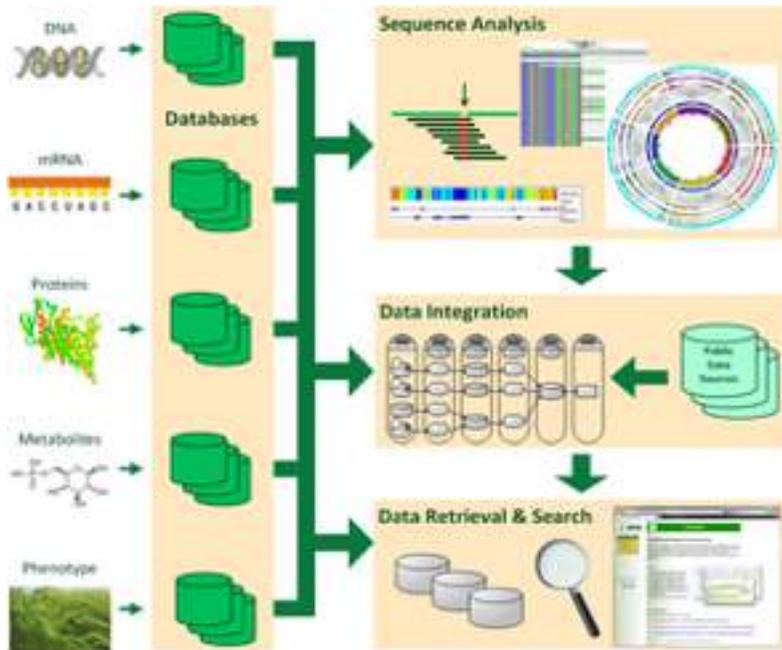


Клиническая биоинформатика
vs
биоинформатика

Alla L Lapidus, Ph.D.

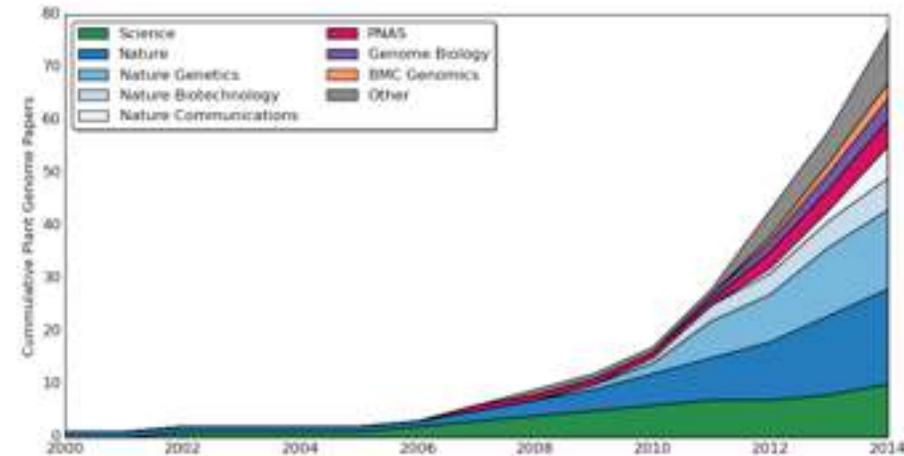
What data does bioinformatics deal with?



DNA, RNA, and protein sequences
Molecular Structures
Expression Data
Data storage
Data transfer
Data Bases

Bibliographic Data

The number of scientific articles has increased dramatically in the last few decades, due to the increasing number of research projects and genome sequencing programs. These articles are organized in public databases available online



Why bioinformatics is important

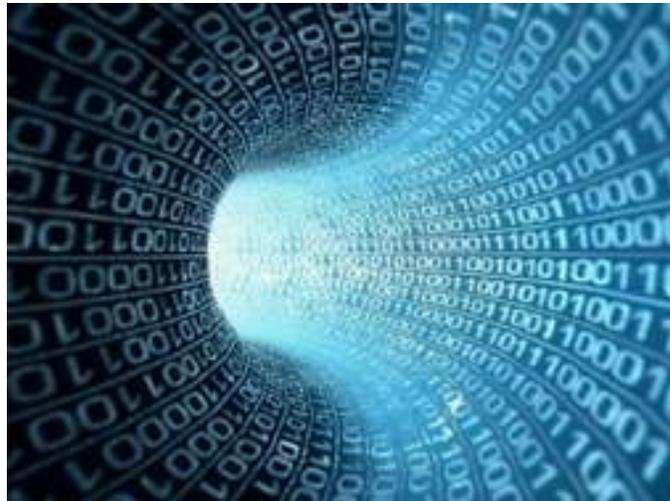
Shift from experimental laboratory bench only to
Incorporation of computers into the research process

The greatest challenges

- to make sense of the wealth of data that has been produced at **genome, transcriptome, proteome and metabalome levels**
- intelligent and efficient storage of huge amount of data generated
- to provide easy and reliable access to this data
- creation of tools to allow the extraction of meaningful biological information
- apply bioinformatics tools to reduce time and cost in molecular marker development, drug development etc

Why bioinformatics is important

The genomic era has seen a massive explosion in the amount of biological information due to huge advances in the fields of *molecular biology and new sequencing technologies*



Depends on access to high performance computing capabilities



Instruments on the market vs applications

Single gene assays

Multiplexed hotspots

Multi gene panels

Whole exome

Whole genome



ABI

Sequenom
MassArray

IonTorrent
PGM

MiSeq
Illumina

HiSeq
Illumina

Instruments on the market vs applications



“A space-based DNA sequencer could identify microbes, diagnose diseases and understand crew member health, and potentially help detect DNA-based life elsewhere in the solar system”

Bioinformatics is being used in the following fields

- ◆ Molecular medicine
- ◆ Personalised medicine
- ◆ Preventative medicine
- ◆ Gene therapy
- ◆ Drug development
- ◆ Antibiotic resistance
- ◆ Microbial genome applications
- ◆ Evolutionary studies
- ◆ Forensic analysis
- ◆ Waste cleanup

- ◆ Climate change Studies
- ◆ Alternative energy sources
- ◆ Bio-weapon creation
- ◆ Biotechnology
- ◆ Crop improvement
- ◆ Insect resistance
- ◆ Improve nutritional quality
- ◆ Development of Drought resistant varieties
- ◆ Veterinary Science

Clinical application of bioinformatics => Clinical Bioinformatics (CBI)

- helps to understand molecular mechanisms and potential therapies for human diseases
- integrates molecular and clinical data to accelerate the translation of knowledge discovery into effective treatment and personalized medicine.
- is aimed at providing methods and tools to support clinicians and researchers

Особенности биоинформатики, исследующей данные человека

Повышенные требования к:

- качеству экспериментальных данных
- уровню ответственности
- простоте и надежности создаваемых программ и аналитических подходов

Этические аспекты

Bioinformatics tools in Medical Research

Help

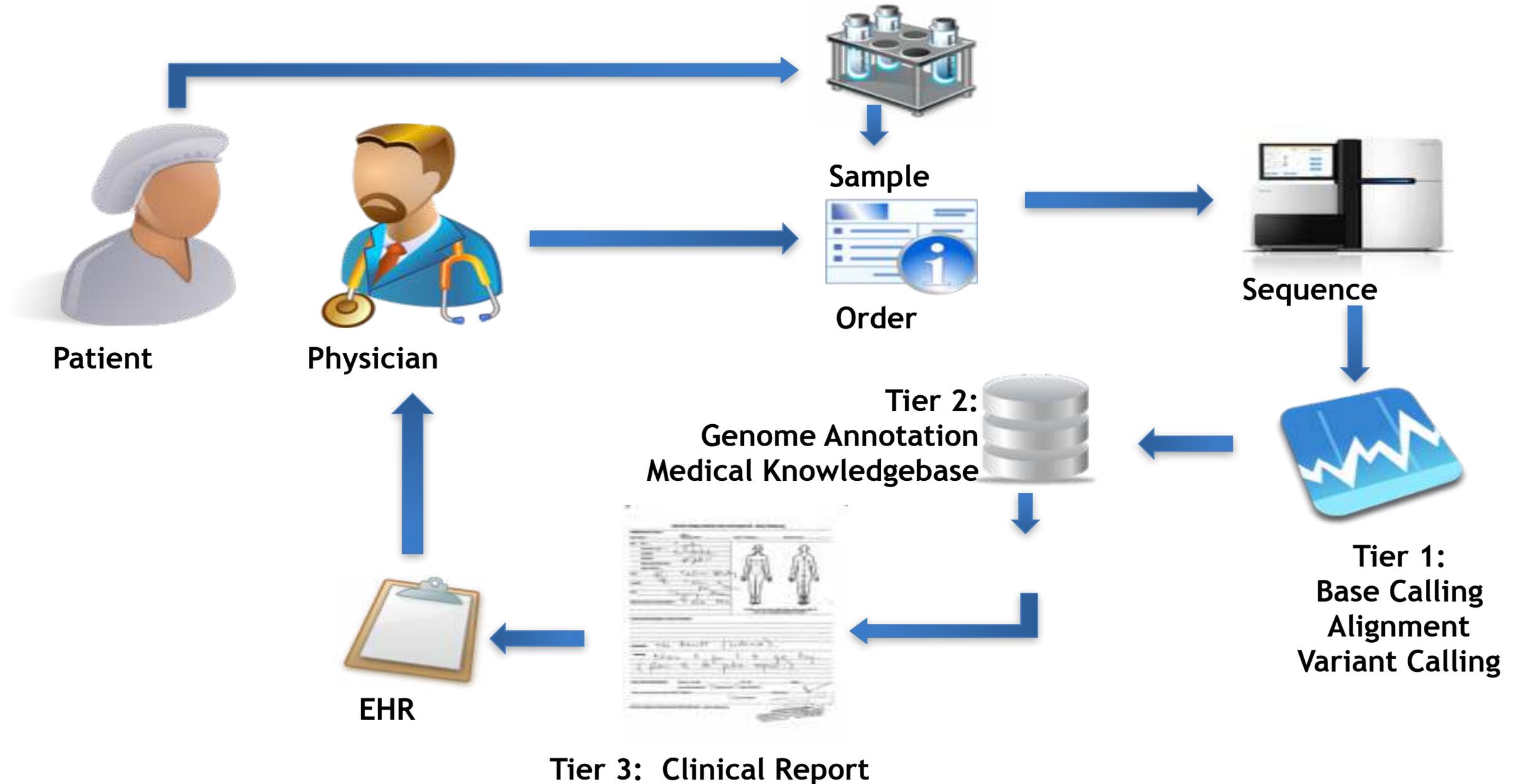
- to compare genetic and genomic data
- with understanding of various molecules that are amenable for the disease
- to analyze and document the biological systems and pathways
- to combine clinical data from patients in Electronic Medical Reports
- can affirm potential biomarkers and clinical phenotypes that allow researchers to develop experimental strategies using selected patient

Bioinformatics and Health Informatics

If bioinformatics is the study of the flow of information in biological sciences,
Health Informatics is the study of the information in patient care



Medicine: Informatics pipeline workflow



Автоматическое составление клинического отчета

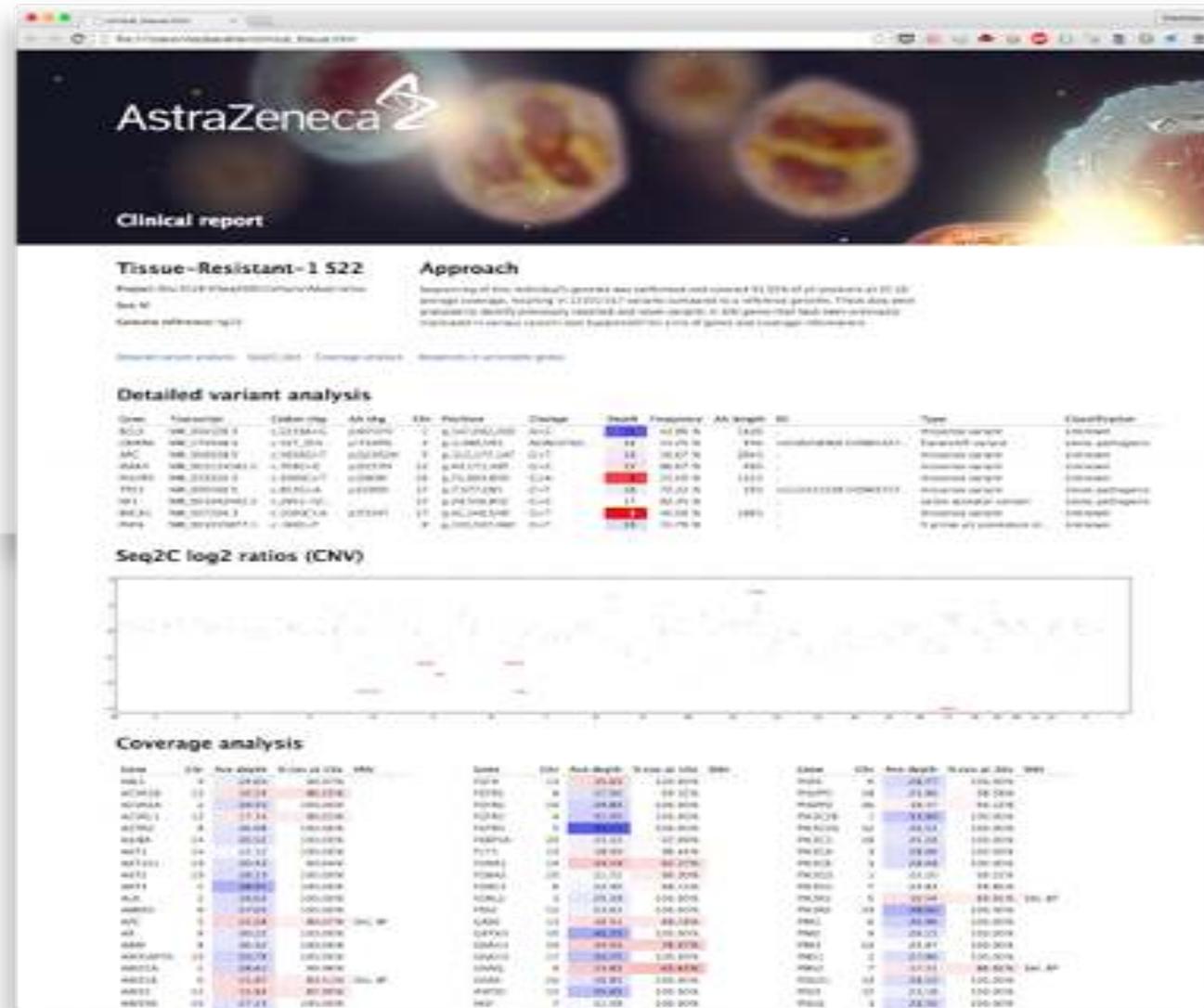
1. Информация о пациенте
2. Метод секвенирования
3. Найденные мутации и их влияние на фенотип
4. Качество покрытия ключевых для онкологии генов
5. Вариация копийности экзонов
6. Открытые исследования, где встречались такие же мутации:

Инф. о пациенте и диагнозе

Лекарства

Существующие клеточные линии

Ассоц. биомаркеры



Tools

Aligners (Bwa, MAQ, Bowtie, Eland, MOSAIK, SHRiMP, SOAP2 etc)

Mutation calling (SamTools, CASAVA, VarScan, etc)

Indel calling (Pindel, GATK, CLC, etc)

Visualization (MapView, Maqview, consed, CLC, Consed)

Annotation (SeattleSNP website, UCSC *Genome* Browser, Ensembl...)

GATK

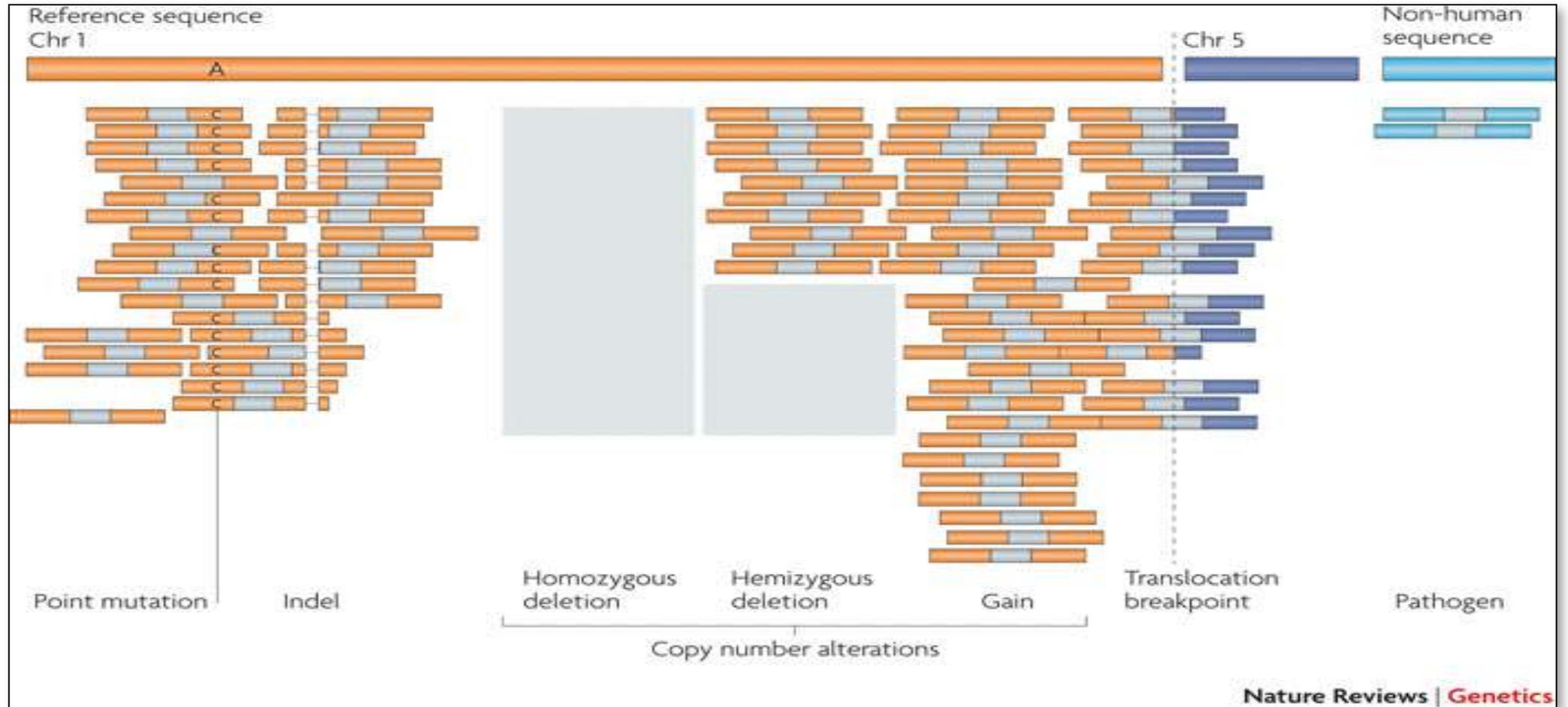
Genome Analysis Toolkit

Variant Discovery in High-Throughput Sequencing Data



Developed by the [Data Science and Data Engineering](#) group at the [Broad Institute](#), the toolkit offers a wide variety of tools with a primary focus on variant discovery and genotyping. Its powerful processing engine and high-performance computing features make it capable of taking on projects of any size.

Genome alterations



Few Mutations Can Make a Big Difference

Different people have slightly different genomes:
on average, roughly 1 mutation in 1000 nucleotides.

The 1 in 1000 nucleotides difference accounts for height, high cholesterol susceptibility, and 1000s of genetic diseases.



```
CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAGCTACC  
ACATCGTAGCTACGATGCATTAGCAAGCTATCGATCGATCGATCGATTAT  
CTACGATCGATCGATCGATCACTATACGAGCTACTACGTACGTACGATC  
GCGGGACTATTATCGACTACAGATAAAACATGCTAGTACAACAGTATACA  
TAGCTGCGGGATACGATTAGCTAATAGCTGACGATATCCGAT
```



```
CTGATGATGGACTACGCTACTACTGCTAGCTGTATTACGATCAGCTACA  
ACATCGTAGCTACGATGCATTAGCAAGCTATCGATCGATCGATCGATTAT  
CTACGATCGATCGATCGATCACTATACGAGCTACTACGTACGTACGATC  
GCGTGACTATTATCGACTACAGATGAAACATGCTAGTACAACAGTATACA  
TAGCTGCGGGATACGATTAGCTAATAGCTGACGATATCCGAT
```



Some challenges in cancer genome sequencing and analysis

Some genomic alterations in cancer are predominantly found at low frequency in clinical samples

Tools with low level of false-positive detections are needed

Normal tissue should be sequenced to identify rare somatic mutations

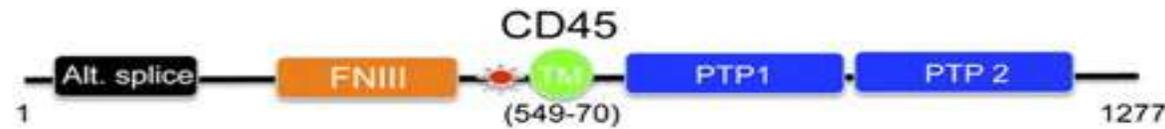
Whole genome sequence and de-novo assemblies are needed to analyze chromosomal rearrangements and somatic mutations in non-coding and not annotated regions (costly!)

Gapped reference for read alignment (genome finishing is in progress)

Cancer genomes are highly diverse and complex and vary significantly from normal genomes and from each other (mutation frequency, genome structure)

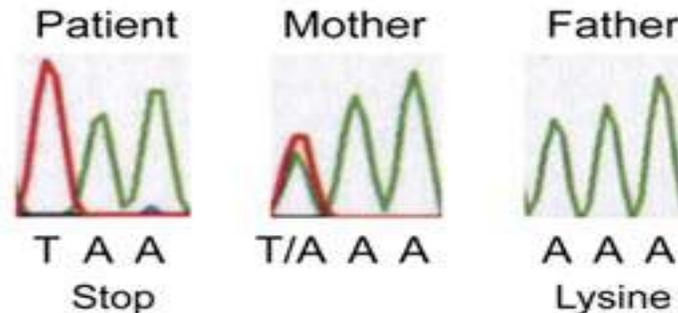
“Noisy” cancer DNA (non-tumor contamination; heterogeneity within the tumor)

Sequence analysis of the CD45 alleles of the patient and parents.



1609 - TTC CGT GTA AAA GAT WT hCD45
537 - F R V K D

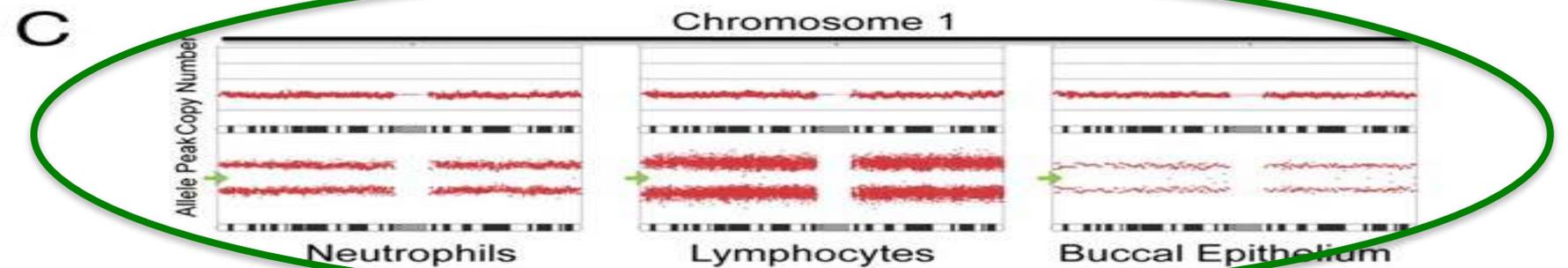
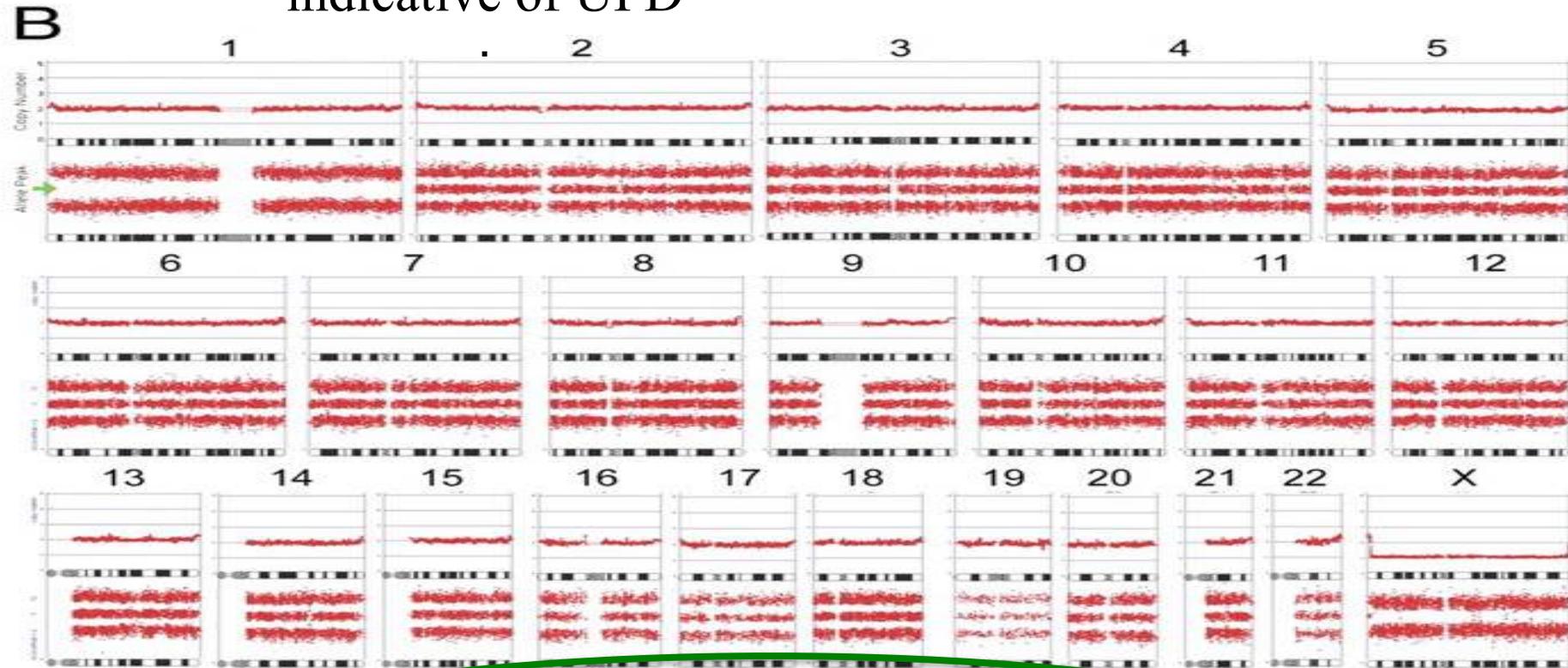
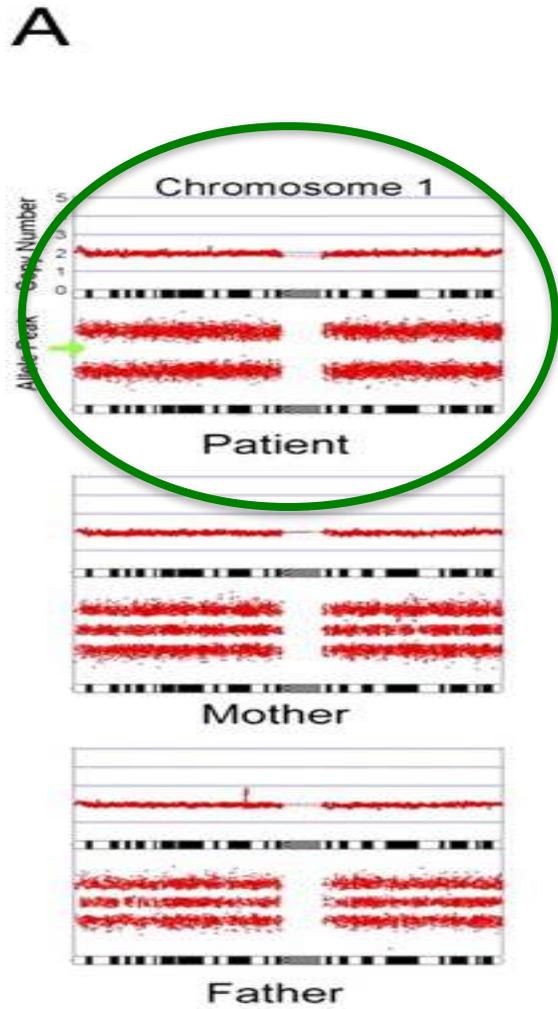
1609 - TTC CGT GTA TAA GAT Mut hCD45
537 - F R V *



CD45 gene codes tyrosine phosphatase

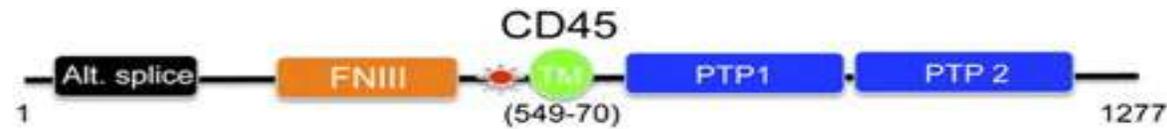
SCID = Combined Immunodeficiency Disease.

The copy number panels reveal two copies of each gene on Chr1, but show only two bands, with loss of the Middle (heterozygous) band (arrow) across the entire chromosome, indicative of UPD



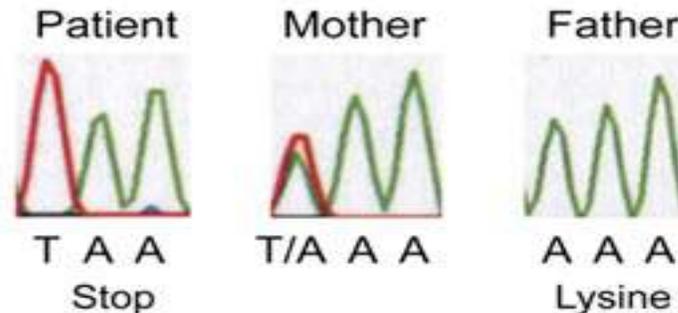
Roberts, J. L., Buckley, R. H., Luo, B., Pei, J., Lapidus, A., Peri, S., ... Wiest, D. L. (2012). CD45-deficient severe combined immunodeficiency caused by uniparental disomy. *PNAS*, 109(26), 10456–10461.

Sequence analysis of the CD45 alleles of the patient and parents



1609 - TTC CGT GTA AAA GAT WT hCD45
537 - F R V K D

1609 - TTC CGT GTA TAA GAT Mut hCD45
537 - F R V *



CD45 gene codes tyrosine phosphatase

SCID = Combined Immunodeficiency Disease.

Search for new antibiotics

Устойчивость к антибиотикам из научной и клинической давно переросла в экономическую. В год в отдельной стране число инфекций, вызванных устойчивыми к антибиотикам бактериями исчисляется миллионами. В одних только Соединенных Штатах 23 000 смертей/год

Современные методы секвенирования ДНК



е маркеры обнаруживаются в течение

нескольких часов/дней



источниками вспышек и помогают отслеживать

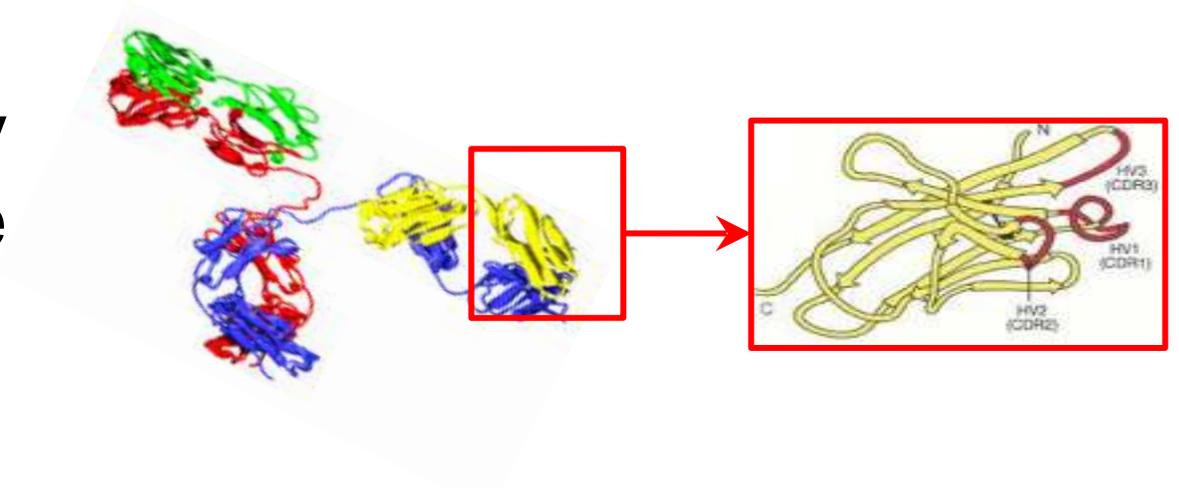
их

распространение

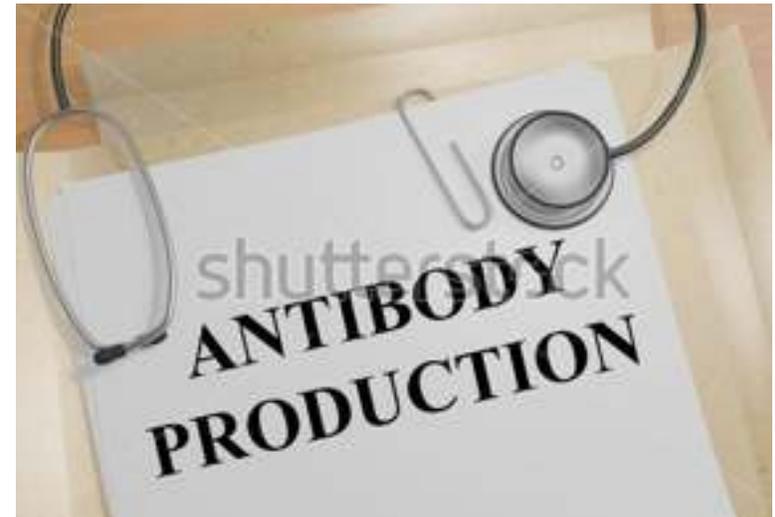
- Острая необходимость в инновационных методах диагностики
- ПОТРЕБНОСТЬ В СОЗДАНИИ АНТИБИОТИКОВ НОВОГО ПОКОЛЕНИЯ

Antibody drugs

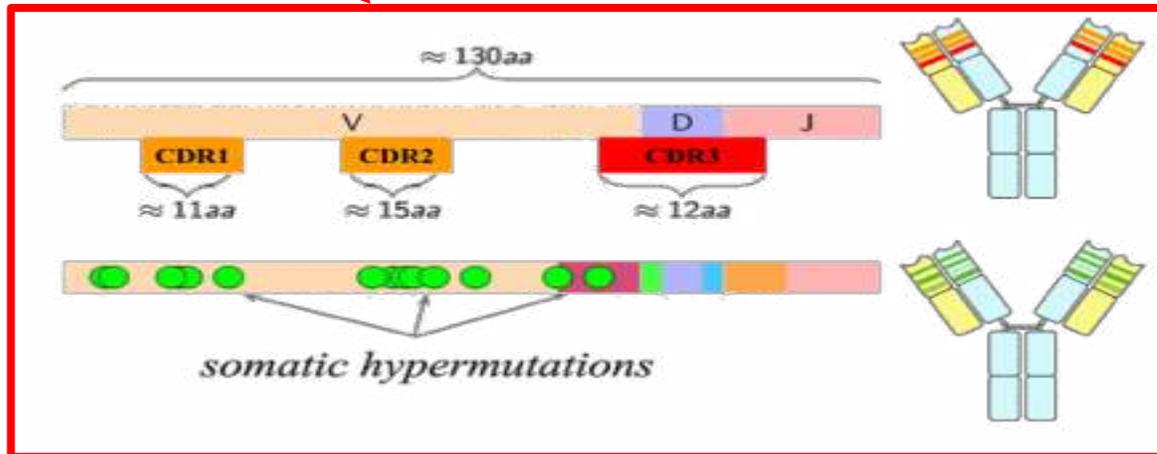
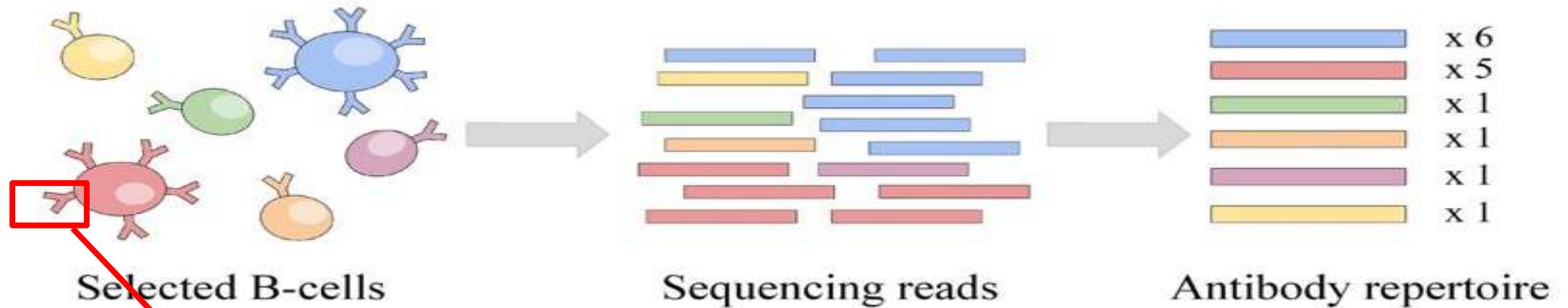
Antibodies are produced by immune system to neutralize antigens



In 2013, global sale revenue for antibody drugs was \$75 billion



Antibody repertoire



Construction of antibody repertoire from NGS data and its analysis are important steps in design of antibody drugs and clinical studies

Immunoinformatics projects (SPbU – Center for Algorithmic Biotechnology)

IgRepertoireConstructor



tool for construction of antibody repertoire and using mass spectra

IgAnalyzer & IgQUAST



tools for “quality” assessment of antibody repertoire

IgSimulator



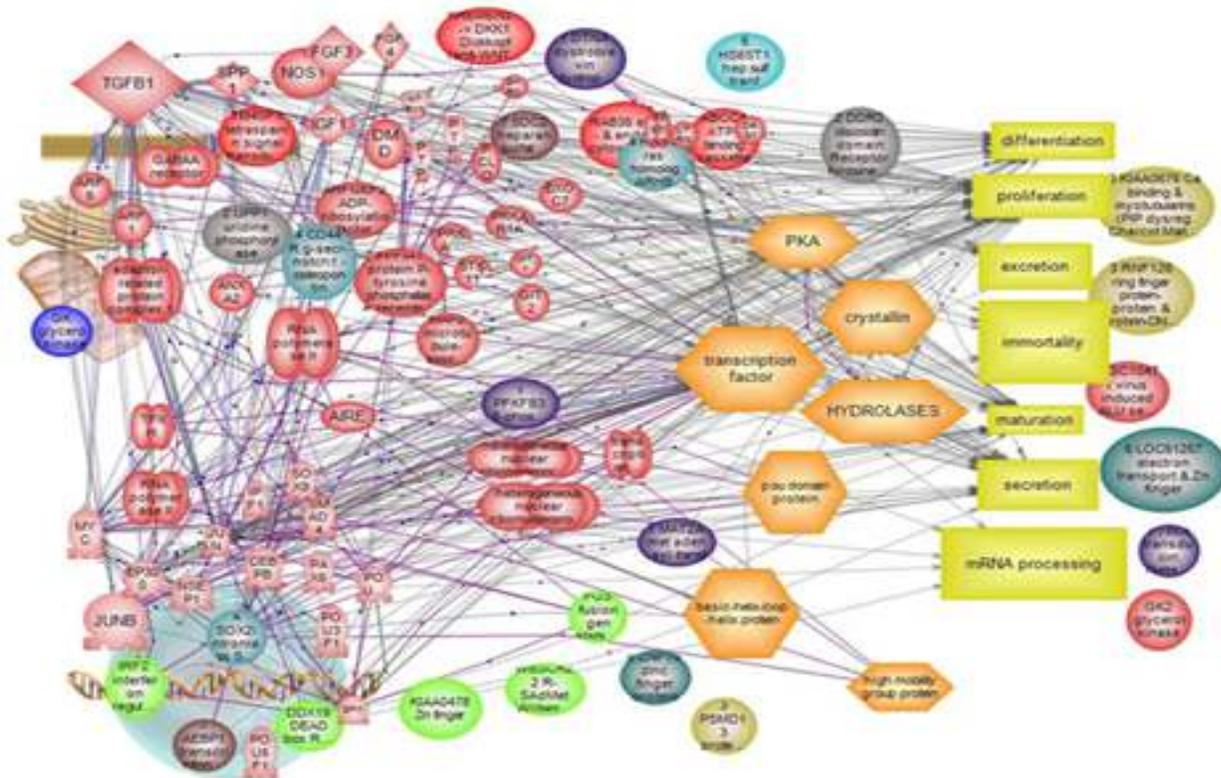
versatile immunosequencing simulator

AntEvolvo



tool for construction of clonal trees and analysis of SHMs

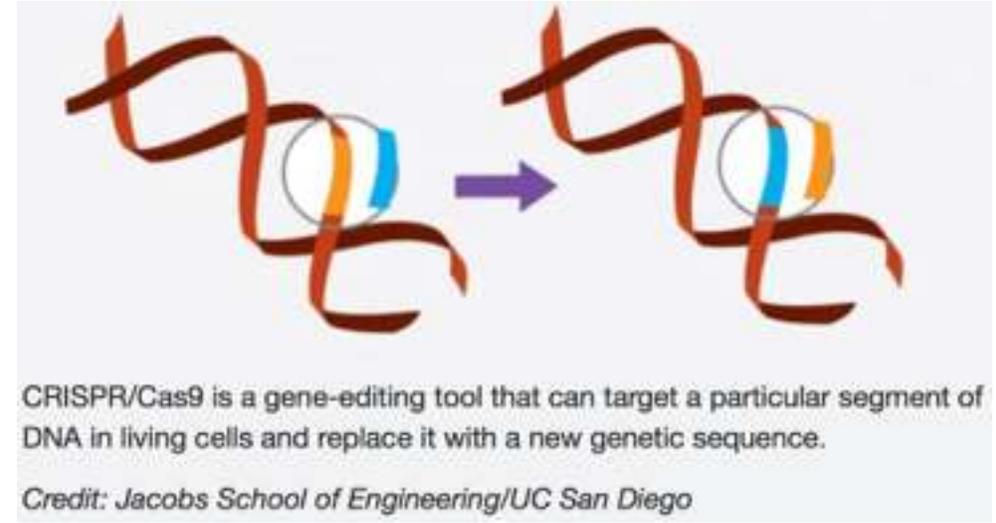
To improve drug discovery we need to discover
(read "develop")
efficient bioinformatics algorithms and approaches for



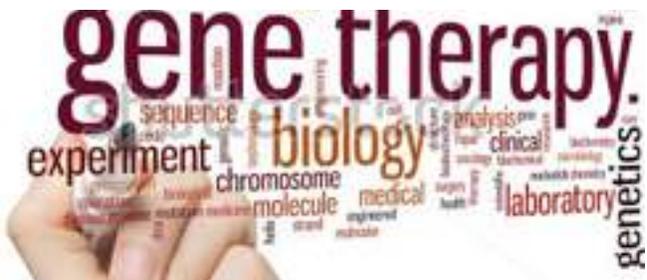
target identification
target validation
lead identification
lead optimization

Gene editing – new path to gene therapy

“**CRISPR** (*Clustered Regularly Interspaced Short Palindromic Repeats*) and its associated protein, **Cas9**, provide sequence-specific adaptive immunity in bacteria and archaea by integrating short viral DNA sequences in the host cell genome, allowing the cell to remember, recognize, and clear infections”

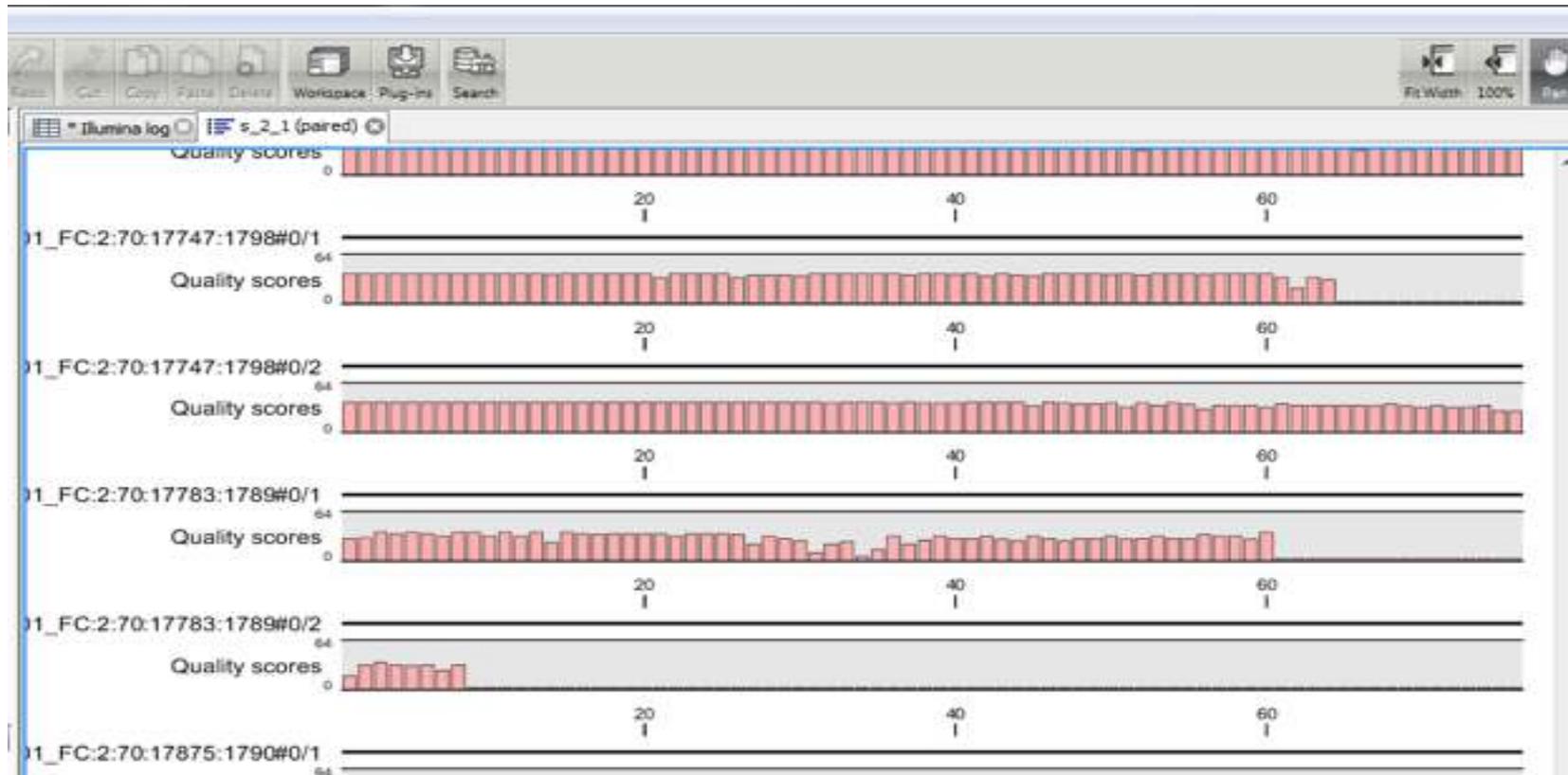


sgRNA Scorer - an interactive software for users to find guide RNAs that are predicted to be highly specific and highly active for their gene targets



QC: visualization

CLC



QC: duplications

The image shows a bioinformatics viewer interface with two windows. The top window, titled "Aligned Reads", displays a reference sequence and multiple aligned reads. A red circle highlights a specific region in the reads, indicating a duplication. The bottom window, titled "CONTIG141.C1", shows a zoomed-in view of a specific region of the reference sequence, with a red box highlighting a segment that appears to be a duplication of a smaller unit.

Aligned Reads Window:

- File: 4085752.reference.fasta.ace
- Contig: CONTIG114.C2
- Search for String: Uncompl
- Compare Cont: Find Main Min
- Err/10kb: 1.00
- Position: 528,650 to 528,740
- Consensus: AGTCCGACACCTCTTTCTATCTCGATAAAACAGAAAGCTTGGTTGAATTTAGACCGTGCCCGCGATGCAGAGGGTGTAGGGCTGTTATGGTGCAGAAA

CONTIG141.C1 Window:

- File: contig141.fasta.ace
- Contig: CONTIG141.C1
- Search for String: Compl Cont
- Compare Cont: Find Main Min
- Err/10kb: 1.00
- Position: 200,970 to 201,030
- Consensus: TTATGGCTGAATTAITGGTAAATACTGGACTTCTTGGGAAT*GGGGGGTCTATGCATATTTTAGT*

QC: Alignment quality = f (read quality)



SNP detection

The screenshot displays a bioinformatics software interface with two main panels. The top panel shows a table of SNP detection results, and the bottom panel shows a read mapping visualization.

Table Settings

- Column width: Manual
- Show column:
 - Mapping
 - Reference position
 - Consensus position
 - Variation type
 - Length
 - Reference
 - Variants
 - Allele variations
 - Frequencies
 - Counts

Table Data

Mapping	Reference ...	Consensus...	Variation t...	Reference	Variants	Allele varia...	Frequencies
chr 1 mapping	59133	20407	SNP	A	1	G	97.7
chr 1 mapping	59374	20648	SNP	A	1	G	98.0
chr 1 mapping	714073	157731	SNP	T	2	A/T	50.0/50.0
chr 1 mapping	714419	157932	SNP	G	2	A/G	40.0/40.0
chr 1 mapping	714427	157940	SNP	C	2	A/C	50.0/50.0
chr 1 mapping	714429	157942	SNP	C	2	C/G	50.0/50.0
chr 1 mapping	714822	158147	SNP	A	2	A/T	50.0/50.0
chr 1 mapping	715219	158496	SNP	G	1	T	75.0
chr 1 mapping	715225	158502	SNP	T	1	G	75.0
chr 1 mapping	715227	158504	SNP	T	1	C	75.0
chr 1 mapping	715237	158514	SNP	C	1	A	75.0
chr 1 mapping	715239	158516	SNP	C	1	G	75.0
chr 1 mapping	715424	158669	SNP	C	1	T	66.7
chr 1 mapping	715430	158675	SNP	G	1	A	80.0
chr 1 mapping	715432	158677	SNP	T	2	T/A	60.0/40.0

Read Mapping Settings

- Read layout:
 - Gather sequences at top
 - Show sequence ends
 - Find Conflict
 - Low coverage threshold: 8
 - Find Low Coverage
- Sequence layout:
 - No spacing
 - Numbers on sequences
 - Relative to: 1
 - Numbers on plus strand
 - Lock top sequence
 - Hide labels
 - Lock labels
 - Sequence label: Name
 - Compactness

Read Mapping Visualization

chr1 iTGTCGGGCATTATGGCTGTC **A** CATGGGGGAATTGGCTTTTCT

SNP
SNP

Consensus iTGTCGGGCATTATGGCTGTC **G** CATGGGGGAATTGGCTTTTCT

x:16193:2409#0/1 TGTCGGGCATTATGGCTGTC **C** CATGGGGGAATTGGCTTTTCT

46:8921:5076#0/1 TGTCGGGCATTATGGCTGTC **C** CATGGGGGAATTGGCTTTTCT

18:5154:3831#0/2 GTCGGGCATTATGGCTGTC **C** CATGGGGGAATTGGCTTTTCT

59:6081:8791#0/1 TCGGCATTATGGCTGTC **C** CATGGGGGAATTGGCTTTTCT

5:19038:8762#0/1 TCGGCATTATGGCTGTC **C** CATGGGGGAATTGGCTTTTCT

9:4985:21128#0/1 TCGGCATTATGGCTGTC **C** CATGGGGGAATTGGCTTTTCT

5:19026:8786#0/1 TCGGCATTATGGCTGTC **C** CATGGGGGAATTGGCTTTTCT

1:8620:8090#0/2 TCGGCATTATGGCTGTC **C** CATGGGGGAATTGGCTTTTCT

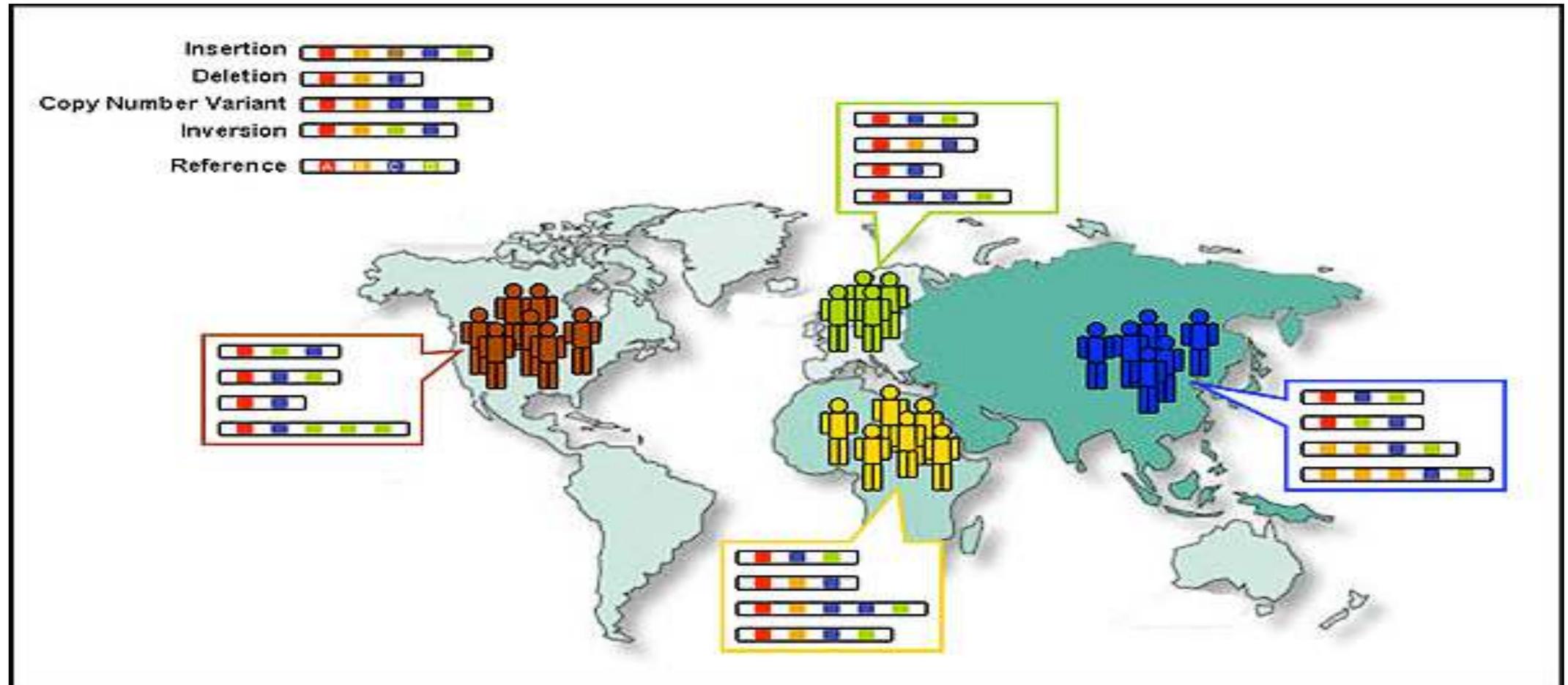
14535:13637#0/1 CGGCATTATGGCTGTC **C** CATGGGGGAATTGGCTTTTCT

15:4167:4407#0/2 CGGCATTATGGCTGTC **C** CATGGGGGAATTGGCTTTTCT

2:10357:6777#0/2 CGGCATTATGGCTGTC **C** CATGGGGGAATTGGCTTTTCT

2:5156:11218#0/2 CGGCATTATGGCTGTC **C** CATGGGGGAATTGGCTTTTCT

1000 Genomes Project



Changes in the number and order of genes (A-D) create genetic diversity within and between populations.

Sequence of Personal Genomes?

2010:

Nicholas Volker became the first human being to be saved by genome sequencing.



Doctors could not diagnose his condition; he went through dozens of surgeries.

Sequencing revealed a rare mutation in a *XIAP* gene linked to a defect in his immune system.

This led doctors to use immunotherapy, which saved the child.

Bottlenecks in bioinformatics

- education of biologists in the use of advanced computing tools
- the recruitment of computer scientists into the field of bioinformatics
- the limited availability of developed databases of biological information
- need for efficient and intelligent search engines for complex databases
- need for innovative NGS pipelines

Key segments of the Bioinformatics market by:

Platforms, Tools, and Services

Platforms

- ✓ Sequence Manipulation
- ✓ Sequence Alignment
- ✓ Structural Analysis
- ✓ Sequence Analysis

Tools

- ✓ General Knowledge
- ✓ Specialized Knowledge

Services

- ✓ Data analysis
- ✓ Sequencing services
- ✓ Database & management
- ✓ Other services

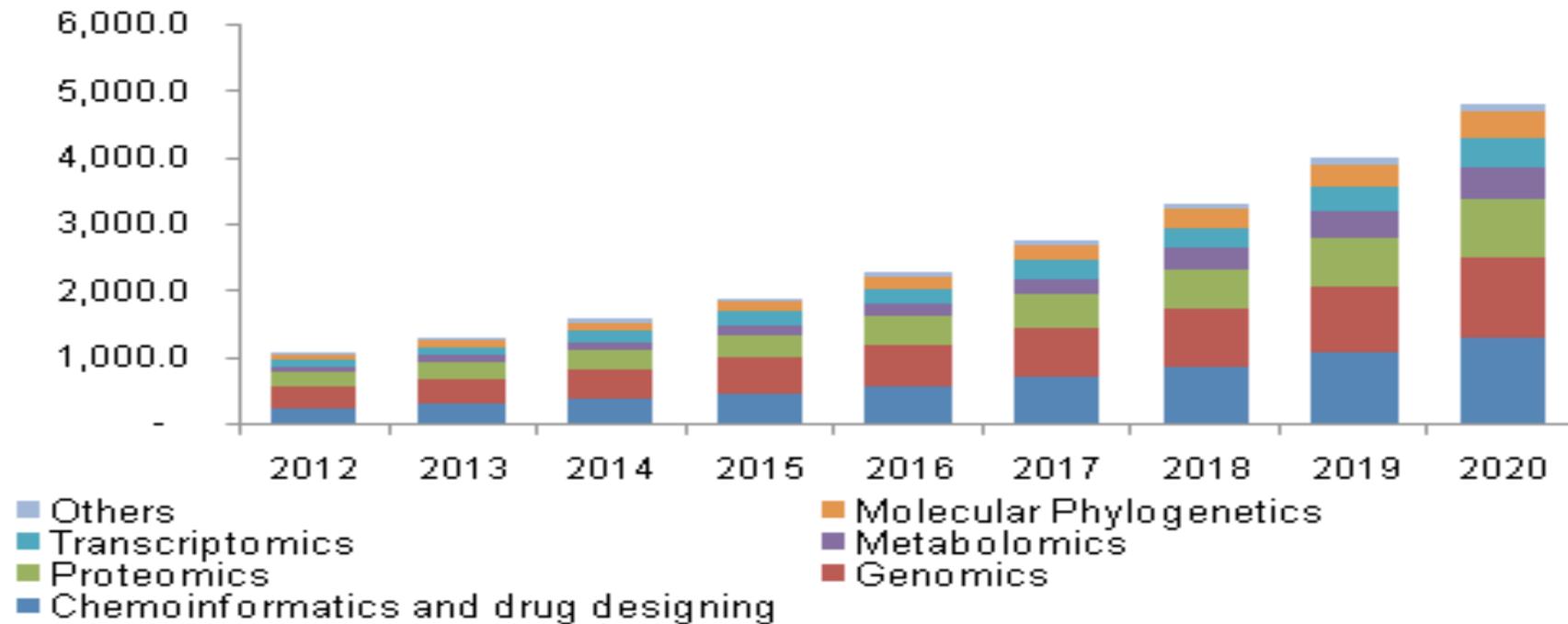
Applications

- ✓ Preventive medicine
- ✓ Molecular medicine
- ✓ Gene therapy
- ✓ Drug development
- ✓ Others

Geography

- ✓ North America (35%)
- ✓ Europe
- ✓ Asia-Pacific
- ✓ Rest of the World

European bioinformatics market, by product, 2012-2020 (USD Million)



Application Insights

Genomic was the leading applications in 2013 accounting for approximately USD 1.0 billion of the global revenue. Growing demand for pharmacogenomics in drug development and sequence screening and introduction of technological advancements aimed at managing large sets of genomic data are key factors accounting for this large share.

Chemo-informatics is expected to gain market share and grow at a CAGR of 23.5% during the forecast period wherein increasing demand for biomarker discovery and development is the primary reason accounting for this expected growth over the forecast period.

Proteomics has experienced extensive R&D investment and bioinformatics plays an integral role in analyzing and managing the resultant data. Also, this makes it simpler to handle heterogeneous and large data sets, introduce new algorithms and improve the knowledge discovery procedure.



THANK YOU!



piterlabs@gmail.com