

*Center for Biomedical Informatics and Information Technology (CBIIT)*  
NCI/Shady Grove Campus





- We are not endorsing, promoting any specific tool during this training.
- We may recommend open source tools based on the literature and our experience.

# *Outline of Exome Sequencing Training*

## *Overview*

- Quick introduction of CGBG Team..
- Introduction to exome sequencing analysis workflow.
- Challenges in Whole Exome Sequencing.
- Exome sequencing pipelines.

## *Hands On Tutorial*

- Get a Galaxy account
- Upload dataset and pre-alignment QC
- Alignment with BWA and post-alignment Summary
- Build pileup for variant call and variant détection and variant annotation
- Visualize alignments and variants on IGV

# DNA sequencing technologies: 2006–2016

Elaine R Mardis

*Nature Protocols* 12, 213–218 (2017) | doi:10.1038/nprot.2016.182

Received 17 August 2016 | Accepted 20 October 2016 | Published online 05 January 2017



## Abstract

[Abstract](#) • [Introduction](#) • [References](#) • [Acknowledgments](#) • [Author information](#)

Recent advances in the field of genomics have largely been due to the ability to sequence DNA at increasing throughput and decreasing cost. DNA sequencing was first introduced in 1977, and next-generation sequencing technologies have been available only during the past decade, but the diverse experiments and corresponding analyses facilitated by these techniques have transformed biological and biomedical research. Here, I review developments in DNA sequencing technologies over the past 10 years and look to the future for further applications.



All Content

[Search](#) [Advanced Search](#)

[< Previous Article](#)

**February 2016** Volume 43, Issue 1, Pages 36–48

[Next Article >](#)

## The promise of omics-based approaches to cancer prevention

Daoud Meerzaman Barbara K. Dunn, Maxwell Lee, Qingrong Chen, Chunhua Yan, Sharon Ross

DOI: <http://dx.doi.org/10.1053/j.seminoncol.2015.09.004>

[Article Info](#)



**Abstract**

[Full Text](#)

[Images](#)

[References](#)

## Abstract

Cancer is a complex category of diseases caused in large part by genetic or genomic, transcriptomic, and epigenetic or epigenomic alterations in affected cells and the surrounding microenvironment. Carcinogenesis reflects the clonal expansion of cells that progressively acquire these genetic and epigenetic alterations—changes that, in turn, lead to modifications at the RNA level. Gradually advancing technology and most recently, the advent of next-generation sequencing (NGS), combined with bioinformatics analytic tools, have revolutionized our ability to interrogate cancer cells. The ultimate goal is to apply these high-throughput technologies to the various aspects of clinical cancer care: cancer-risk assessment, diagnosis, as well as target identification for treatment and prevention. In this article, we emphasize how the knowledge gained through large-scale omics-oriented approaches, with a focus on variations at the level of nucleic acids, can inform the

# Computational Genomics & Bioinformatics Branch (GCB)



Computational Genomics &  
Bioinformatics Branch

Computational Genomics &  
AI

Scientific & Tools Consulting  
and training

Provide state-of-art data analysis on a variety of multifaceted cancer-related research projects including proteogenomic, AI, and clinical high-throughput data.

- NGS and Integrated analysis of proteogenomic data for Intramural NCI and moonshot projects (APOLLO, CIMAC-CIDC and SeQC2).
- Multimodal AI and ML approaches for patient outcome prediction. Deep learning algorithms, CNNs.
- Generate and implement a web-based workflow in CRDC (CGC) and OmicCircos and (MOGSA)
- Widely used tools for NGS analyses and Visualization tools
- ITCR Scientific software
- CBIIT Summer Internship Program



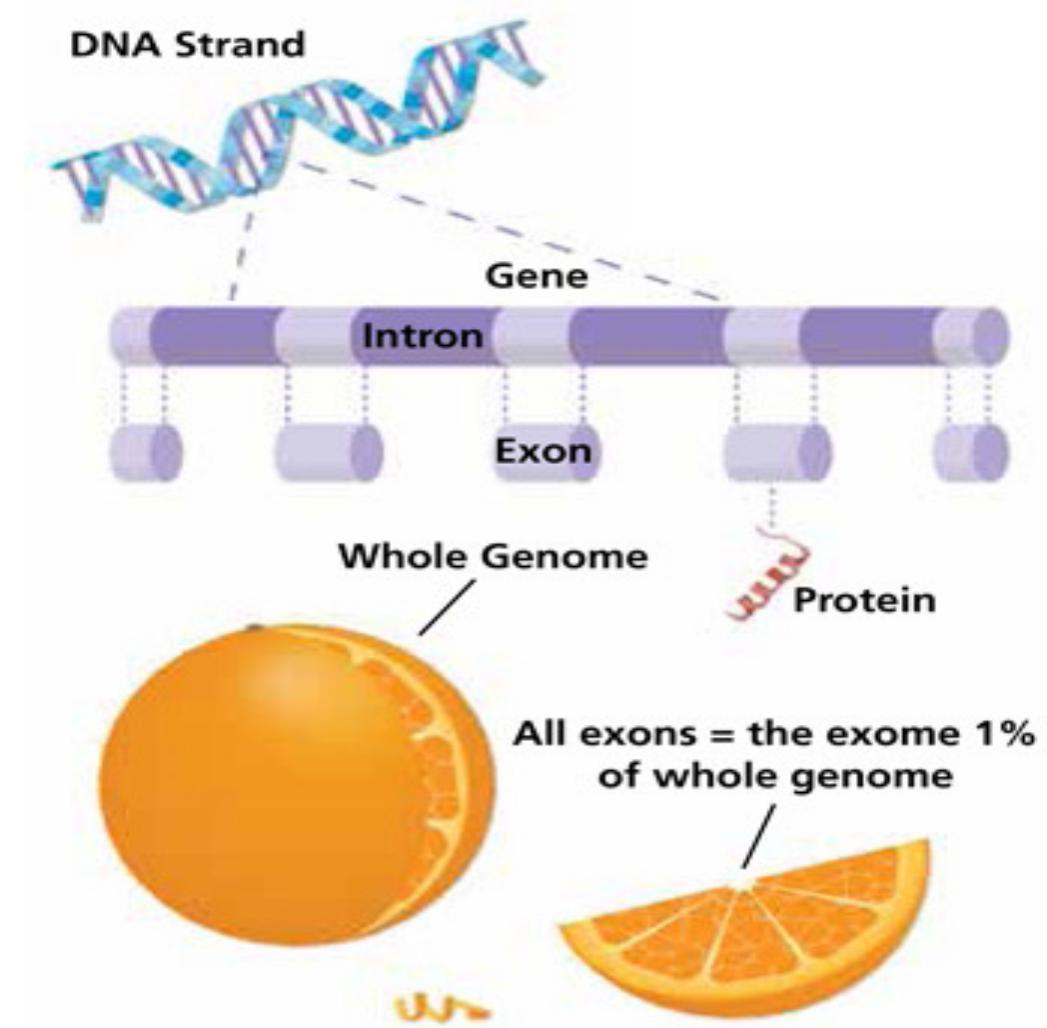
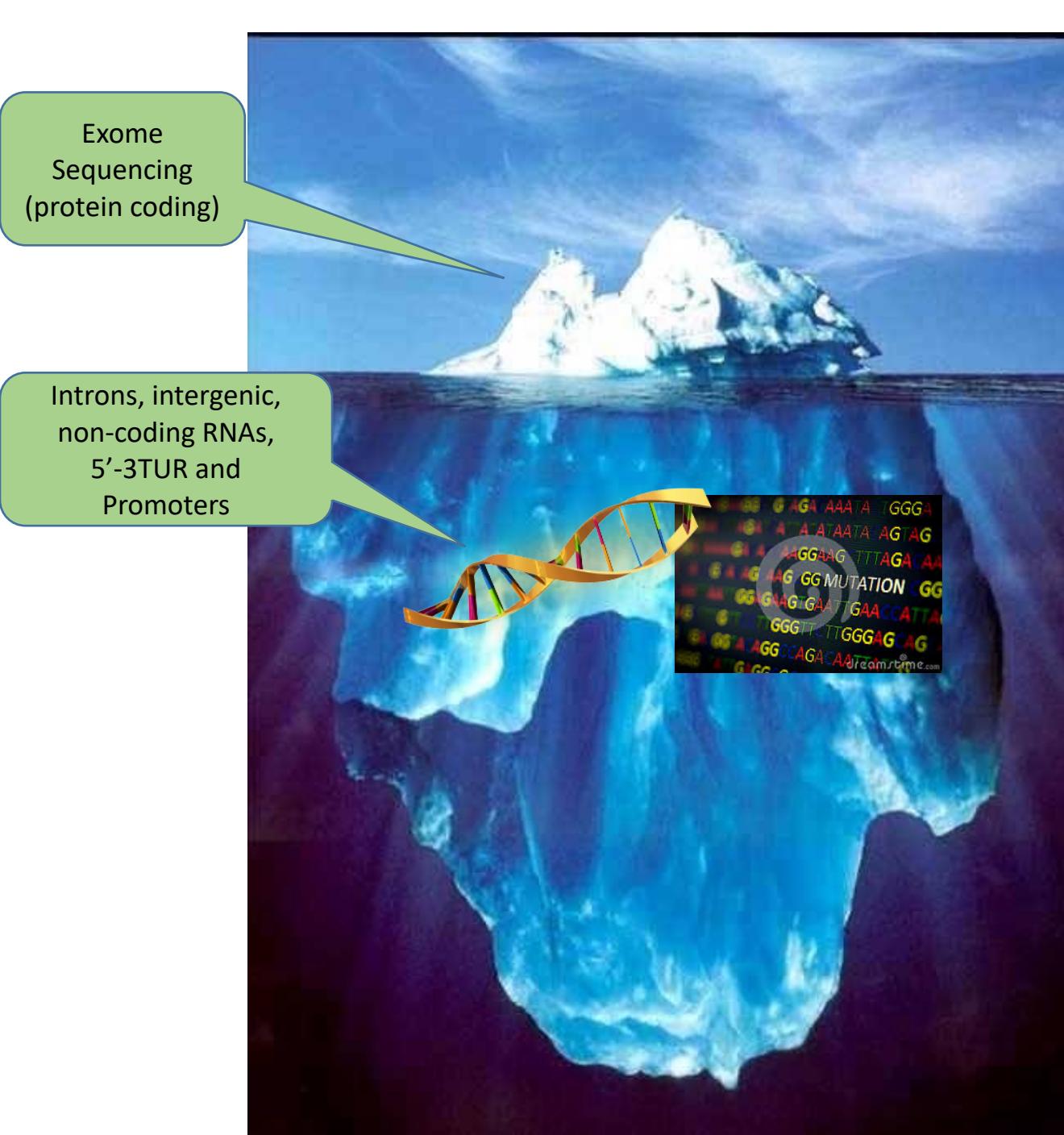
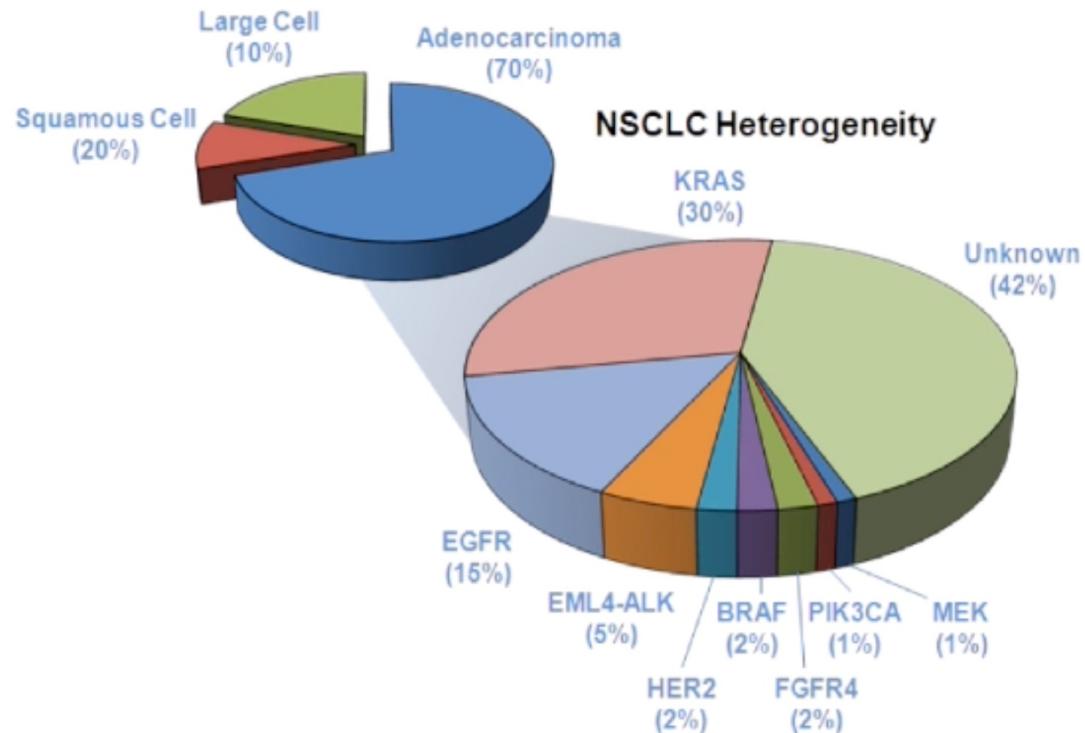


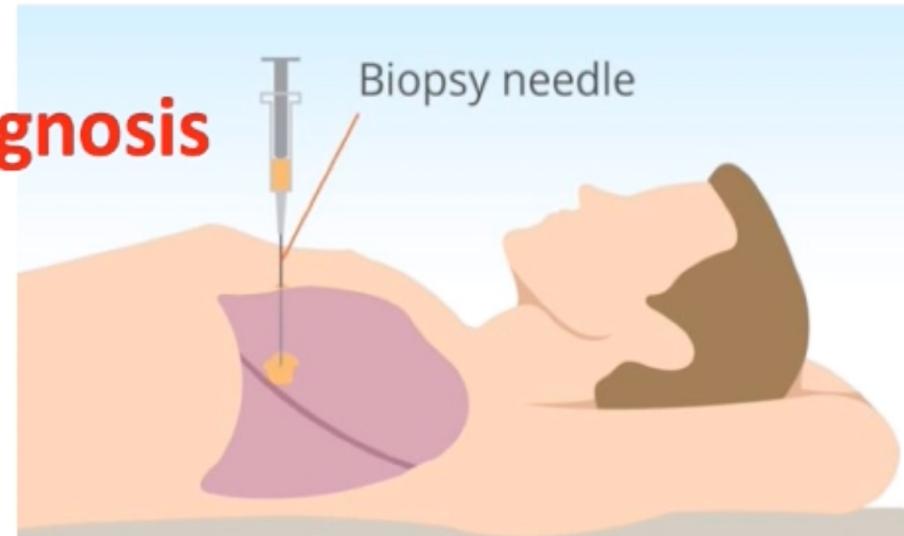
ILLUSTRATION BY MEAHGAN HARRIGAN

# Adult cancers: genomics provides a treatment playbook

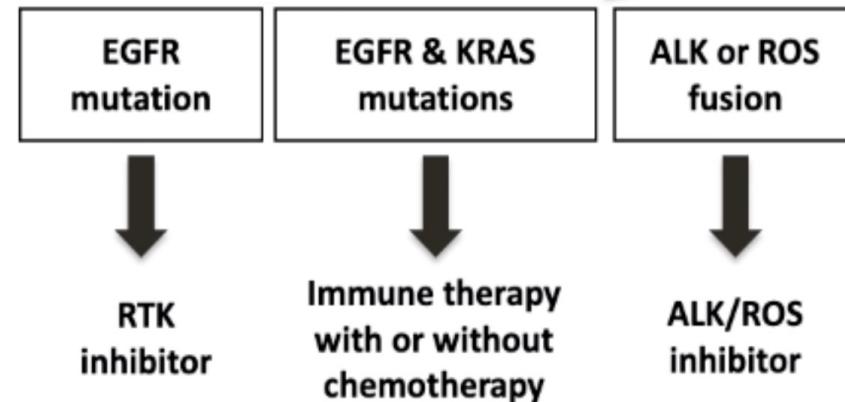
## discovery



## diagnosis



## “Genome analysis”



## Why Whole Exome Sequencing?

Exons are part of the genome we understand best

~ 85% of known mutations in Mendelian diseases affect the exome: Nonsense, missense, splice, indel mutations

Small fraction of the whole genome

~1% of human genomes  
~200,000 coding exons in  
~20,000 genes  
~35-60Mb

Less Resources

- Whole exome Seq is 1/6 the cost of whole genome
- 1/15 the amount of data
- More resources
- More data interpretation

# Comparing different Sequencing technologies

Pros

Sanger	Targeted	Exome	Whole Genome
<ul style="list-style-type: none"><li>• Accurate</li><li>• Cheap per exon</li><li>• High turn-around</li></ul>	<ul style="list-style-type: none"><li>• Optimization possible</li><li>• Low chance of incidental findings</li><li>• Easy analysis</li></ul>	<ul style="list-style-type: none"><li>• No bias for genes</li><li>• Standardized workflow</li><li>• Re-use of performed exomes to interpret new ones</li></ul>	<ul style="list-style-type: none"><li>• No sequencing bias</li><li>• Detect SVs and SNVs</li></ul>
<ul style="list-style-type: none"><li>• Low diagnostic yield for genetically heterogeneous diseases</li></ul>	<ul style="list-style-type: none"><li>• Design and re-design required</li><li>• Different designs for different disorders</li></ul>	<ul style="list-style-type: none"><li>• No non-coding regions</li><li>• Sequencing bias</li><li>• Incidental findings</li></ul>	<ul style="list-style-type: none"><li>• Data analysis bottleneck</li><li>• Interpretation of non-coding regions</li><li>• Expensive, time-consuming</li></ul>

Cons

# Factors Impacting Cancer Genomic Data

Bio-samples  
Fresh frozen tissue,  
Cell lines  
FFPE tissue

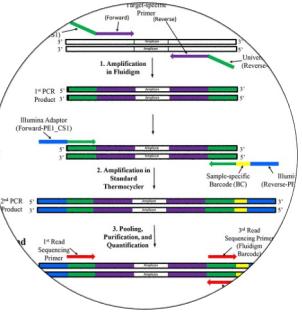
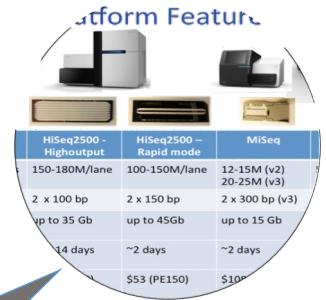
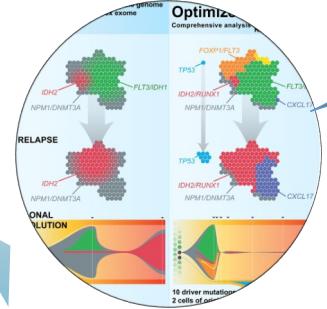
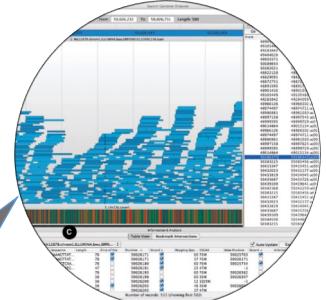
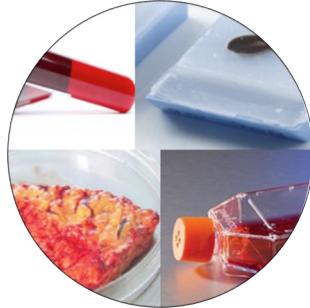
Tumor purity,  
immune cells,  
fibroblasts and blood  
vessels.

Algorithm  
Reads alignment  
Mapping quality score  
Quality filter  
Mutation allele frequency

DNA Extraction  
Fragmentation  
Library prep

System errors  
Different platform

Different  
sequencing sites



# *Whole exome DNA sources*

## Tumor DNA

- Fresh frozen (FF)
- Formalin-Fixed Paraffin-Embedded (FFPE)
- Cell line
- Xenograft

## Normal tissue

- Blood
- Tumor adjacent tissue

OPEN  ACCESS Freely available online



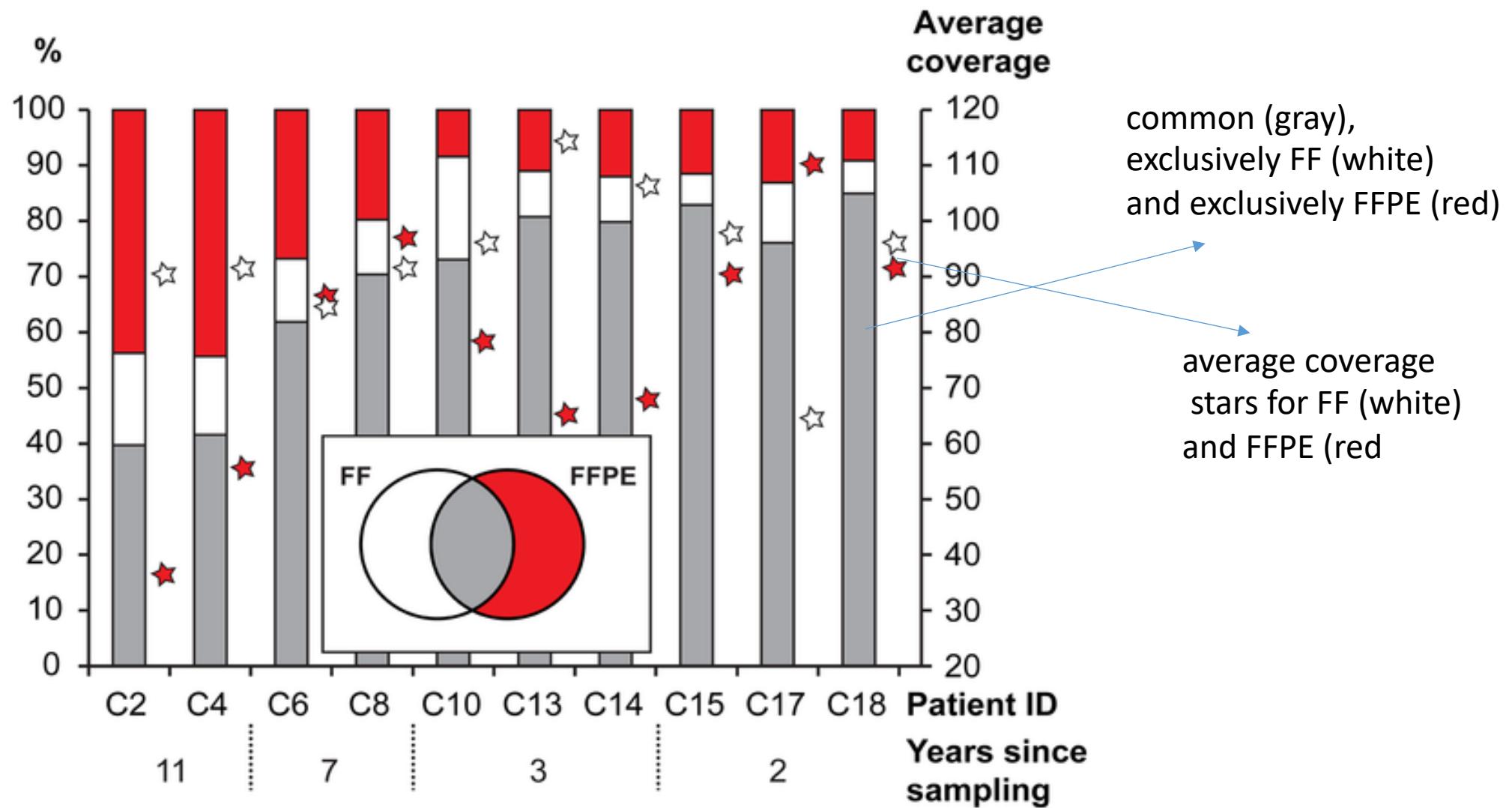
## **Next-Generation Sequencing of RNA and DNA Isolated from Paired Fresh-Frozen and Formalin-Fixed Paraffin-Embedded Samples of Human Cancer and Normal Tissue**

**Jakob Hedegaard<sup>1\*</sup>, Kasper Thorsen<sup>1</sup>, Mette Katrine Lund<sup>2</sup>, Anne-Mette K. Hein<sup>3</sup>, Stephen Jacques Hamilton-Dutoit<sup>4</sup>, Søren Vang<sup>1</sup>, Iver Nordentoft<sup>1</sup>, Karin Birkenkamp-Demtröder<sup>1</sup>, Mogens Kruhøffer<sup>2</sup>, Henrik Hager<sup>4</sup>, Bjarne Knudsen<sup>3</sup>, Claus Lindbjerg Andersen<sup>1</sup>, Karina Dalsgaard Sørensen<sup>1</sup>, Jakob Skou Pedersen<sup>1</sup>, Torben Falck Ørntoft<sup>1</sup>, Lars Dyrskjøt<sup>1</sup>**



**1** Department of Molecular Medicine (MOMA), Molecular Diagnostic Laboratory, Aarhus University Hospital, Skejby, Aarhus, Denmark, **2** AROS Applied Biotechnology A/S, Science Park Skejby, Aarhus, Denmark, **3** CLC bio, Aarhus, Denmark, **4** Institute of Pathology, Aarhus University Hospital, Aarhus, Denmark

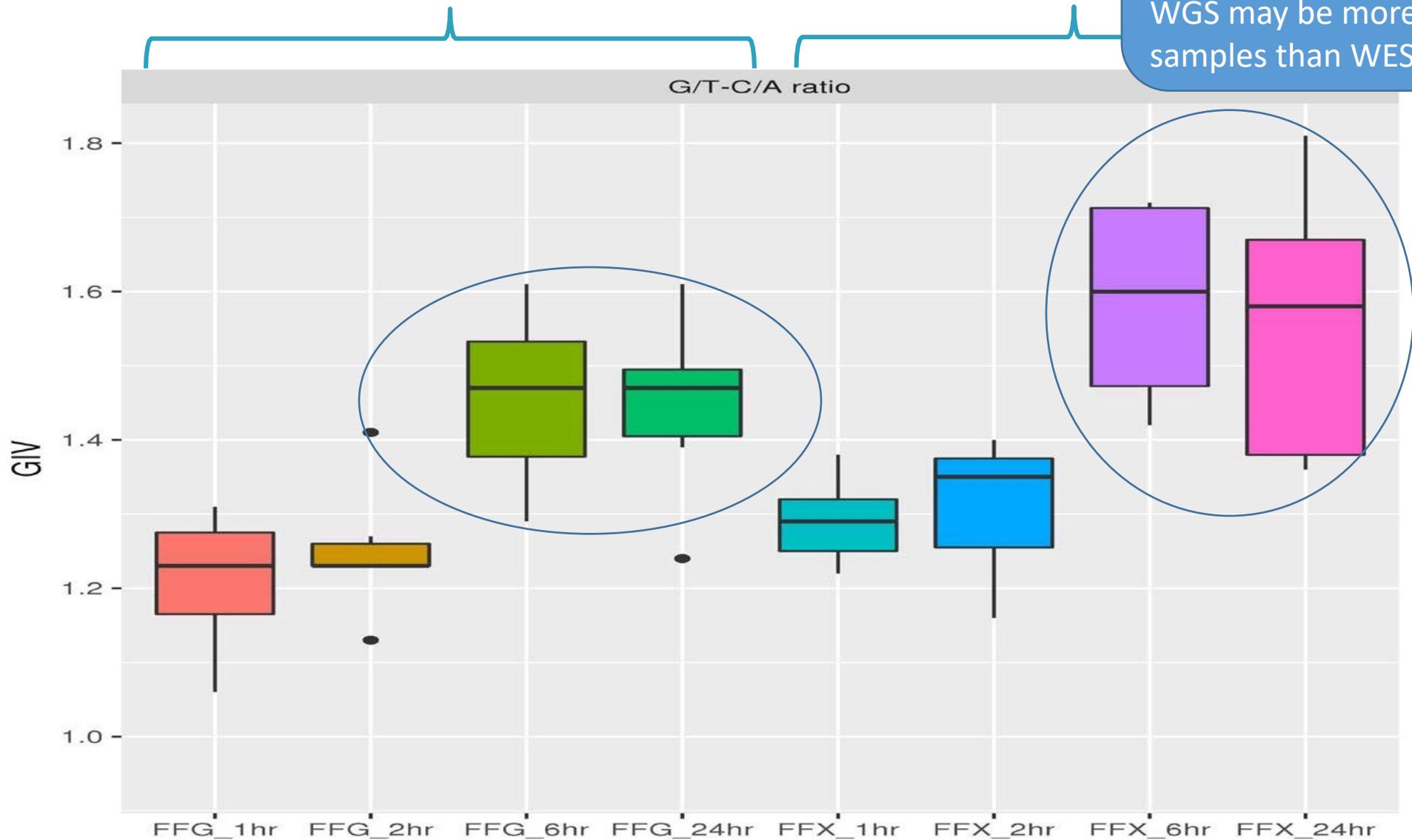
# *Variants detected in exome sequencing data from the paired FF/FFPE samples*



# FFPE artifacts

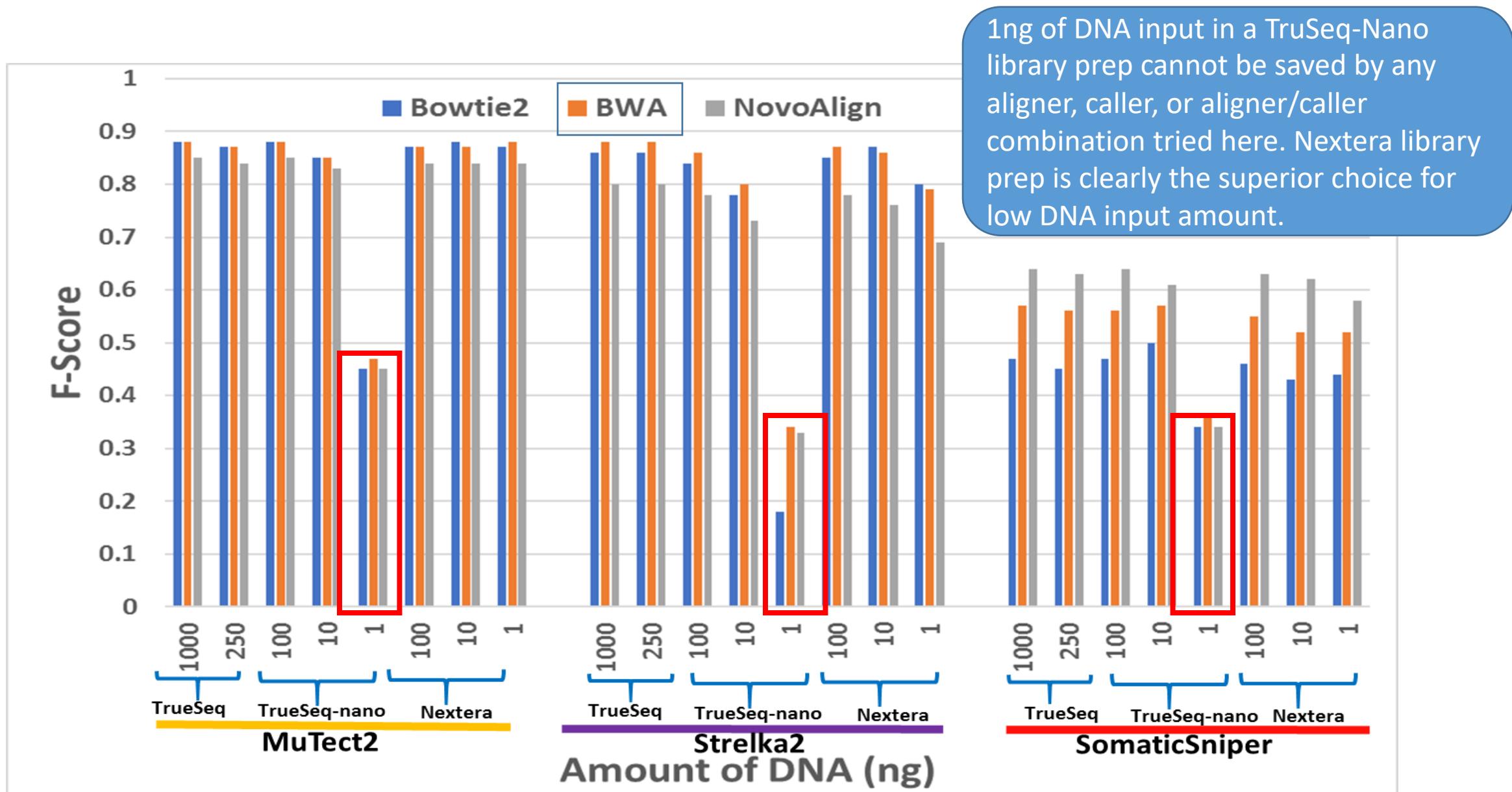
WGS

WES

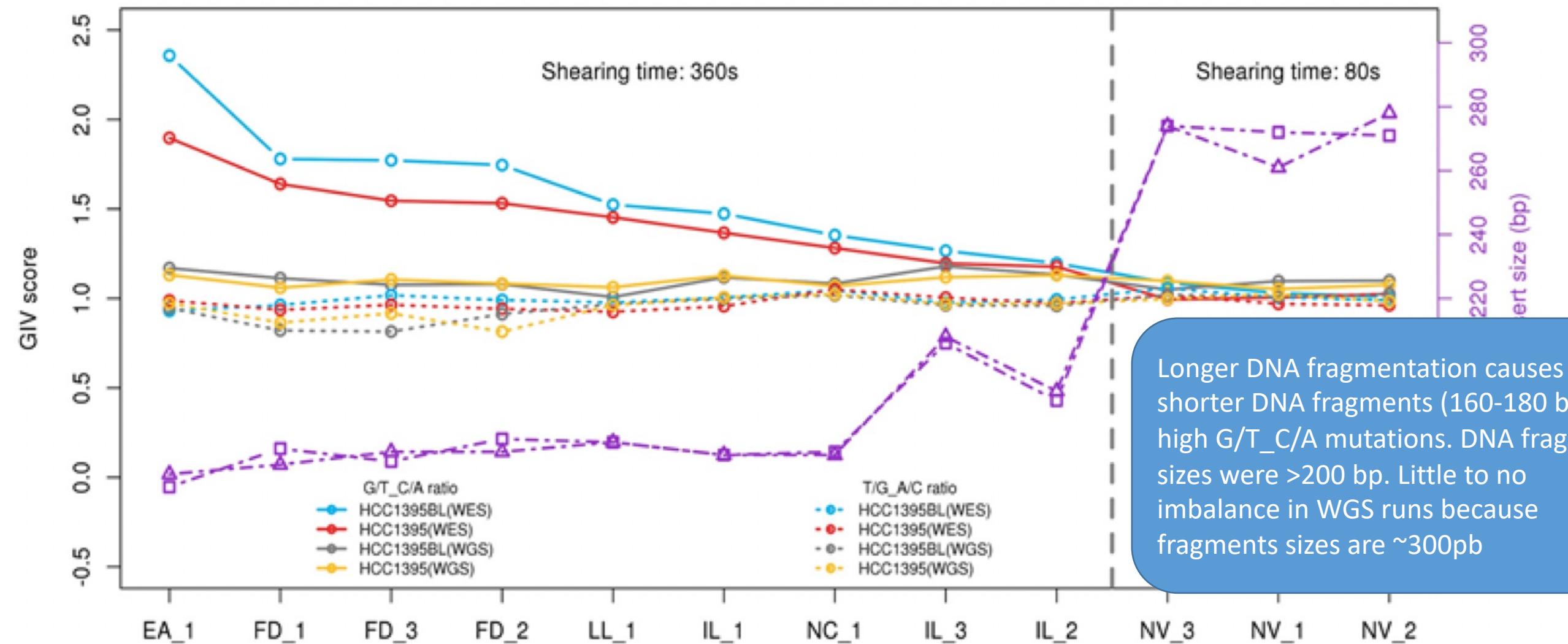


Formaldehyde also causes the deamination of guanine. formaldehyde fixing have an additive DNA damage effect. WGS may be more suitable for FFPE samples than WES.

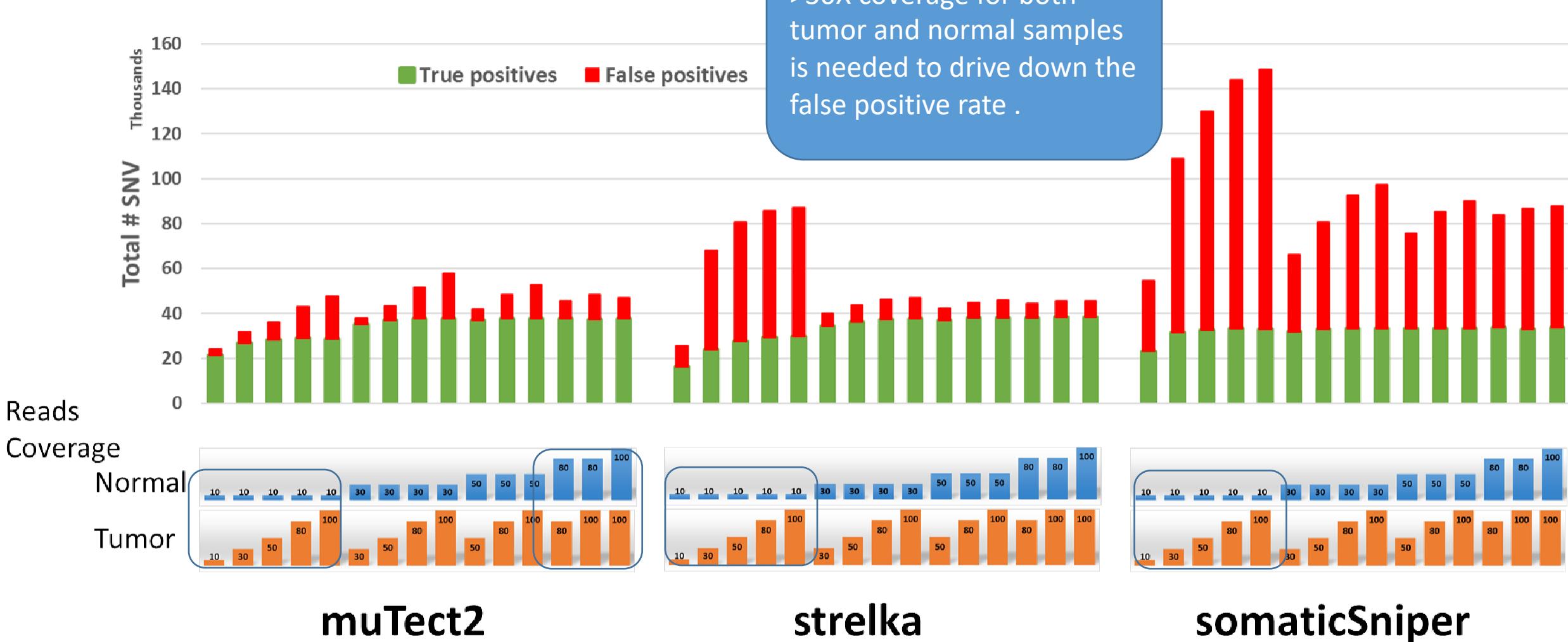
## *Effect of DNA inputs (How much DNA to start with)*



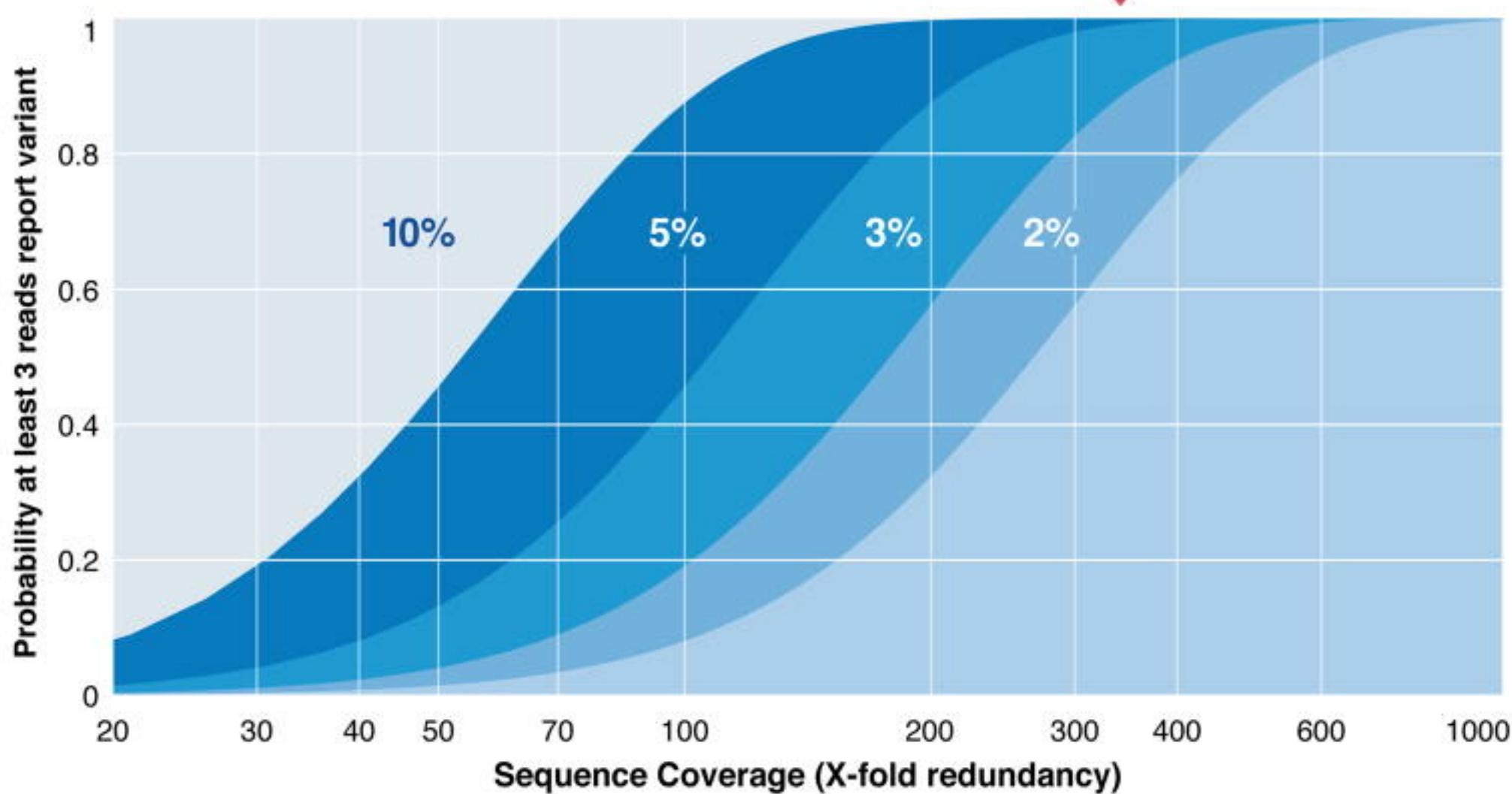
# *G/T\_C/A mutation caused by oxidative DNA damage*



# *reads coverages (Non-analytical factors)*



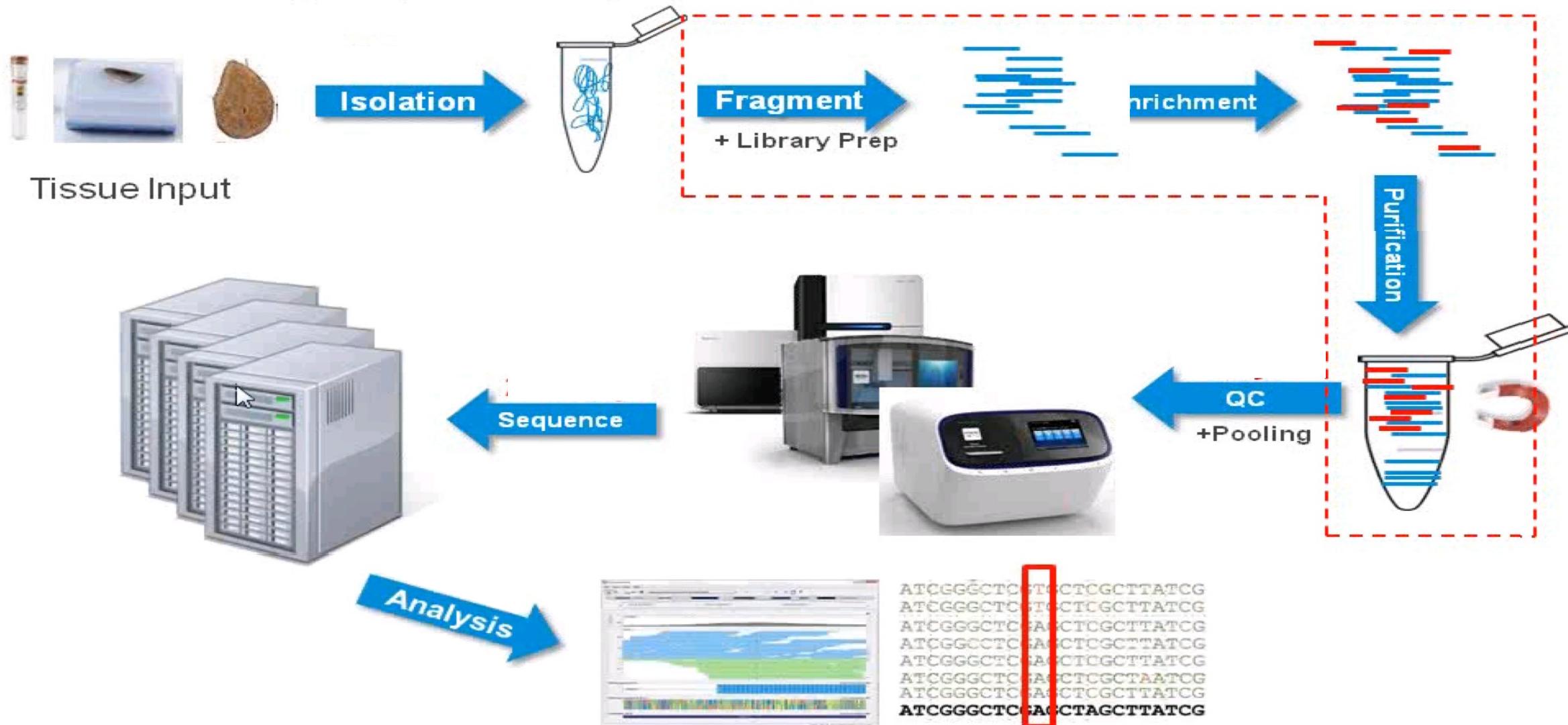
*High coverage is needed for low tumor*



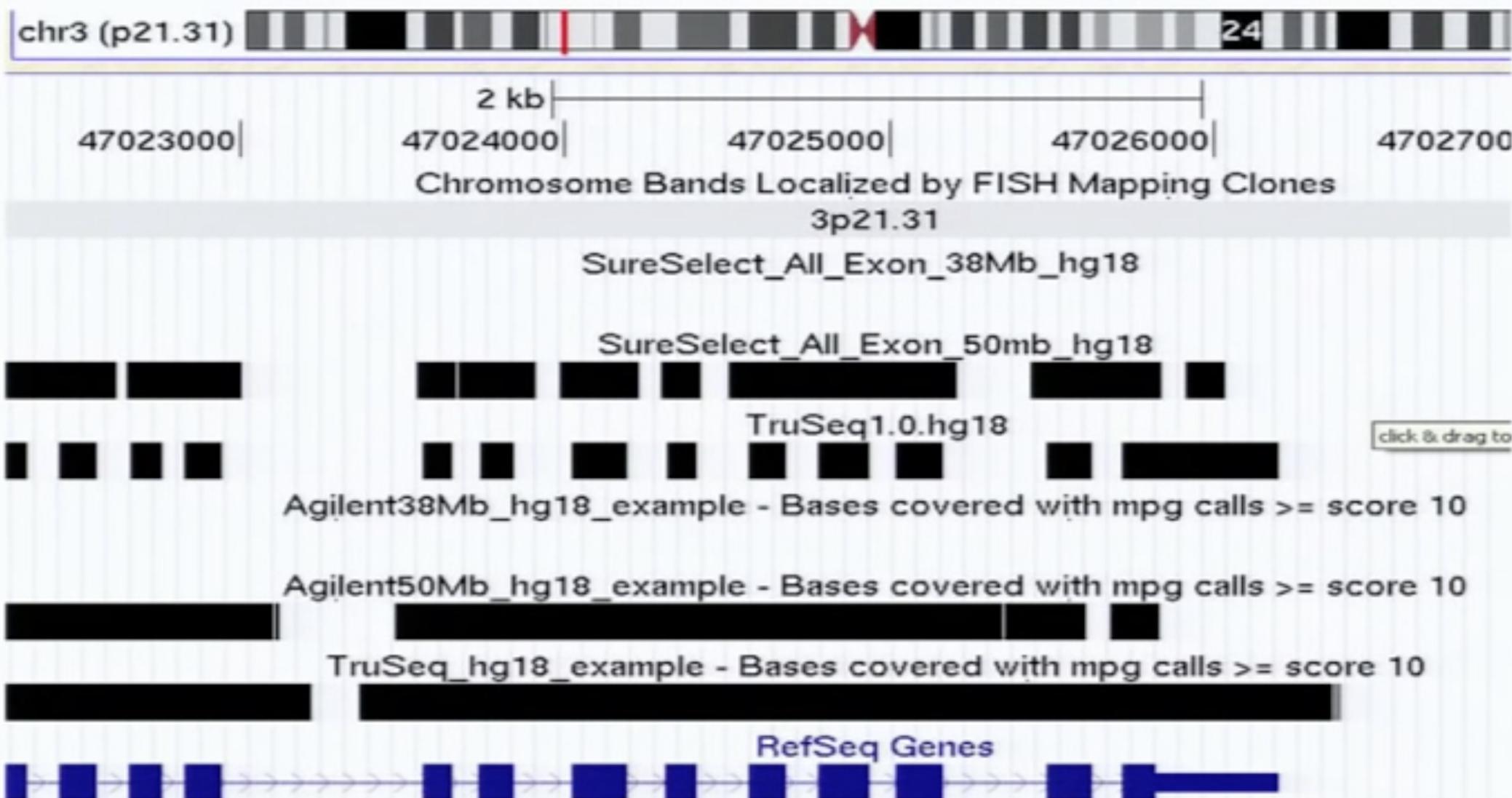
# Exome Sequencing Workflow



Times can always vary at each step depending on reference lab



## NBEAL2 5'



# Performance comparison of four exome capture systems for deep sequencing

Chandra Sekhar Reddy Chilamakuri<sup>1,3\*</sup>, Susanne Lorenz<sup>1,3,4</sup>, Mohammed-Amin Madou<sup>1,4</sup>, Daniel Vodák<sup>1,5</sup>, Jinchang Sun<sup>1,3,4</sup>, Eivind Hovig<sup>1,2,3,5</sup>, Ola Myklebost<sup>1,3</sup> and Leonardo A Meza-Zepeda<sup>1,3,4\*</sup>

## Abstract

**Background:** Recent developments in deep (next-generation) sequencing technologies are significantly impacting medical research. The global analysis of protein coding regions in genomes of interest by whole exome sequencing is a widely used application. Many technologies for exome capture are commercially available; here we compare the performance of four of them: NimbleGen's SeqCap EZ v3.0, Agilent's SureSelect v4.0, Illumina's TruSeq Exome, and Illumina's Nextera Exome, all applied to the same human tumor DNA sample.

**Results:** Each capture technology was evaluated for its coverage of different exome databases, target coverage efficiency, GC bias, sensitivity in single nucleotide variant detection, sensitivity in small indel detection, and technical reproducibility. In general, all technologies performed well; however, our data demonstrated small, but consistent differences between the four capture technologies. Illumina technologies cover more bases in coding and untranslated regions. Furthermore, whereas most of the technologies provide reduced coverage in regions with low or high GC content, the Nextera technology tends to bias towards target regions with high GC content.

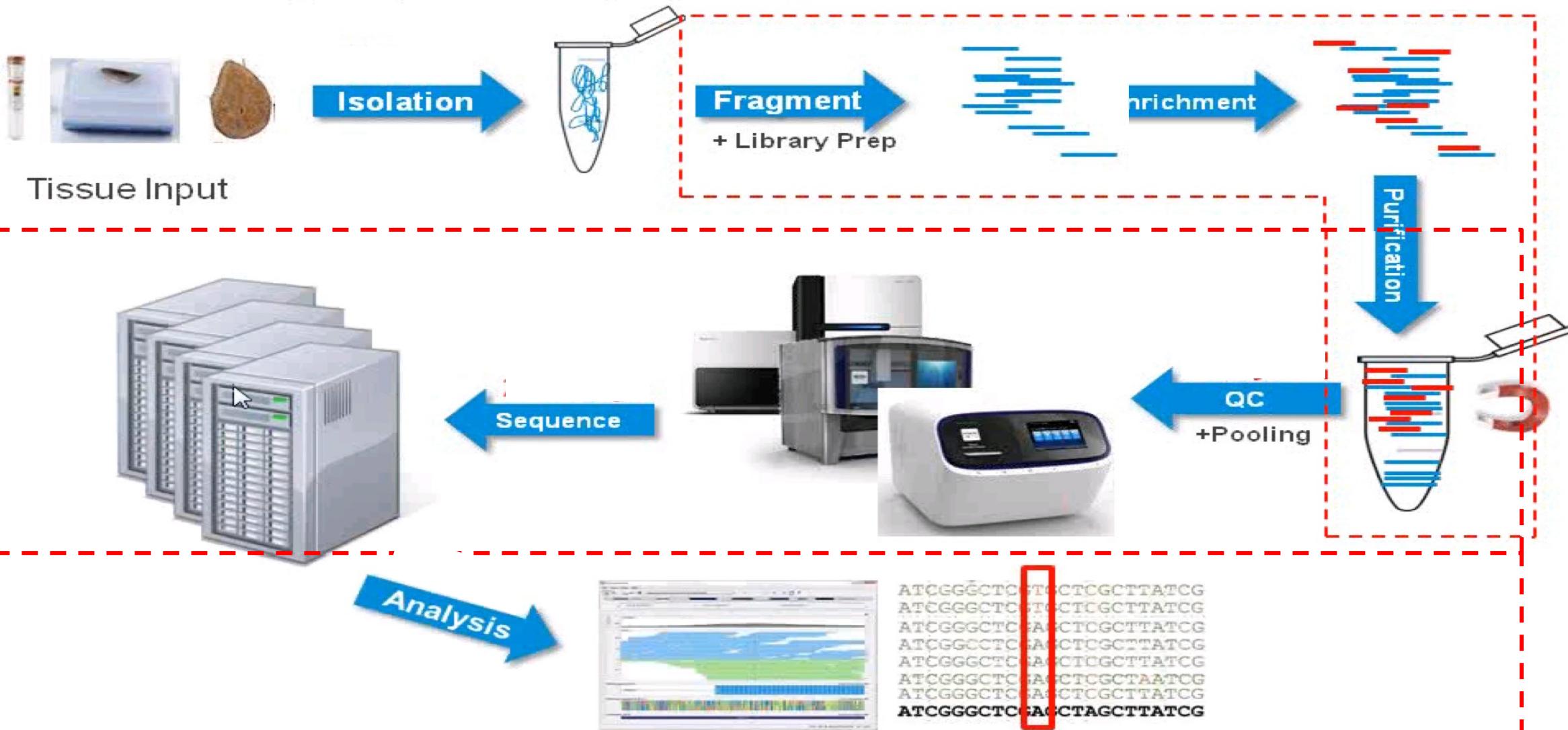
**Conclusions:** We show key differences in performance between the four technologies. Our data should help researchers who are planning exome sequencing to select appropriate exome capture technology for their particular application.

**Keywords:** Exome capture technology, Next-generation sequencing, Coverage efficiency, Enrichment efficiency, GC bias, Single nucleotide variant, Indel

# Exome Sequencing Workflow



Times can always vary at each step depending on reference lab



# Sequencing systems for every lab



## KEY APPLICATIONS

Production-Scale Whole-Genome Sequencing

ON THE HISEQ X FIVE SYSTEM

Population-Scale Whole-Genome Sequencing

ON THE HISEQ X TEN SYSTEM

Population power for whole-genome sequencing.

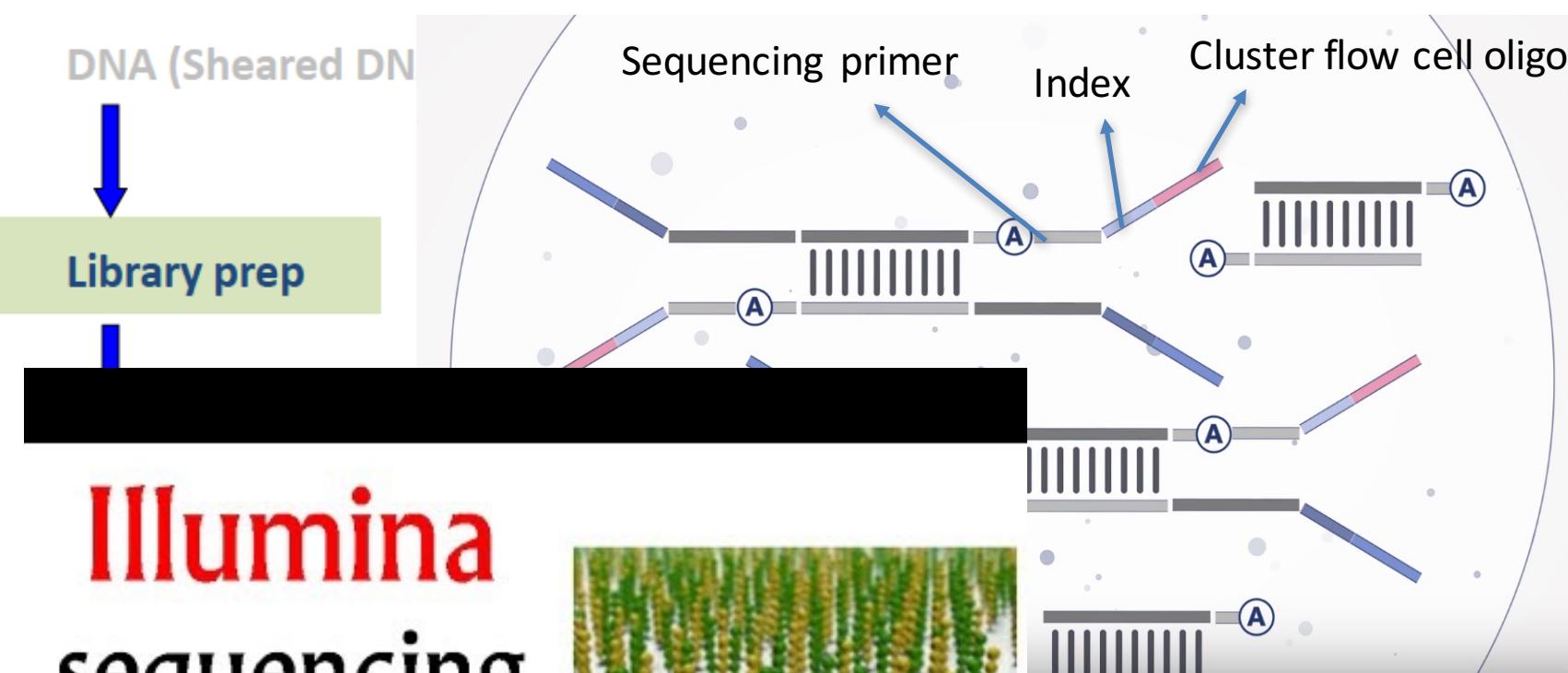
[Explore HiSeq X](#)

[Explore All Sequencing Platforms](#)

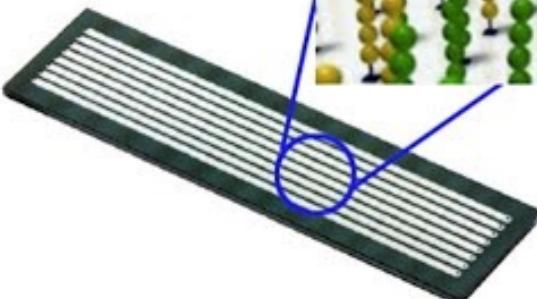
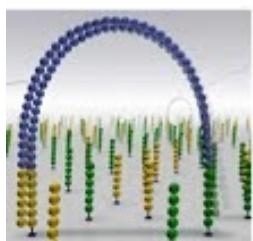
<https://www.illumina.com/systems.html>

# CLINICAL EXOME SEQUENCING

## Work flow :



Illumina  
sequencing

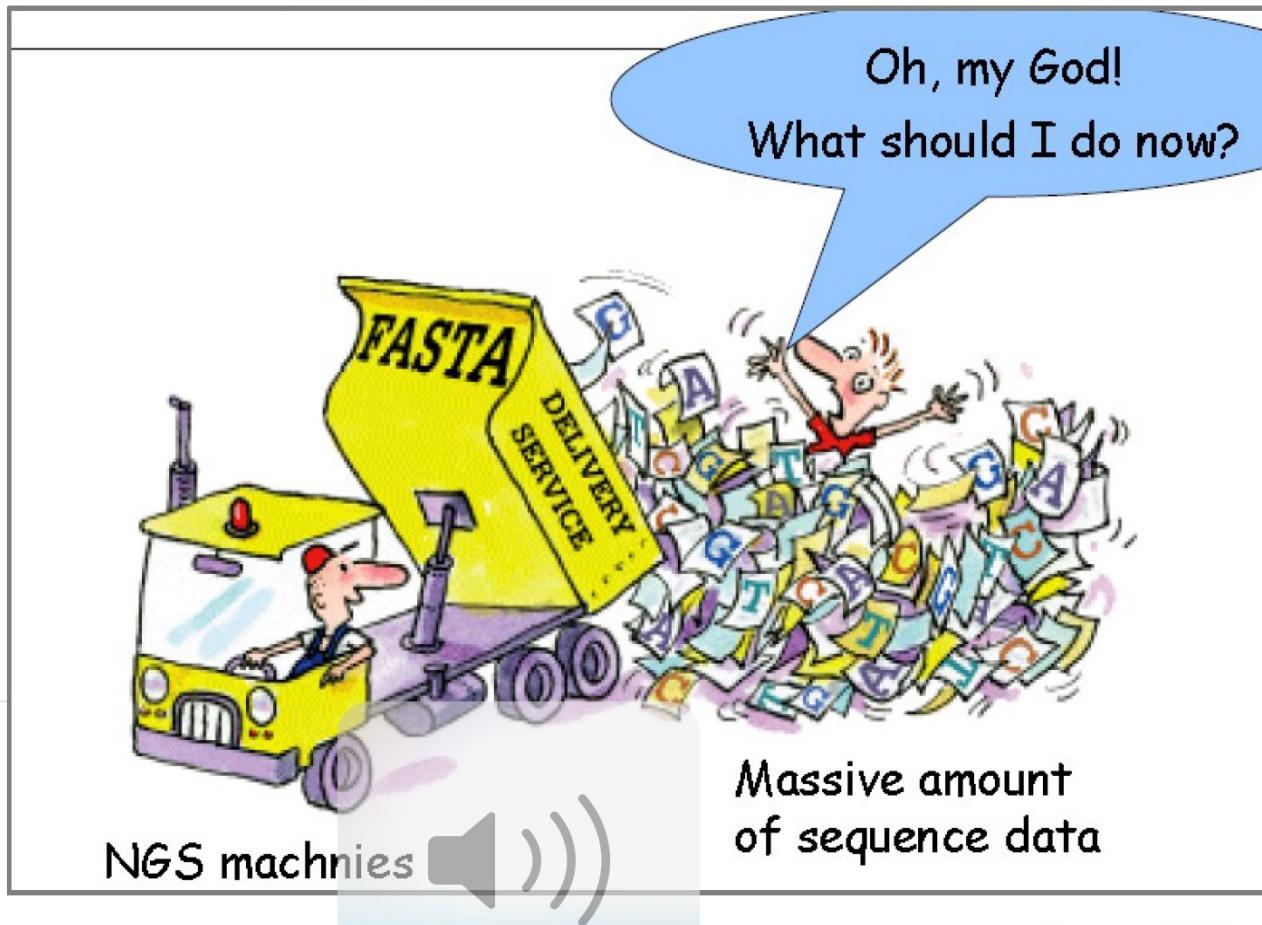


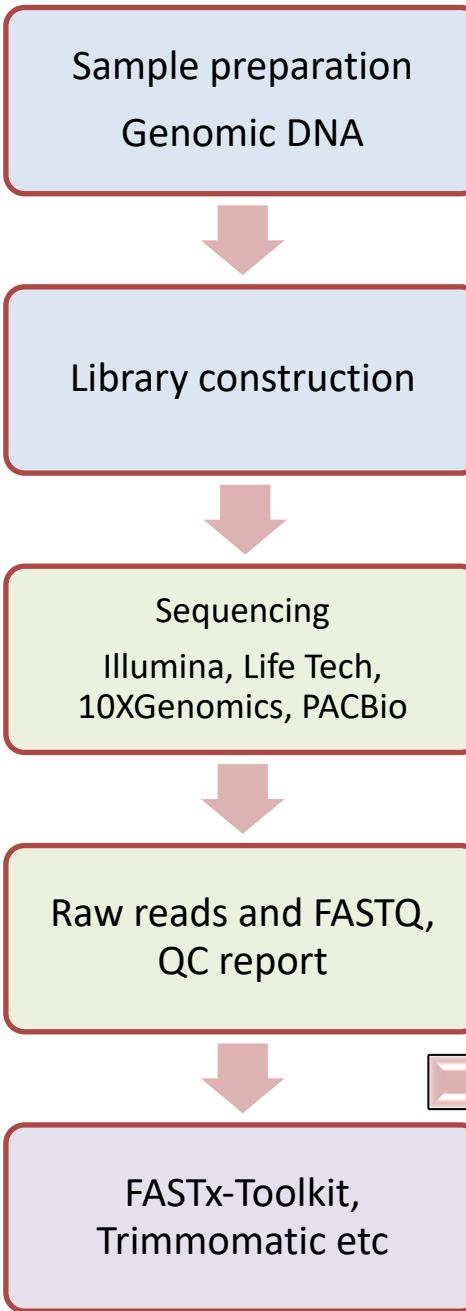
Adapters attach flow cells for  
cluster formation

Animation of exome sequencing

<https://www.youtube.com/watch?v=womKfikWIxM>

# I have my sequences/images. Now what?





## Quality Control and Preprocessing

### FASTQ format “Raw Reads”

Diagram illustrating the structure of a FASTQ file:

- Header**: @HWI-BR001:994:B809UWABXX:1:1101:13501:2240 1:N:0:CTTGTA
- Sequence**: TGAAACCAGTGTCTTAATTGGCATTACACACACACAGAATTAAAAAAAATCAAAGG
- Q-score**: +55>7 ; : : BDADDD@EE88DCD?DFFFEFFECBE6666BB=B;<;<-34 : ; <CB51>=BBEE>EE?
- ASCII**: CCAAAACATTAAGTAACTCTTAAATGGCACACAGGTTAAAGCTATTGGTTTCCCTTAACCT

The sequence continues with multiple lines of header, sequence, Q-score, and ASCII data, representing a single forward or reverse read from a sequencer.

12 million to 3 Billion in a run

# *Phred Score*

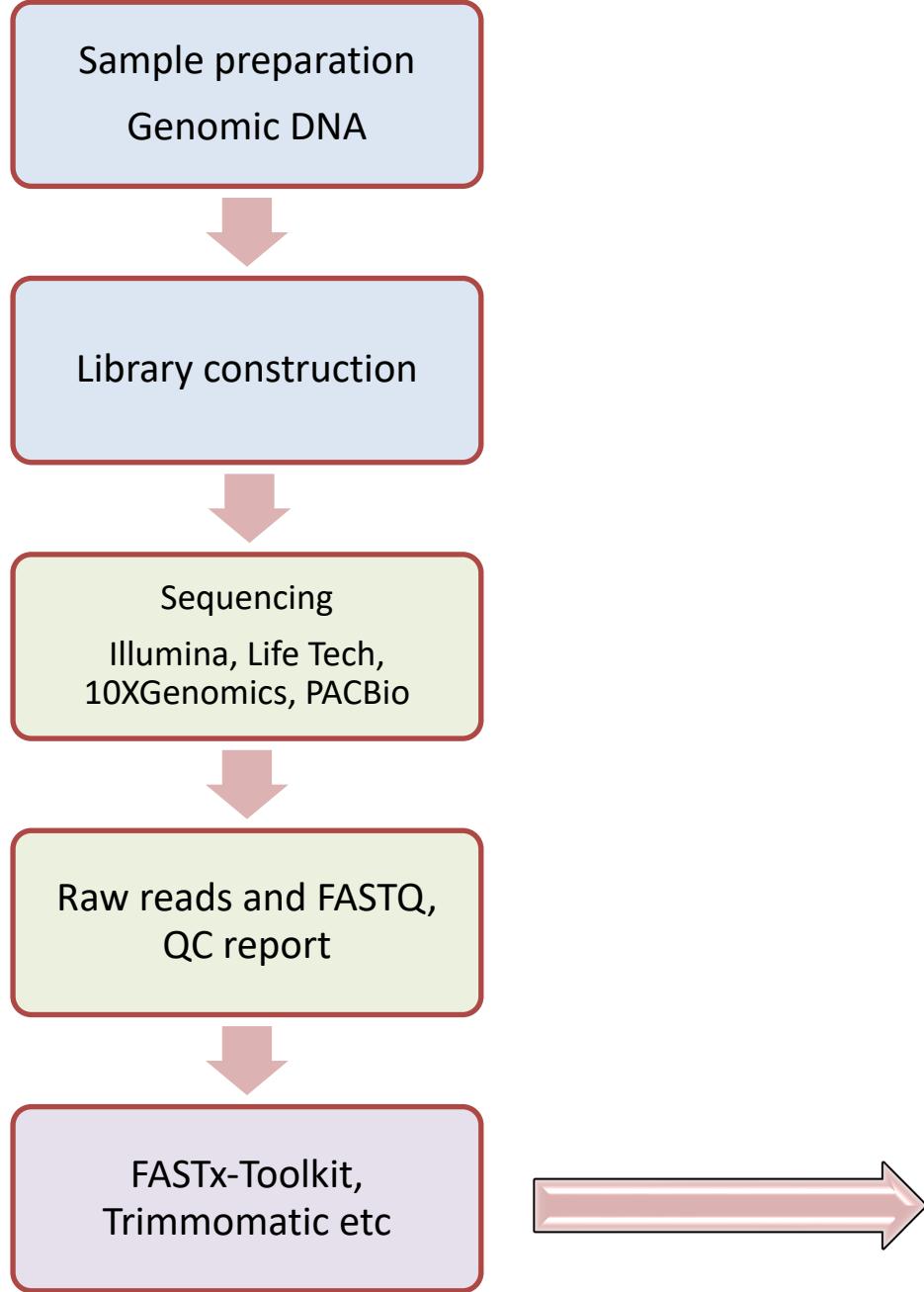
- Phred quality scores are defined as a property which is logarithmically **related to the base-calling error probabilities**

## **Phred qualities**

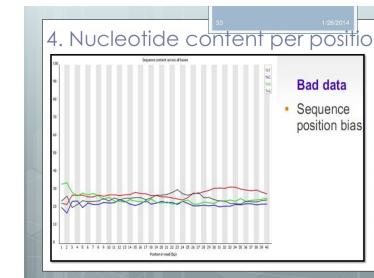
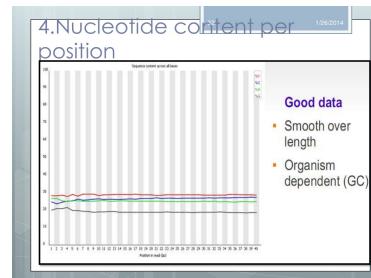
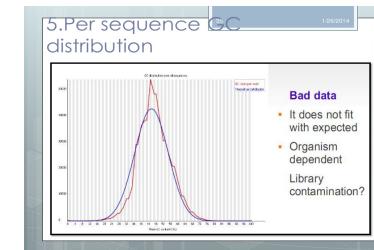
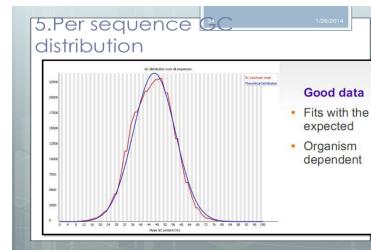
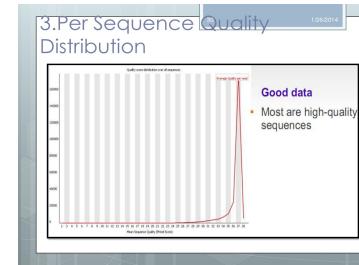
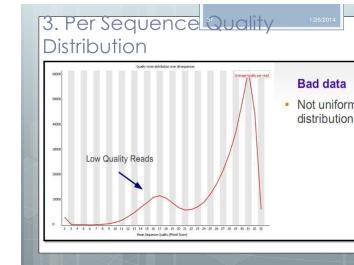
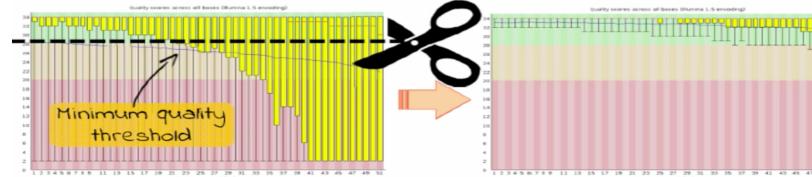
Quality value	Chance it is wrong	Accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%

- $Q = -10 \log_{10} P \iff P = 10^{-Q/10}$ 
  - $Q$  = Phred quality score
  - $P$  = probability of base call being incorrect

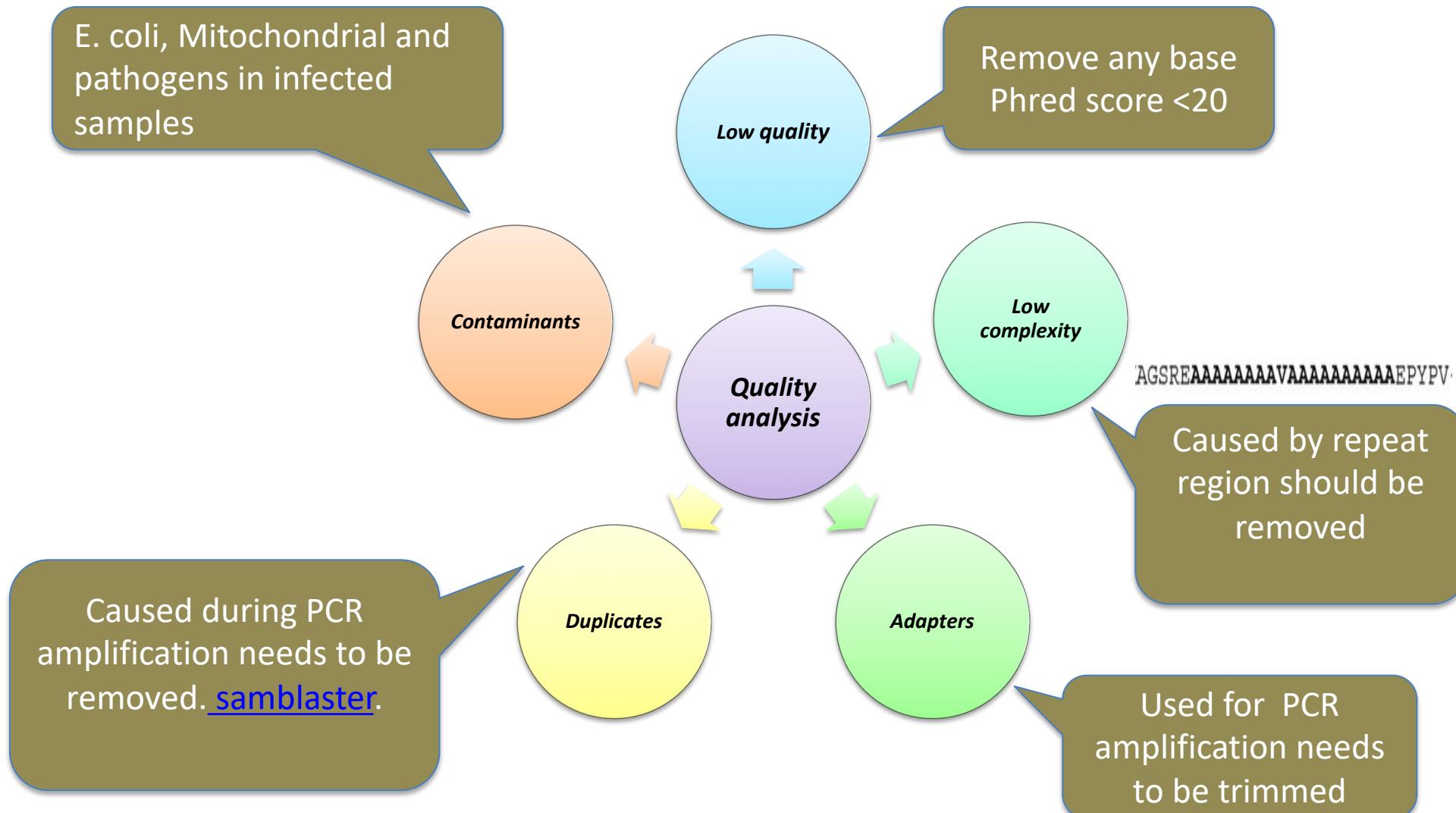
# Quality Control and Preprocessing



## Filtering and trimming



# *Quality Control and Preprocessing*



## SAM/BAM Format files

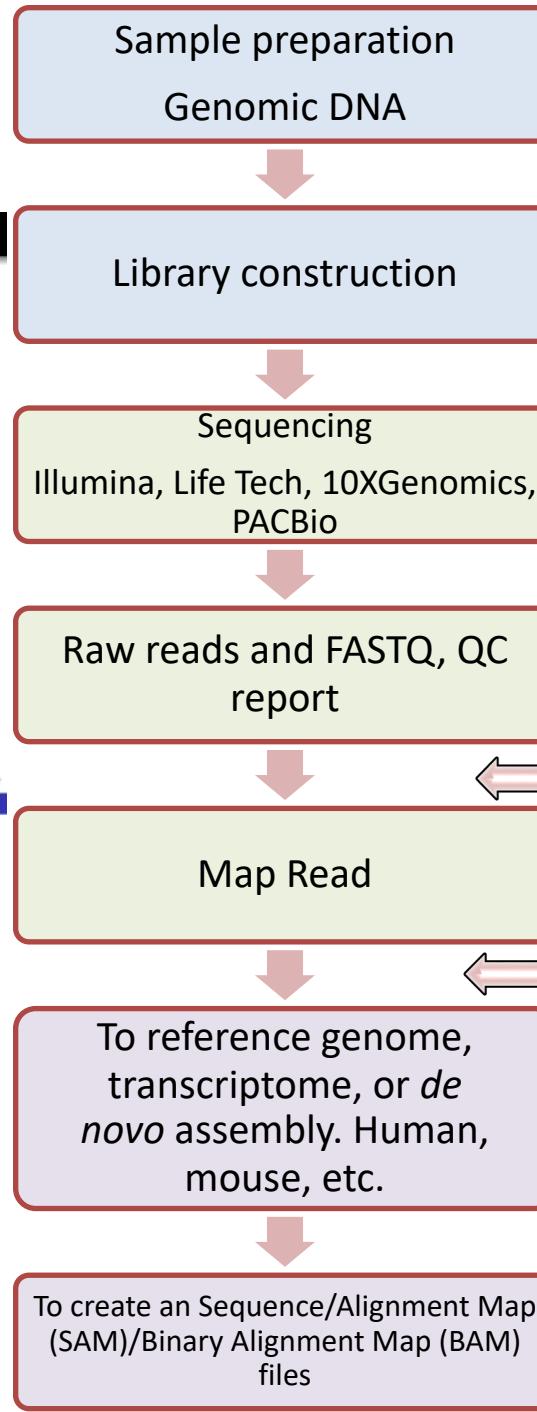
### SAM Format

- Is flexible enough to store all the alignment information generated by various alignment programs;
- Is simple enough to be easily generated by alignment programs or converted from existing alignment formats;
- Is compact in file size;
- Allows most of operations on the alignment to work on a stream without loading the whole alignment into memory;
- Allows the file to be indexed by genomic position to efficiently retrieve all reads aligning to a locus.

Pierre Lindenbaum@yokofakun pierre.lindenbaum@univ-nantes.fr Next Generation Sequencing File Formats.

### BAM File

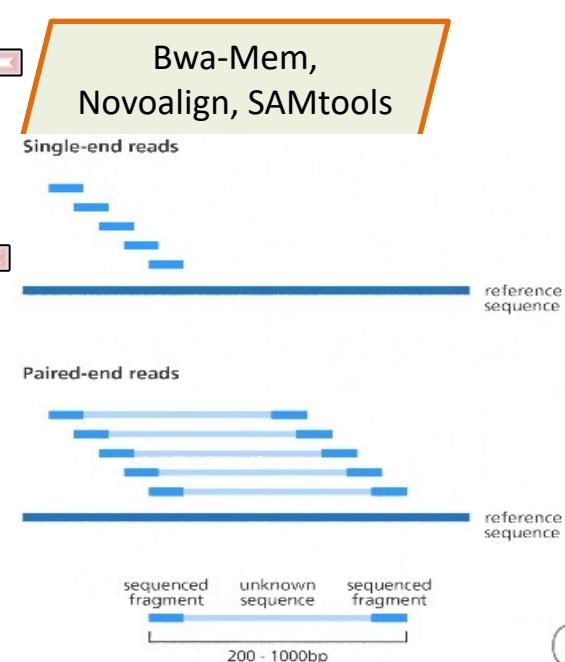
- SAM (Sequence Alignment/Map) files are typically converted to Binary Alignment/Map (BAM) format
- They are binary
  - Can NOT be opened like a text file
  - Compressed = Smaller (great for server storage)
  - Sorted and indexed = Faster
- Once you have a file in a BAM format you can delete your aligned read files
  - You can recover the FASTQ reads from the BAM format
  - Delete the aligned reads – trimmed/filtered reads
  - Do not delete the RAW reads – zip them and save them



## Alignments and Mapping of raw reads

### Algorithms

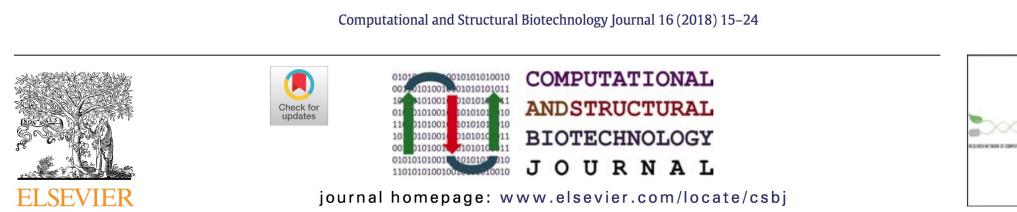
Bwa-Mem,  
Novoalign, SAMtools



# *SAM/BAM Files Format*

- SAM (Sequence Alignment/Map) format
  - Single unified format for storing read alignments to a reference genome
- BAM (Binary Alignment/Map) format
  - Binary equivalent of SAM
  - Advantages
    - Supports indexing
    - Compact size

# Variant Calling Algorithms



A review of somatic single nucleotide variant calling algorithms for next-generation sequencing data

Chang Xu

Life Science Research and Foundation, Qiagen Sciences, Inc., 6951 Executive Way, Frederick, Maryland 21703, USA

## ARTICLE INFO

Article history:  
Received 8 September 2017  
Received in revised form 20 January 2018  
Accepted 28 January 2018  
Available online 6 February 2018

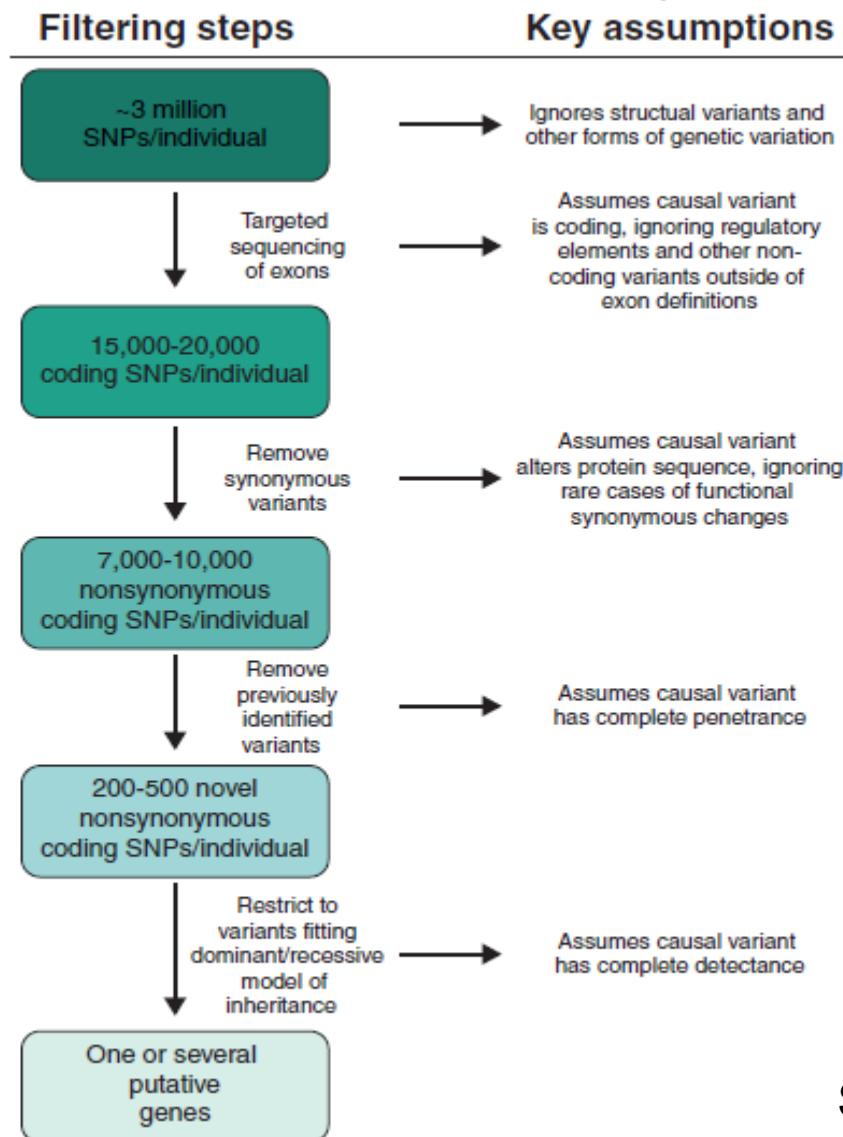
Keywords:  
Variant calling  
Somatic mutation  
Unique molecular identifier  
Low-frequency mutation  
Benchmarking

## ABSTRACT

Detection of somatic mutations holds great potential in cancer treatment and has been a very active field in the past few years, especially since the breakthrough of the next-generation sequencing technology. A collection of variant calling pipelines have been developed with different underlying models, filter data requirements, and targeted applications. This review aims to enumerate these unique features of state-of-the-art variant callers, in the hope to provide a practical guide for selecting the appropriate pipeline for specific applications. We will focus on the detection of somatic single nucleotide variants, ranging from traditional variant callers based on whole genome or exome sequencing of paired tumor-normal samples to recent low-frequency variant callers designed for targeted sequencing protocols with unique identifiers. The variant callers have been extensively benchmarked with inconsistent performances in these studies. We will review the reference materials, datasets, and performance metrics that have been used in the benchmarking studies. In the end, we will discuss emerging trends and future directions of variant calling algorithms.

Variant caller	Type of variant	Single-sample mode	Type of core algorithm
BAYSIC [48]	SNV	No	Machine learning (ensemble caller)
CaVEMan [34]	SNV	No	Joint genotype analysis
deepSNV [38]	SNV	No	Allele frequency analysis
EBCall [37]	SNV, indel	No	Allele frequency analysis
FaSD-somatic [31]	SNV	Yes	Joint genotype analysis
FreeBayes [44]	SNV, indel	Yes	Haplotype analysis
HapMuC [42]	SNV, indel	Yes	Haplotype analysis
JointSNVMix2 [30]	SNV	No	Joint genotype analysis
LocHap [43]	SNV, indel	No	Haplotype analysis
LoFreq [36]	SNV, indel	Yes	Allele frequency analysis
LoLoPicker [39]	SNV	No	Allele frequency analysis
MutationSeq [45]	SNV	No	Machine learning
MuSE [40]	SNV	No	Markov chain model
MuTect [35]	SNV	Yes	Allele frequency analysis
SAMtools [8]	SNV, indel	Yes	Joint genotype analysis
Platypus [41]	SNV, indel, SV	Yes	Haplotype analysis
qSNP [24]	SNV	No	Heuristic threshold
RADIA [26]	SNV	No	Heuristic threshold
Seurat [33]	SNV, indel, SV	No	Joint genotype analysis
Shimmer [25]	SNV, indel	No	Heuristic threshold
SNooPer [47]	SNV, indel	Yes	Machine learning
SNVSniffer [32]	SNV, indel	Yes	Joint genotype analysis
SOAPsnv [27]	SNV	No	Heuristic threshold
SomaticSeq [46]	SNV	No	Machine learning (ensemble caller)
SomaticSniper [28]	SNV	No	Joint genotype analysis
Strelka [17]	SNV, indel	No	Allele frequency analysis
TVC [97]	SNV, indel, SV	Yes	Ion Torrent specific

# Identifying causal variants: filtering



# Analysis of protein-coding genetic variation in 60,706 humans

Monkol Lek<sup>1,2,3,4</sup>, Konrad J. Karczewski<sup>1,2\*</sup>, Eric V. Minikel<sup>1,2,5\*</sup>, Kaitlin E. Samocha<sup>1,2,5,6\*</sup>, Eric Banks<sup>2</sup>, Timothy Fennell<sup>2</sup>, Anne H. O'Donnell-Luria<sup>1,2,7</sup>, James S. Ware<sup>2,8,9,10,11</sup>, Andrew J. Hill<sup>1,2,12</sup>, Beryl B. Cummings<sup>1,2,5</sup>, Taru Tukiainen<sup>1,2</sup>, Daniel P. Birnbaum<sup>2</sup>, Jack A. Kosmicki<sup>1,2,6,13</sup>, Laramie E. Duncan<sup>1,2,6</sup>, Karol Estrada<sup>1,2</sup>, Fengmei Zhao<sup>1,2</sup>, James Zou<sup>2</sup>, Emma Pierce-Hoffman<sup>1,2</sup>, Joanne Berghout<sup>14,15</sup>, David N. Cooper<sup>16</sup>, Nicole Deflaux<sup>17</sup>, Mark DePristo<sup>18</sup>, Ron Do<sup>19,20,21,22</sup>, Jason Flannick<sup>2,23</sup>, Menachem Fromer<sup>1,6,19,20,24</sup>, Laura Gauthier<sup>18</sup>, Jackie Goldstein<sup>1,2,6</sup>, Namrata Gupta<sup>2</sup>, Daniel Howrigan<sup>1,2,6</sup>, Adam Kiezun<sup>18</sup>, Mitja I. Kurki<sup>2,25</sup>, Ami Levy Moonshine<sup>18</sup>, Pradeep Natarajan<sup>2,26,27,28</sup>, Lorena Orozco<sup>29</sup>, Gina M. Peloso<sup>2,27,28</sup>, Ryan Poplin<sup>18</sup>, Manuel A. Rivas<sup>2</sup>, Valentin Ruano-Rubio<sup>18</sup>, Samuel A. Rose<sup>6</sup>, Douglas M. Ruderfer<sup>19,20,24</sup>, Khalid Shakir<sup>18</sup>, Peter D. Stenson<sup>16</sup>, Christine Stevens<sup>2</sup>, Brett P. Thomas<sup>1,2</sup>, Grace Tiao<sup>18</sup>, Maria T. Tusie-Luna<sup>30</sup>, Ben Weisburd<sup>2</sup>, Hong-Hee Won<sup>31</sup>, Dongmei Yu<sup>6,25,27,32</sup>, David M. Altshuler<sup>2,33</sup>, Diego Ardiissino<sup>34</sup>, Michael Boehnke<sup>35</sup>, John Danesh<sup>36</sup>, Stacey Donnelly<sup>2</sup>, Roberto Elosua<sup>37</sup>, Jose C. Florez<sup>2,26,27</sup>, Stacey B. Gabriel<sup>2</sup>, Gad Getz<sup>18,26,38</sup>, Stephen J. Glatt<sup>39,40,41</sup>, Christina M. Hultman<sup>42</sup>, Sekar Kathiresan<sup>2,26,27,28</sup>, Markku Laakso<sup>43</sup>, Steven McCarroll<sup>6,8</sup>, Mark I. McCarthy<sup>44,45,46</sup>, Dermot McGovern<sup>47</sup>, Ruth McPherson<sup>48</sup>, Benjamin M. Neale<sup>1,2,6</sup>, Aarno Palotie<sup>1,2,5,49</sup>, Shaun M. Purcell<sup>19,20,24</sup>, Danish Saleheen<sup>50,51,52</sup>, Jeremiah M. Scharf<sup>2,6,25,27,32</sup>, Pamela Sklar<sup>19,20,24,53,54</sup>, Patrick F. Sullivan<sup>55,56</sup>, Jaakko Tuomilehto<sup>57</sup>, Ming T. Tsuang<sup>58</sup>, Hugh C. Watkins<sup>44,59</sup>, James G. Wilson<sup>60</sup>, Mark J. Daly<sup>1,2,6</sup>, Daniel G. MacArthur<sup>1,2</sup> & Exome Aggregation Consortium†

# *Variants Annotation and Interpretation*

## Level I

Annotation and analysis of individual genetic alterations

## Example tools

SNPeff VEP

MuSiC MutSig

ANNOVAR SIFT

Oncodrive TieDIE

PolyPhen2 CHASM

HotNet PathScan

MutationAssessor

Dendrix MEMo

ActiveDriver

PARADIGM

## Level II

Population-based analysis of genetic alterations and identification of significant alterations, genes, pathways and networks

Ding L et al. (2014) Nat Rev Genet.

## Sample preparation

Genomic DNA

## Library construction

## Sequencing

Illumina, Life Tech, 10XGenomics, PACBio

## Raw reads and FASTQ, QC report

## Map Read

## Possible mutations and indels

## Annotated and filtered functional variants interpretation

## Algorithms

Bwa-Mem, Novoalign, SAMtools

MuTect2, Muse, VarScan, GATK

Pathogenicity Pred. Disease DBs

PolyPhen OMIM

SIFT HGMD

Splice Site prediction Gene Tests

# *Variant Annotation*

- Gene/intergenic annotation of the variants
- Exonic, intronic, 3'UTR, 5'UTR, promoter region, conserved transcription factor binding sites, conserved intergenic region, and active transcription region identified using ChIP-Seq experiments.
- Silent, missense, non-sense, frame-shift, in-frame annotation of genic variants.
- Comparison with OncoMD somatic mutation and germline collection.
- Comparison with OMIM, SNPedia and other relevant databases
- Comparison with common SNP databases: 1000-genomes project,
- dbSNP 135, and personal genome variants.

# **Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology**

Sue Richards, PhD<sup>1</sup>, Nazneen Aziz, PhD<sup>2,16</sup>, Sherri Bale, PhD<sup>3</sup>, David Bick, MD<sup>4</sup>, Soma Das, PhD<sup>5</sup>, Julie Gastier-Foster, PhD<sup>6,7,8</sup>, Wayne W. Grody, MD, PhD<sup>9,10,11</sup>, Madhuri Hegde, PhD<sup>12</sup>, Elaine Lyon, PhD<sup>13</sup>, Elaine Spector, PhD<sup>14</sup>, Karl Voelkerding, MD<sup>13</sup> and Heidi L. Rehm, PhD<sup>15</sup>; on behalf of the ACMG Laboratory Quality Assurance Committee

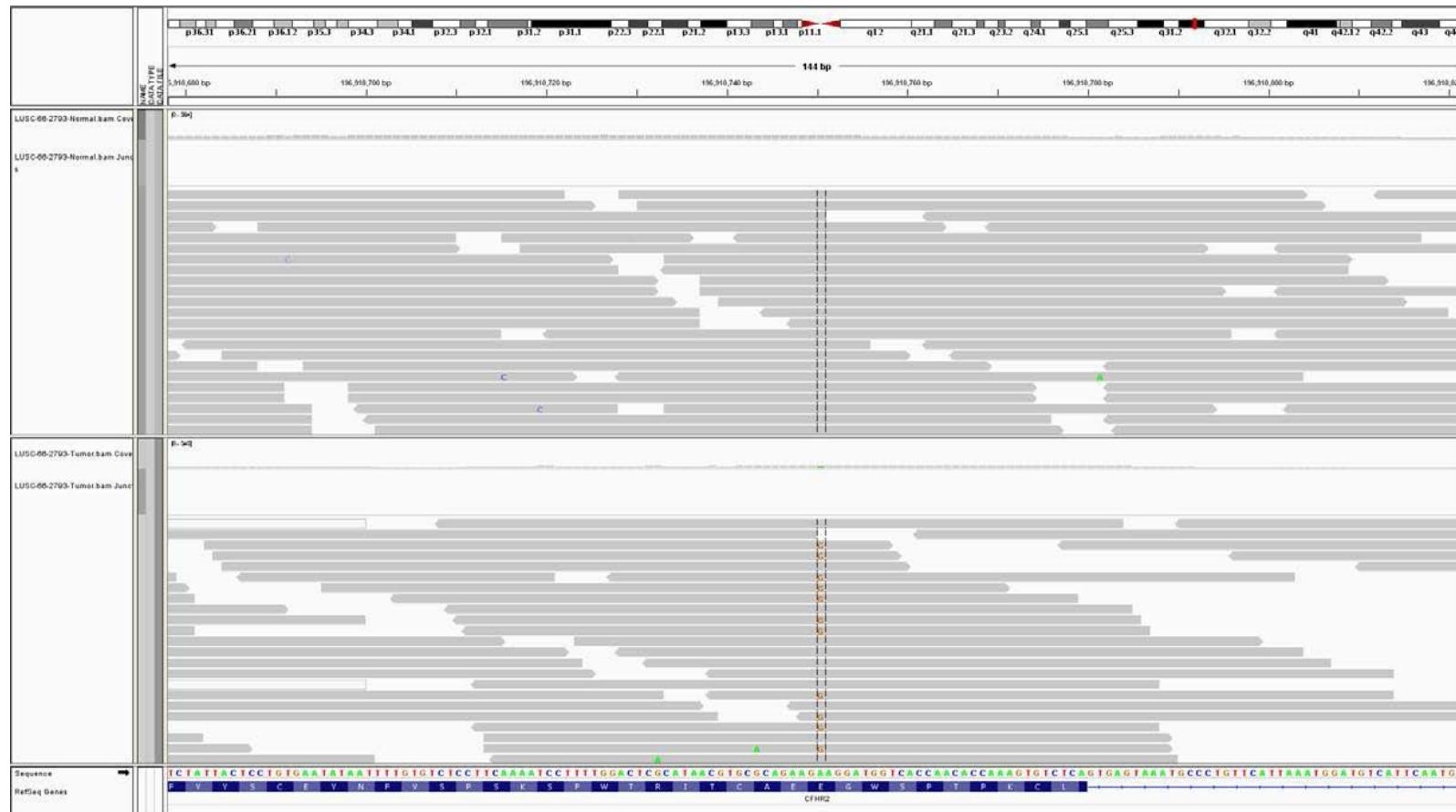
---

**Disclaimer:** These ACMG Standards and Guidelines were developed primarily as an educational resource for clinical laboratory geneticists to help them provide quality clinical laboratory services. Adherence to these standards and guidelines is voluntary and does not necessarily assure a successful medical outcome. These Standards and Guidelines should not be considered inclusive of all proper procedures and tests or exclusive of other procedures and tests that are reasonably directed to obtaining the same results. In determining the propriety of any specific procedure or test, the clinical laboratory geneticist should apply his or her own professional judgment to the specific circumstances presented by the individual patient or specimen. Clinical laboratory geneticists are encouraged to document in the patient's record the rationale for the use of a particular procedure or test, whether or not it is in conformance with these Standards and Guidelines. They also are advised to take notice of the date any particular guideline was adopted and to consider other relevant medical and scientific information that becomes available after that date. It also would be prudent to consider whether intellectual property interests may restrict the performance of certain tests and other procedures.

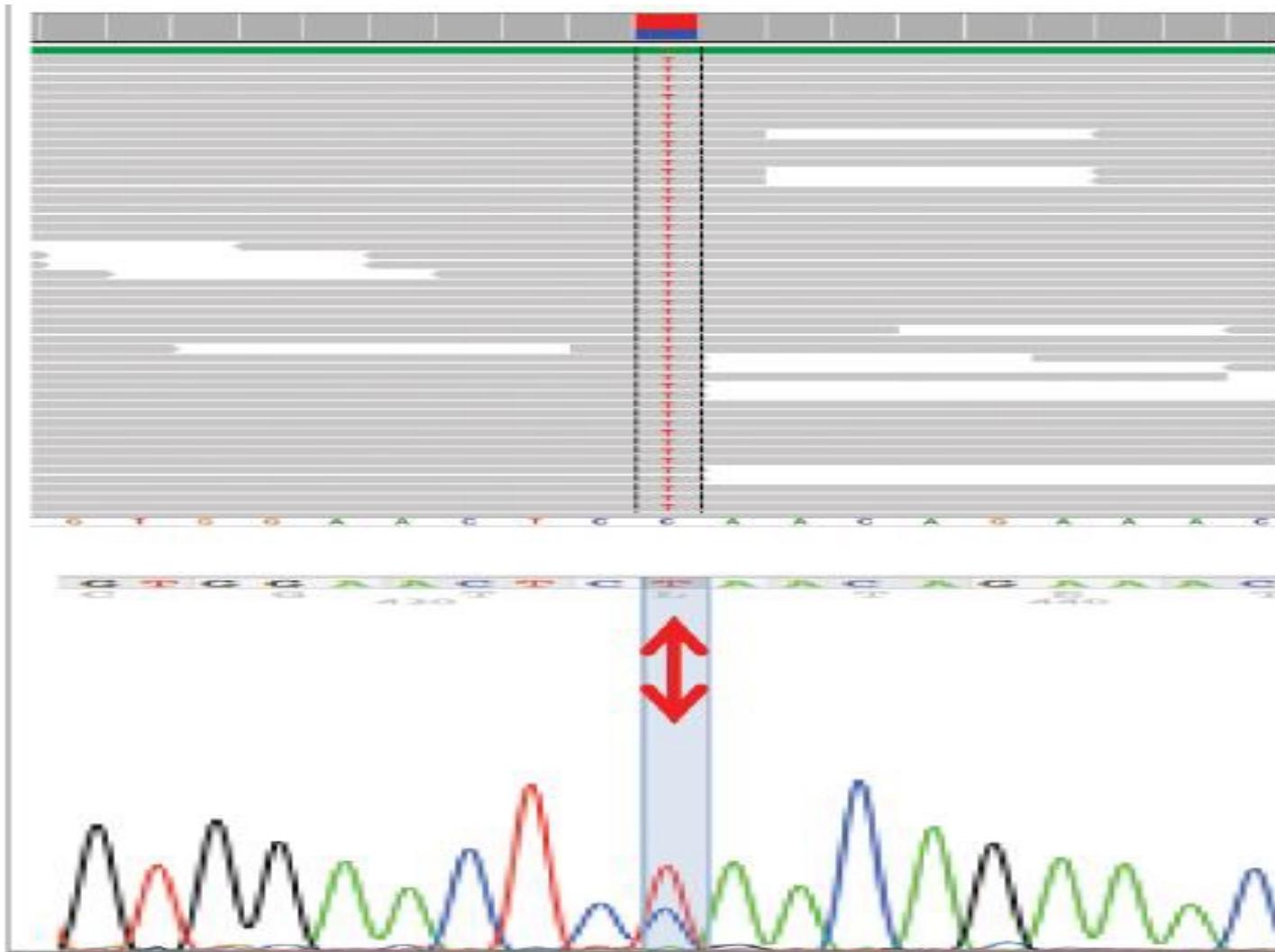
## False Negativity/Positivity

- Most false negatives are due to lack of coverage
- False positives are due to multiple reasons, including:
  - Variant is only called on one strand
  - Variant is only called at the end of the read
  - Coverage of the matched normal at that locus is poor
  - Gene has a pseudogene/paralog and the reads are mis-mapped
  - High sensitivity variant calling algorithms have elevated false positive rates to achieve detection of subclonal variants and low false negative rates

# Mutation calls from paired samples



# *Variant Calling and Validation*



Differences to the reference

Reference: C  
Sample: T

# *Gene Set Enrichment Analysis*

- Examine whether multiple genes in pre-defined sets have more mutations than expected



Ding, Nat Rev Genet. 2014

