance that can inform the development and implementation of targeted therapeutics. Whole genome sequencing (WGS) and whole exome sequencing (WES) allow the detection of SNVs, indels, CNVs, and rearrangements. However, the relatively low coverage of WGS and WES, as currently implemented in most of the sequencing laboratories (100–250×), may have limited ability to cost-effectively detect aberrations that are present in a subpopulation of tumor cells while identifying a myriad of aberrations of unknown clinical significance (2). Somatic aberrations present at low allele frequencies across different types of tumors (3, 4) can potentially impact patient prognosis or response (5) and thus are important to detect reliably. Targeted sequencing to a depth that allows detection of relatively low mutant allele frequency (MAF) may represent an alternative or a complement to WGS and WES to detect clinically relevant alterations. Additionally, in most clinical and research settings, the amount of DNA that can be isolated from tumor samples is limited and the DNA is often damaged owing to fixation and storage procedures such as those used with formalin-fixed paraffin-embedded (FFPE) samples. Therefore a multiplexed targeted platform that can generate reliable data with high sensitivity from limited amounts of DNA from FFPE samples is needed. Several targeted sequencing panels have been successfully implemented (6, 7). However, the details of a platform's design and parameterization will influence the precision and reliability of the molecular profiling results, impacting both translational research and clinical decision-making. Thus, it is of great value to explore multiple potential solutions in a real patient care environment until a community-wide solution is established, validated, and well accepted.

1.3.1 Target Enrichment Methods

1.3.2 Illumina Sequencing Chemistry

1.4 NGS Data Analysis

1.4.1 GATK Best Practices

1.5 Practical Implications in the Laboratory

1.6 Aims of the Thesis

2 Material and Methods

- How was the data analyzed ?
- Present econometric/statistical estimation method and give reasons why it is suitable to answer the given problem.
- Allows the reader to judge the validity of the study and its findings.
- Depending on the topic this section can also be split up into separate sections.

2.1 Library Preparation

2.1.1 Patients

Melanoma Non-Small Cell Lung Carcinoma (NSCLC) metastatic colorectal cancer (mCRC) Chronic lymphocytic leukemia (CLL) were extracted from blood and did not undergo FFPE treatment -> were used as some kind of good quality samples to see if there are really more C₅T variants in FFPE samples.

2.1.2 DNA Extraction, Quantification and Quality Control

DNA Extraction Kit from Qiagen

Quantification Qubit fluorometer, either High Sensitivity kit or Broad Range

Quality Control Illumina Infinium FFPE QC Assay kit

2.1.3 Agilent Haloplex ClearSeq Cancer

2.1.4 Illumina TruSight Tumor 15

2.2 Bioinformatic Analysis

2.2.1 Agilent SureCall

With alignment algorithms installed Windows System, 3.00GHz, 16GB RAM\$ modifiable

2.2.2 Illumina BaseSpace TruSight Tumor 15 App

Cloud-based parameters not modifiable

2.2.3 Custom In-House Pipeline (Velona)

Linux System

2.2.4 Variant Calling Algorithms

Tested on Linux Ubuntu 14.04.4 LTS Trusty Tahr installed on VMware virtual box with of 12GB RAM

GATK HaplotypeCaller

VarScan 2

Mutect1.1.7 [2]

SomVarIUS

3 Results

3.1 Sample Preparation

Before pooling the adaptor-ligated and indexed sequencing libraries, the success of library preparation is validated using the Agilent Bioanalyzer instrument. 1a and 1b show representative electropherograms of a sample that has been processed using both kits. The expected DNA products should be detected at 175-600 bp for Haloplex CSC and 200-400 for TST15.

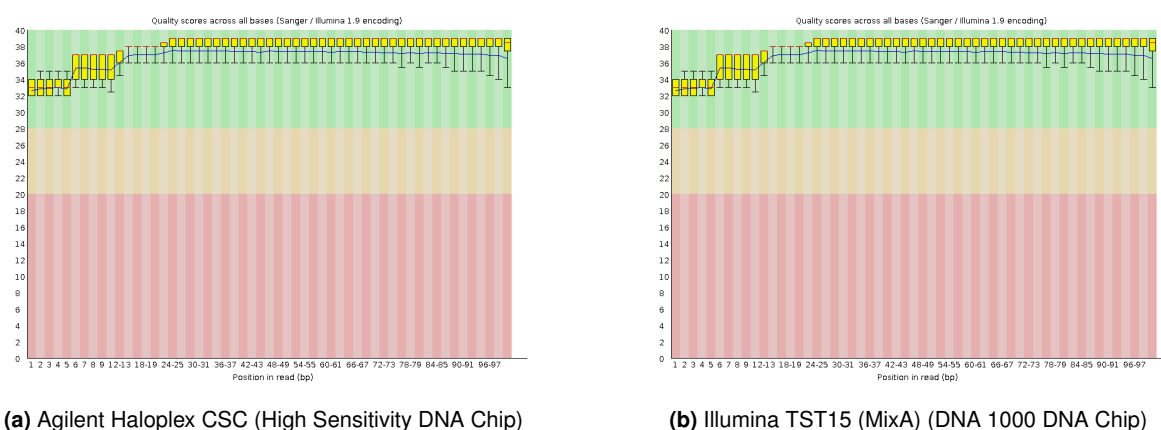


Figure 1: Electropherograms of representative sequencing libraries prepared by Agilent Haloplex ClearSeq Cancer and Illumina TruSight Tumor 15. (*) represents the lower marker, (**) represents the upper marker

Using the blablabla software, the concentration, molarity and total peak area (TPA) of the expected sequencing libraries were calculated.

Maybe: relationship between dCt and TPA?

Maybe: the enrichment of one or the other kit is more affected by bad quality samples

3.2 NGS Data Quality

Sequencing run parameters were calculated by the Illumina Sequencing Viewer software. Table XXX shows the averaged run parameters of runs with Haloplex CSC and TST15 sample preparation.

TST15 has a higher cluster density and therefore a higher total yield, but has lower reads passing a phred-score threshold of Q30 than Haloplex. This is due to the different chemistries used. TST15 uses v3 chemistry, while Haloplex uses v2. v2 generally has lower cluster density and output, but

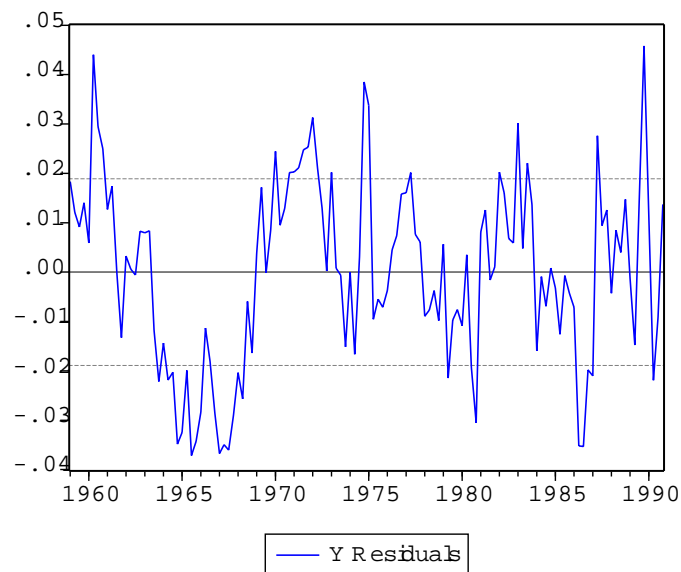


Figure 2: Scatter plot of the corrected peak area (X axis) of the regions corresponding to the sequencing libraries defined in the blabla software and the dCt (Y axis). Agilent Haloplex ClearSeq Cancer data are represented as blue dots, Illumina TruSight Tumor 15 data are represented as red dots.

Table 1: Comparison of Run Parameters (Averaged) of Sequencing Runs with Haloplex CSC & TST15 Sample Preparation

Parameter	Halo CSC	TST15
Yield total (Gb)	3.7	7.37
% >Q30	93.8	82.355
Cluster Density PF (k/mm2)	1084	1180
Cluster Density PF (%)	85.95	<u>79.95</u>

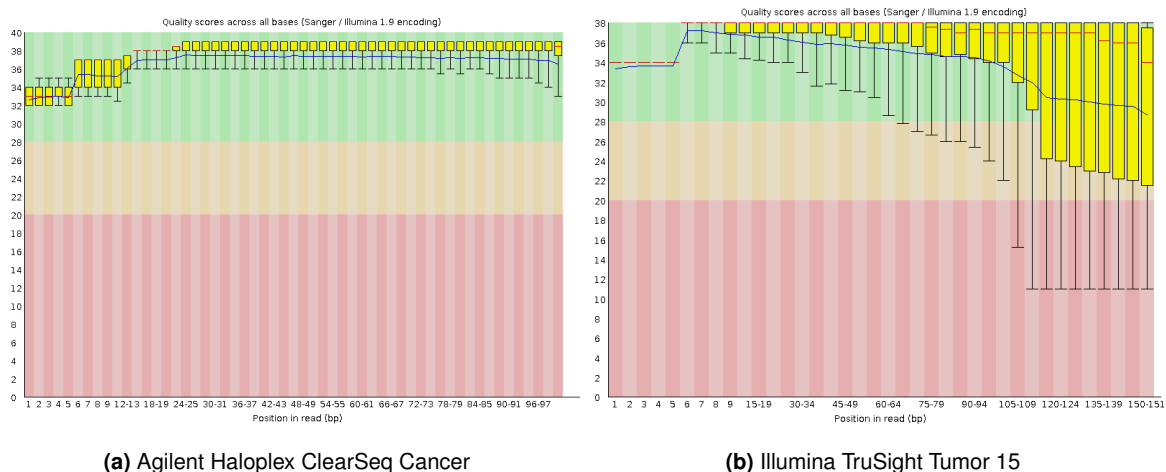


Figure 3: Comparison of coverage distributions per amplicon as reported by FastQC

therefore better quality.

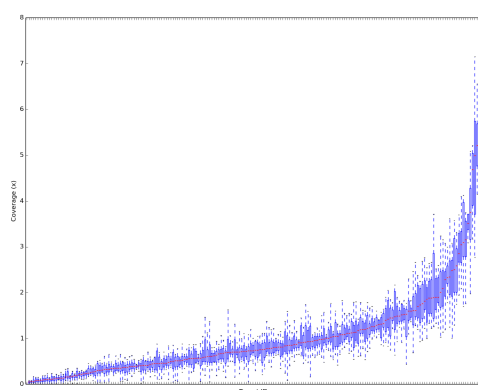
Table XXX shows the boxplot representations of the read qualities per position of two representative FASTQ files as reported by FastQC. Both workflows yield high quality data, yet Haloplex CSC data have more narrow distributions and are of higher quality. This is in direct relationship with the sequencing chemistry kit used.

Table 2: Blablabla

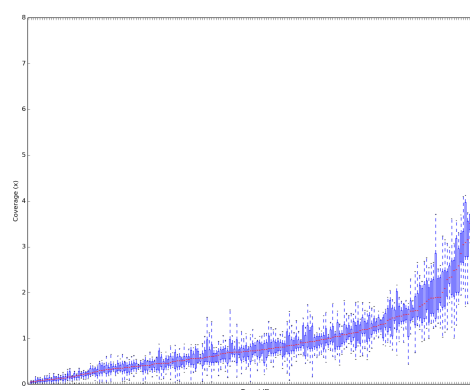
Parameter	Haloplex SureCall A	Haloplex SureCall B	Haloplex Velona	TST15 BaseSpace	TST15 Velona
% mapped	—	—	—	62.6	—
% paired	—	—	—	58.7	—
% singletons	—	—	—	3.8	—

The Samtools Flagstat command was used to determine some basic BAM statistics of BAM files of samples prepared with the respective library preparations and processed with the mentioned bioinformatic pipelines. ?? shows the averaged result of these statistics.

Considering the recommended pipelines, Haloplex CSC data, analyzed with Agilent's SureCall software, has a higher percentage (91%) of mapped reads when compared to Illumina's BaseSpace TruSight Tumor 15 App (62.9%). Data analysis with the recommended SureCall design includes a steps where mates are fixed, but they are not stitched together. Therefore no reads are considered as being paired. The TST15 app in contrast includes a read stitching step and 58% are considered as properly paired. This means that of the 62.6% of mapped reads, 4.2% are not properly paired.



(a) Agilent Haloplex ClearSeq Cancer



(b) Illumina TruSight Tumor 15

Figure 4: Comparison of Coverage Distributions per Amplicon

3.8% of reads processed with the TST15 online App are considered to be singletons, whereas only 1.8% of reads processed with the SureCall software are considered as singletons.

3.3 Coverage Analysis

Coverage Distribution TST15 vs Haloplex

Coverage Distribution per Patient (check if correlation IQR with dCt)

Coverage Distribution per Amplicon (check if some have always lower coverage, check if some failed)

Failed Amplicon Counter

Table 3: Blablabla

Amplicon	1x	50x	250x	500x	1000x
ATM_14	12	0	0	0	0
MAP2K1_2	1	0	0	1	4
ATM_5	0	1	1	0	2
ATM_11	0	1	1	2	3
PTEN_1	0	1	1	0	2
KIT_6	0	1	0	1	2
FGFR3_1	0	0	1	2	2

Table 4: Blablabla

Amplicon	1x	50x	250x	500x	1000x
1METxxxE16TF031SR031	0	11	0	1	3
2KITxxxE09TF003SR003	0	5	6	0	0
2KITxxxE09TF003SR003	0	5	6	0	0
27qxxxxExxTF034SR034	0	0	2	4	5
17qxxxxExxTF018SR018	0	1	1	4	4
1KRASxxE04TF002SR002	0	1	1	3	6
1TP53xxE02TF034SR034	0	0	3	5	9
1EGFRxxE19TF032SR032	0	0	1	5	5
1EGFRxxE21TF035SR035	0	0	1	2	2
2TP53xxE02TF033SR033	0	0	0	1	3

Table XXX and table XXX show how often a given amplicon failed a given coverage threshold. The number of failed amplicons is low for Haloplex CSC as well as for TST15. Most amplicons were amplified efficiently in most samples, which is also confirmed by figure XXX (amplicon distributions). Some amplicons however fail coverage thresholds in several samples.

- Amplicon ATM_14 in Haloplex CSC was never amplified. This amplicon is defined in the BED file but obviously is not part of the kit. (negative control?)
- Amplicons ATM_5, ATM_10, MAP2K1, PTEN_1, KIT_6 and FGFR3_1 in Haloplex CSC data failed a coverage threshold of 1000x in several samples
- In TST15 data, more amplicons fail the respective thresholds. Several amplicons in genes MET, KIT, TP53, KRAS and EGFR fail the 1000x coverage threshold, and often even the required threshold of 500x, which is required by the TruSight Tumor 15 App.

The fact that several amplicons in the genes EGFR and KRAS in TST15 data repetitively fail the required coverage thresholds is problematic. This is especially the case for amplicon 1EGFRxxE21TF035SR035 as it includes the well-known EGFR L858R variant, which confers increases sensitivity for EGFR tyrosine kinase inhibitors.

Fragmentation γ - γ Coverage?

(GATK CallableLoci) (GATK CountLoci???) (GATK FindCoveredIntervals)

On-off target; Enrichment Efficiency TST15 vs Haloplex

Coverage across genome, check where there is coverage

Strandedness?

GATK DepthOfCoverage???

3.4 Variant Calling Algorithm Comparison

3.4.1 Detection of Known Single Nucleotide Variants and Deletions

Table XXX shows known variants in the analyzed samples and the variant frequency reported by the recommended pipelines. There is a high concordance between the results of both kits and previously known variants.

TST15 could not detect EGFR del19 in patient F due to low coverage. The corresponding region was inspected in IGV and the deletion was present. The amplicon was not amplified efficiently and has only a coverage of 167x. The TST15 BaseSpace App however applies a coverage threshold at 500x. Variants with a lower depth are not reported. The same sample was sequenced twice and the deletion was never detected. This is probably due to the fragmentation induced by the FFPE fixation.

Haloplex CSC did not find KRAS p.Gly12Val in patient G, also due to low coverage for this amplicon. The corresponding region was inspected in IGV: the region has a coverage of only 180x. Only one read showed the expected C→A variant. Sample G is also the sample with the worst dCt (2.85). The high fragmentation in this sample is probably responsible for the bad amplification. The same sample will be re-sequenced in a later run to check if the problem is related to the fragmentation of the sample or if library preparation was bad.

Additional previously unknown variants were found (TODO: put a table into the appendix)

3.4.2 Sensitivity Analysis

BRAF Mut and WT samples from Horizon were analyzed. The BRAF Mut sample was sequenced purely, as well as in a 1/3 and 2/3 dilution with the BRAF WT sample.

The observed variant frequencies as detected by the TST15 BaseSpace App are in line with the expected variant frequencies. D816V variant in cKIT could not be observed as the position of this

Table 5: Blablabla

Sample	Context	Tissue	Known variant	Halo CSC Freq (%)	TST15 Freq (%)
A	NSCLC	FFPE	EGFR L858R	24.7	21.4
B	NSCLC	FFPE	EGFR L858R	24.3	13.2
C	NSCLC	FFPE	EGFR L858R	32.7	29.8
D	Melanoma	FFPE	BRAF V600E	44.4	47.6
E	Melanoma	FFPE	BRAF V600E	18.4	21.4
F	NSCLC	FFPE	EGFR del19	not found	50
G	mCRC	FFPE	KRAS CD 12_13	p.Gly12Val (3.4)	not found
H	mCRC	FFPE	KRAS CD 12_13	p.Gly12Asp (37.4)	G12D (31.4)
I	mCRC	FFPE	KRAS CD 12_13	p.Gly13Asp (9.7)	G13D (8.7)
J	mCRC	FFPE	NRAS p.Gly12Asp	25.9	27.2
K	Melanoma	FFPE	BRAF V600E	66.2	59.5
L	mCRC	FFPE	NRAS p.Gly13Val	6.6	5
M	mCRC	FFPE	KRAS and NRAS WT	WT	WT
N	mCRC	FFPE	KRAS and NRAS WT	WT	WT
O	Melanoma	FFPE	BRAF V600E	37.4	

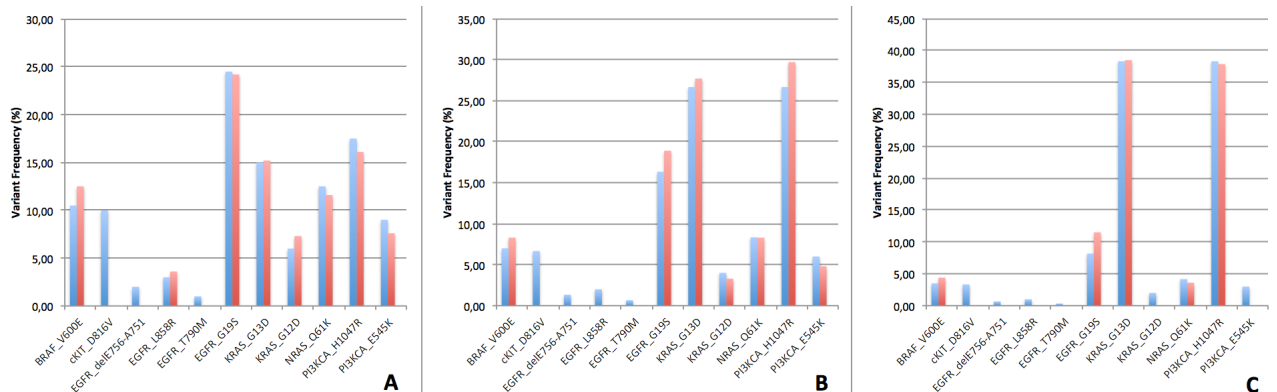


Figure 5: Blablabla

Table 6: Blablabla

Gene	Variant	Exp 100%	Obs Halo CSC	Obs TST 15	Exp at 66%	Obs Halo CSC	Obs TST15	Exp at 33%	Obs Halo CSC	Obs TST15
BRAF	V600E	10.5	–	12.5	7	–	8.3	3.5	–	4.4
cKIT	D816V	10	–	–	6.67	–	–	3.33	–	–
EGFR	delE756- A751	2	–	–	1.33	–	–	0.67	–	–
EGFR	L858R	3	–	3.6	2	–	–	1	–	–
EGFR	T790M	1	–	–	0.67	–	–	0.33	–	–
EGFR	G719S	24.5	–	24.2	16.33	–	18.9	8.17	–	11.5
KRAS	G13D	15	–	15.2	26.67	–	27.7	38.33	–	38.5
KRAS	G12D	6	–	7.3	4	–	3.3	2	–	–
NRAS	Q61K	12.5	–	11.6	8.33	–	8.3	4.17	–	3.6
PIK3CA	H1047R	17.5	–	16.1	26.67	–	29.7	38.33	–	37.9
PIK3CA	E545K	9	–	7.6	6	–	4.8	3	–	–

variant is not covered by the TST15 kit.

<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3986649/>

TODO: do the same with Haloplex CSC

TODO: call variants with MuTect 1.1.7, VarScan 2, GATK HaplotypeCaller, SomVarIUS?, Free-bayes?, Vardict??????? TODO: compare results

Among the variants detected by MuTect, 40–50 % are C₂T variants. This has been reported in several studies. TODO: do this for all samples, check if this is really statistically significant or only happened in a few samples

Tools that may be of use somehow: GATK SelectVariants; GATK VariantFiltration; GATK VariantEval; GATK ValidateVariants

4 Conclusions

References

- [1] M. Berg, “EGFR and downstream genetic alterations in KRAS/BRAF and PI3K/AKT pathways in colorectal cancer – implications for targeted therapy,” *Discovery Medicine*, vol. 14, no. 76, pp. 207–214, 2012.
- [2] K. Cibulskis, M. S. Lawrence, S. L. Carter, A. Sivachenko, D. Jaffe, C. Sougnez, S. Gabriel, M. Meyerson, E. S. Lander, and G. Getz, “Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples,” *Computational Biology*, 2013.
- [3] J. Ferlay, I. Soerjomataram, R. Dikshit, S. Eser, C. Mathers, M. Rebelo, D. M. Parkin, D. Forman, and F. Bray, “Cancer incidence and mortality worldwide: Sources, methods and major patterns in globocan 2012,” *International Journal of Cancer*, 2015.
- [4] R. d. d. Direction de la Santé, Service des statistiques, “Statistiques des causes de décès pour l’année 2014,” 2014.
- [5] F. Mertes, A. E. Sharawy, S. Sauer, J. M. L. M. van Helvoort, P. van der Zaag, A. Franke, M. Nilsson, H. Lehrach, and A. J. Brookes, “Targeted enrichment of genomic dna regions for next-generation sequencing,” *Briefings in Functional Genomics*, vol. 10, no. 6, pp. 374–386, 2011.