

Genomics/DB Course

Yeoun Jin Kim, Ph.D.
CRP-Sante

- Goals
 - To understand current genomics technologies and database structure.
 - To be able to link genomics information to proteomics research.
- Schedule
 - June 3rd afternoon session.
 - June 4th afternoon session.
- Prior Knowledge
 - Molecular biology.
 - Mass spectrometry based proteomics .

Genomics/DB Course

- Basic concept
 - Molecular biology (genome structure, cloning, sequencing)
 - Bioinformatics (data analysis)
- Human Genome Project
 - What is done and what is left
 - Technologies used in HGP
- Database landscape
 - NCBI (National Center for Biotechnology Information): www.ncbi.nlm.nih.gov
 - EBI (European Bioinformatics Institute): www.ebi.ac.uk
 - UCSC Genome bioinformatics: genome.ucsc.edu
 - UniProtKB: www.uniprot.org
- Sequencing
 - G-I: Standard Sanger method
 - G-II: Automated sequencing method
 - G-III: Massively parallel sequencing methods
 - Pyrosequencing
 - Sequencing by synthesis
 - Sequencing by ligation
- Disease Association

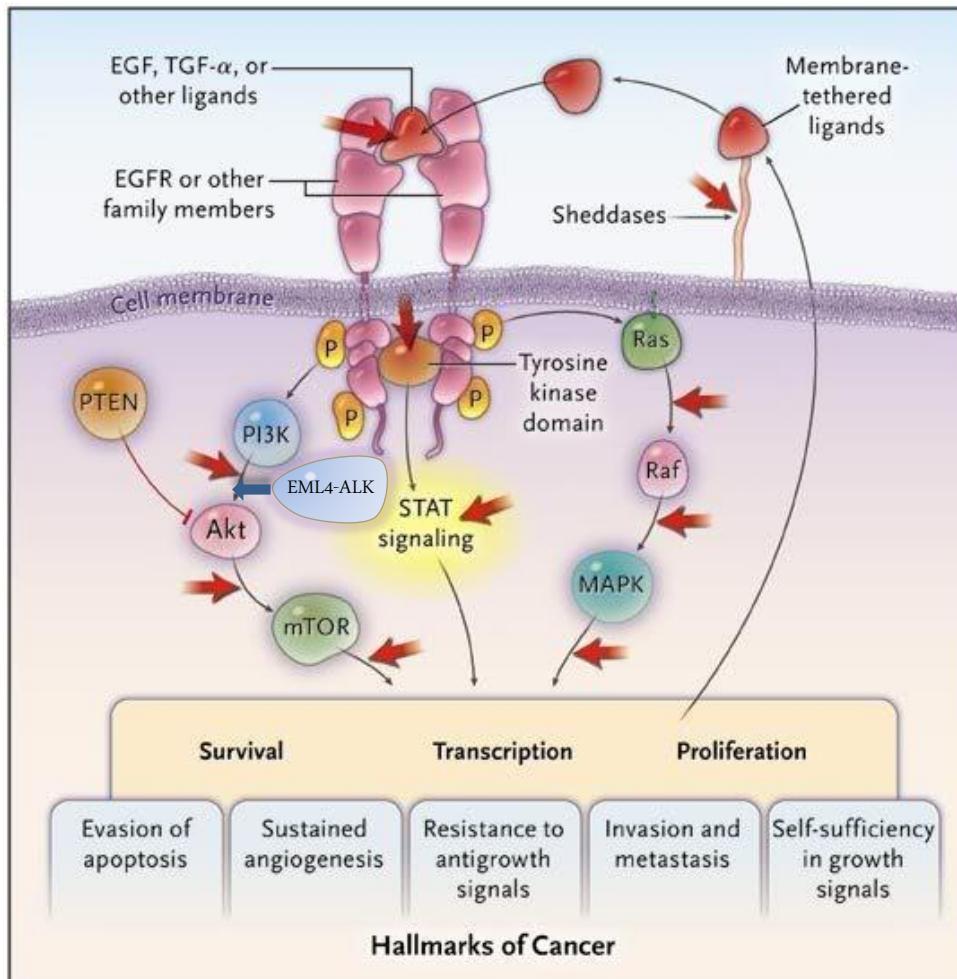
Case I

Case II

Case III

Before we start...

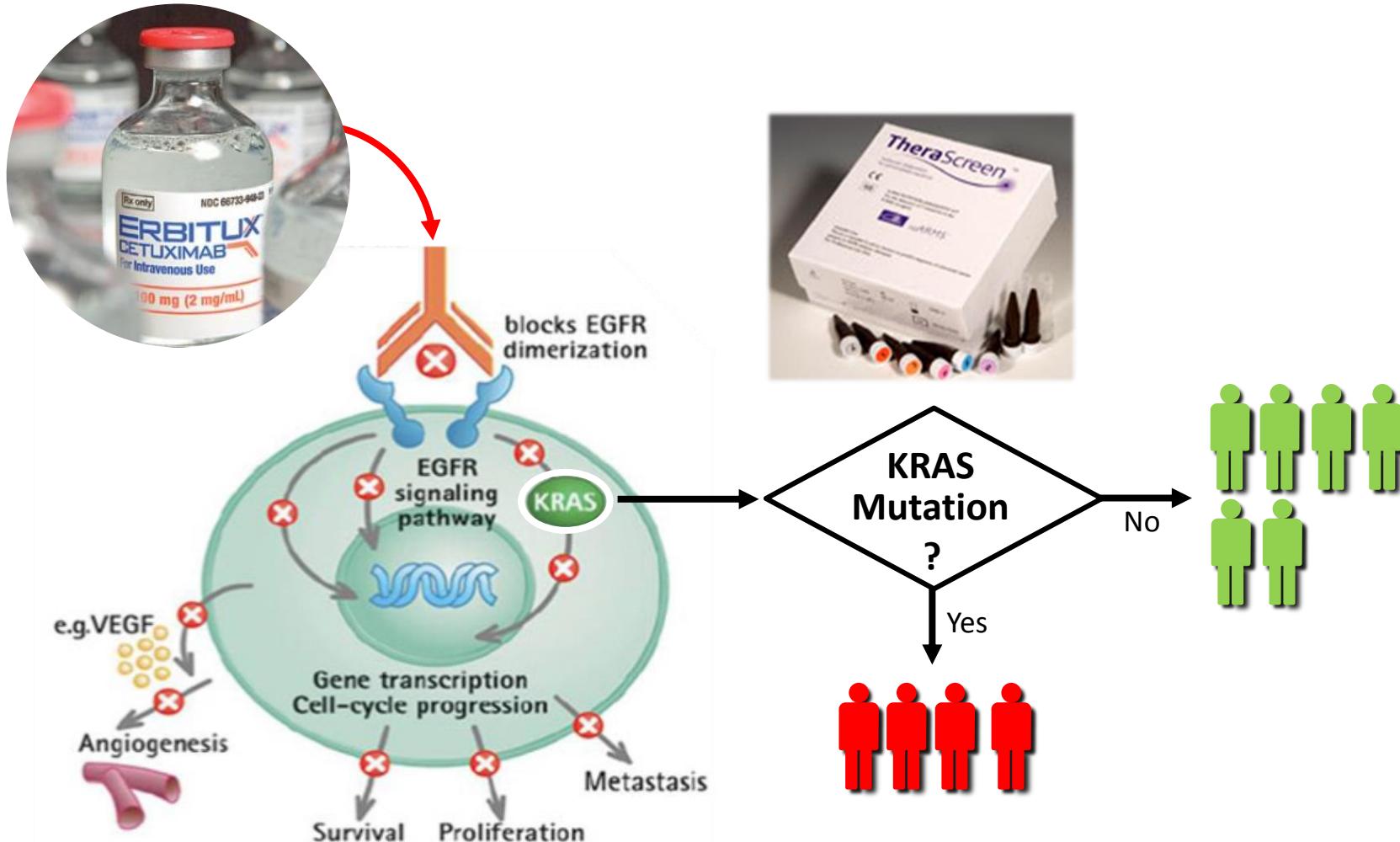
Why we study Genomics to Proteomics?



Before we start...

Why we study Genomics to Proteomics?

Erbitux (Cetuximab) BMS: In 2009, the FDA issued a class labeling change for Erbitux



Case of re-labeling

Erbitux (Cetuximab) BMS

HIGHLIGHTS OF PRESCRIBING INFORMATION

These highlights do not include all the information needed to use ERBITUX safely and effectively. See full prescribing information for ERBITUX.

ERBITUX® (cetuximab)
injection, for intravenous infusion
Initial U.S. Approval: 2004

WARNING: SERIOUS INFUSION REACTIONS and CARDIOPULMONARY ARREST

See full prescribing information.

- Serious infusion reactions, some fatal, occurred in 1% to 2% of patients. (5.1)
- Cardiopulmonary arrest and/or sudden death occurred in 2% of patients with squamous cell carcinoma of the head and neck treated with Erbitux in combination with radiation therapy and in 3% of patients with squamous cell carcinoma of the head and neck treated with Erbitux in combination with platinum-based therapy with 5-fluorouracil (5-FU). Closely monitor serum electrolytes, including serum magnesium, potassium, and calcium, during and after Erbitux administration. (5.2, 5.6)

Limitation of Use: Erbitux is not indicated for treatment of *K-Ras* mutation-positive colorectal cancer.

Limitation of Use: Erbitux (cetuximab) is not indicated for treatment of *K-Ras* mutation-positive colorectal cancer. (5.7, 14.2)

DOSAGE AND ADMINISTRATION

- Premedicate with an H1 antagonist. (2.3)
- Administer 400 mg/m² initial dose as a 120-minute intravenous infusion followed by 250 mg/m² daily infused over 60 minutes. (2.1, 2.2)
- Initiate Erbitux one week prior to initiation of radiation therapy. Complete Erbitux therapy at least 1 week prior to initiation of platinum-based therapy with 5-FU (2.1) and FOLFIRI (2.2).
- Reduce the infusion rate by 50% for NCI CTC Grade 1 or 2 infusion reactions and permanently discontinue for NCI CTC Grade 3 infusion reaction. (2.4)
- Permanently discontinue for serious, persistent acneiform rash. Reduce dose for recurrent, severe rash. (2.4)

RAMS AND STRENGTHS

- 100 mg/50 mL, single-use vial (3)
- 200 mg/100 mL, single-use vial (3)

CONTRAINDICATIONS

None. (4)

WARNINGS AND PRECAUTIONS

- Infusion Reactions: Immediately stop and permanently discontinue Erbitux for serious infusion reactions. Monitor patients following infusion. (5.1)
- Cardiopulmonary Arrest: Closely monitor serum electrolytes during and after Erbitux. (5.2, 5.6)
- Hypomagnesemia: Periodically monitor during and for at least 8 weeks following the completion of Erbitux. Replace electrolytes as necessary. (5.6)

ADVERSE REACTIONS

The most common adverse reactions (incidence ≥25%) are: cutaneous adverse reactions (including rash, pruritus, and nail changes), headache, diarrhea, and infection. (6)

To report SUSPECTED ADVERSE REACTIONS, contact Bristol-Myers Squibb at 1-800-721-5072 or FDA at 1-800-FDA-1088 or www.fda.gov/medwatch

USE IN SPECIFIC POPULATIONS

- **Pregnancy:** Administer Erbitux to a pregnant woman only if the potential benefit justifies the potential risk to the fetus. (8.1)

RECENT MAJOR CHANGES

Indications and Usage 07/2012

Colorectal Cancer (1.2)

Dosage and Administration

Colorectal Cancer (2.2)

Warnings and Precautions

Use of Erbitux in Combination With Radiation and Cisplatin (5.5) 03/2013

K-Ras Testing in Metastatic or Advanced Colorectal Cancer (5.6)

INDICATIONS AND USAGE

Erbitux® is an epidermal growth factor receptor (EGFR) monoclonal antibody. It is indicated for the treatment of:

Head and Neck Cancer

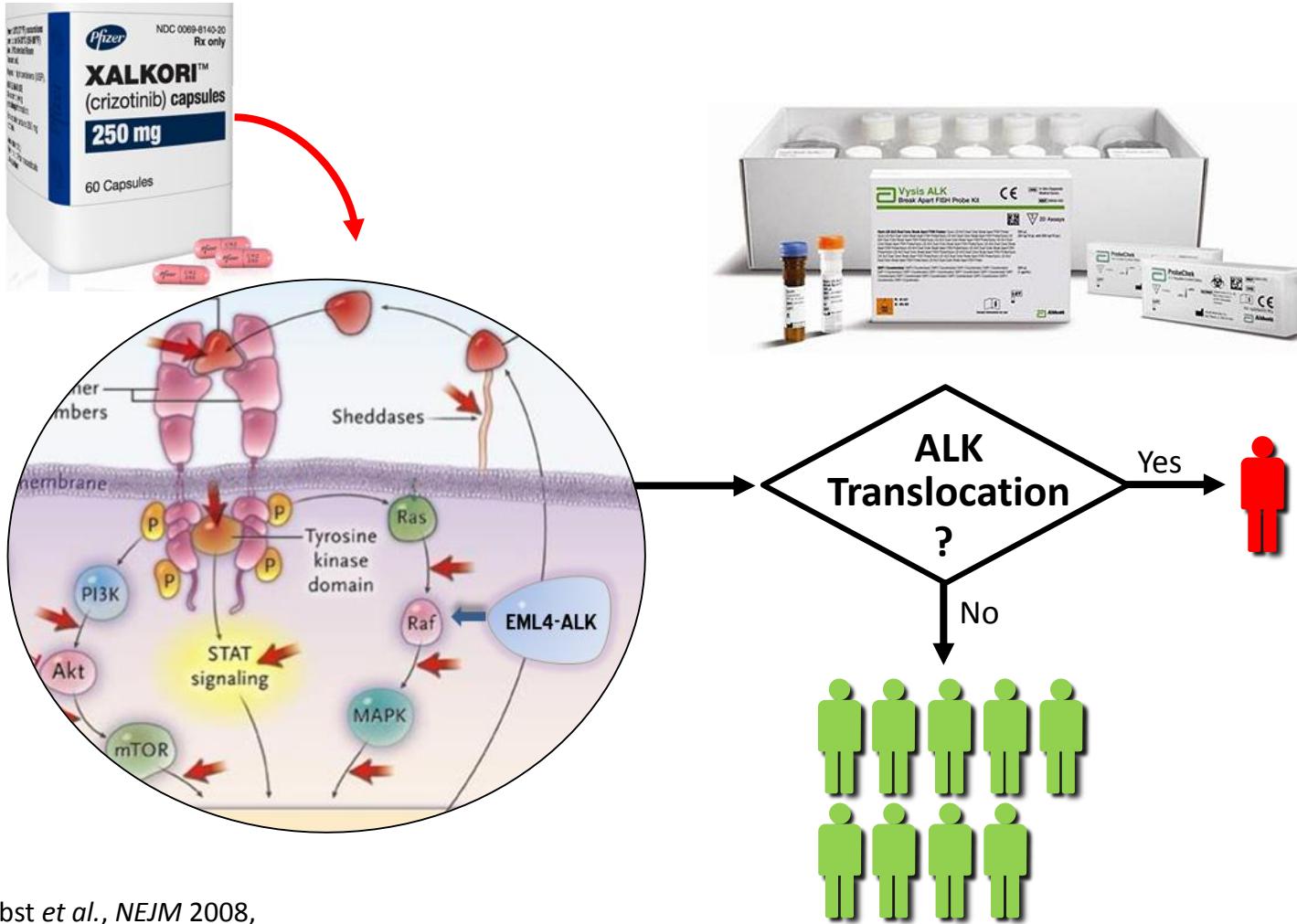
- Locally or regionally advanced squamous cell carcinoma of the head and neck in combination with radiation therapy. (1.1, 14.1)
- Recurrent locoregional disease or metastatic squamous cell carcinoma of the head and neck in combination with platinum-based therapy with 5-FU. (1.1, 14.1)
- Recurrent or metastatic squamous cell carcinoma of the head and neck progressing after platinum-based therapy. (1.1, 14.1)

Colorectal Cancer

K-Ras mutation-negative (wild-type), EGFR-expressing, metastatic colorectal cancer as determined by FDA-approved tests

Xalkori (Crizotinib) Pfizer

In 2011, the FDA 2011, granted accelerated approval to Xalkori with an concurrent approval of ALK Break-Apart FISH Probe Kit (Abbott).



Xalkori (Crizotinib) Pfizer

withhold XALKORI until recovery to less than or equal to Grade 1, then resume XALKORI at 250 mg once daily. Permanently discontinue XALKORI if Grade 3 QTc prolongation recurs [*see Dosage and Administration (2.2) and Clinical Pharmacology (12.4)*].

5.4 ALK Testing

Detection of ALK-positive NSCLC using an FDA-approved test, indicated for this use, is necessary for selection of patients for treatment with XALKORI [*see Clinical Studies (14)*].

Assessment for ALK-positive NSCLC should be performed by laboratories with demonstrated proficiency in the specific technology being utilized. Improper assay performance can lead to unreliable test results.

Refer to an FDA-approved test's package insert for instructions on the identification of patients eligible for treatment with XALKORI.

5.5 Pregnancy

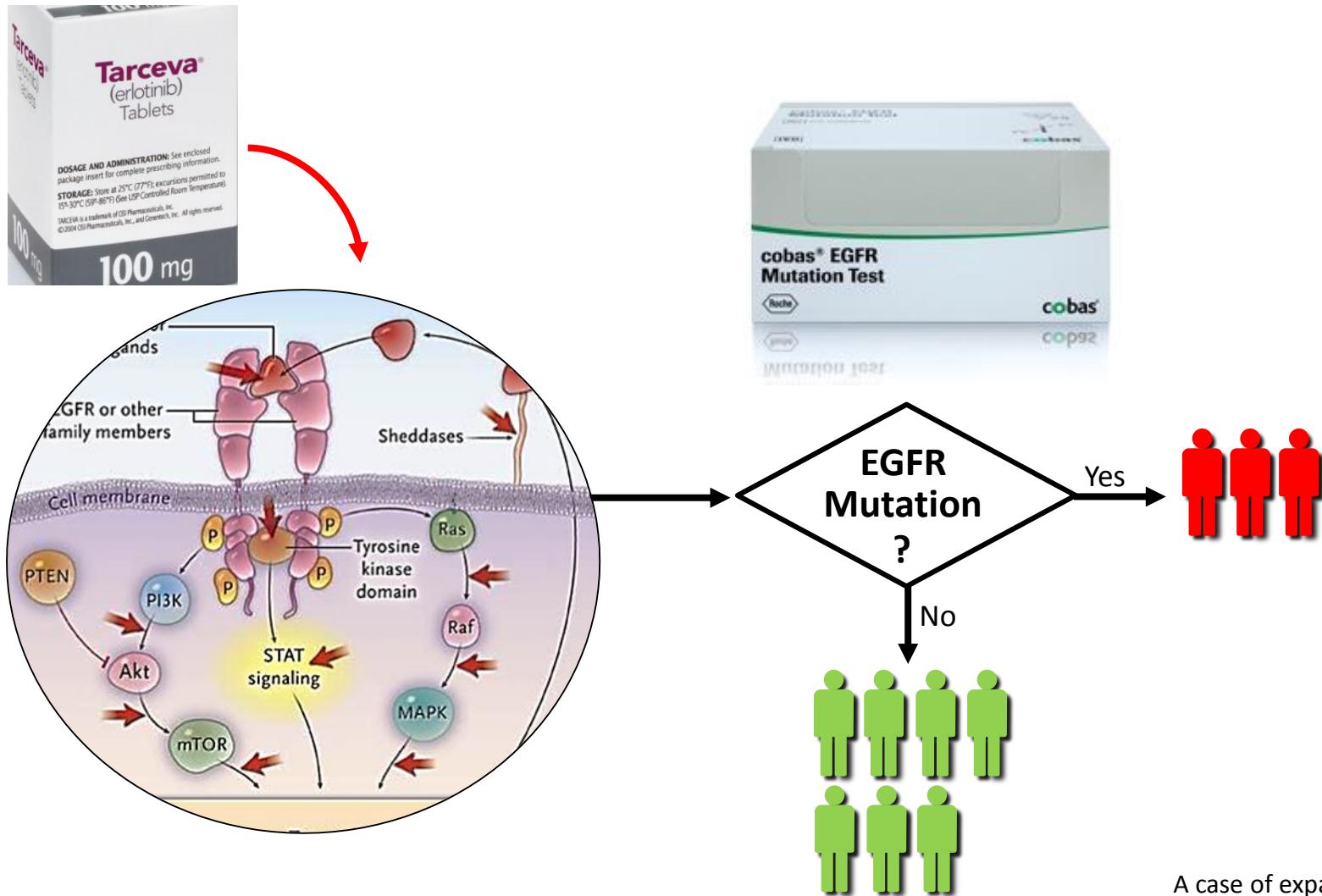
XALKORI can cause fetal harm when administered to a pregnant woman based on its mechanism of action. In nonclinical studies in rats, crizotinib was embryotoxic and fetotoxic at exposures similar to and above those observed in humans at the recommended clinical dose of 250 mg twice daily. There are no adequate and well-controlled studies in pregnant women using XALKORI. If this drug is used during pregnancy, or if the patient becomes pregnant while taking this drug, apprise the patient of the potential hazard to a fetus [*see Use in Specific Populations (8.1)*].

6. ADVERSE REACTIONS

Because clinical trials are conducted under widely varying conditions, adverse reaction rates observed in the clinical trials of a drug cannot be directly compared to rates in the clinical trials of another drug and may not reflect the rates observed in practice.

Tarceva (Erlotinib) Roche

In 2013, the FDA approved Tarceva for the first-line treatment of metastatic NSCLC with EGFR mutation test (Cobas).



A case of expanding market

HIGHLIGHTS OF PRESCRIBING INFORMATION

These highlights do not include all the information needed to use ERBITUX safely and effectively. See full prescribing information for ERBITUX.

ERBITUX® (cetuximab)
injection, for intravenous infusion
Initial U.S. Approval: 2004

WARNING: SERIOUS INFUSION REACTIONS and CARDIOPULMONARY ARREST

See full prescribing information for complete boxed warning.

- Serious infusion reactions, some fatal, occurred in approximately 3% of patients. (5.1)
- Cardiopulmonary arrest and/or sudden death occurred in 2% of patients with squamous cell carcinoma of the head and neck treated with Erbitux and radiation therapy and in 3% of patients with squamous cell carcinoma of the head and neck treated with cetuximab in combination with platinum-based therapy with 5-fluorouracil (5-FU). Closely monitor serum electrolytes, including serum magnesium, potassium, and calcium, during and after Erbitux administration. (5.2, 5.6)

RECENT MAJOR CHANGES

Indications and Usage

Colorectal Cancer (1.2) 07/2012

Dosage and Administration

Colorectal Cancer (2.2) 07/2012

Warnings and Precautions

Use of Erbitux in Combination With Radiation and Cisplatin (5.5) 03/2013
K-Ras Testing in Metastatic or Advanced Colorectal Cancer Patients (5.7) 07/2012

INDICATIONS AND USAGE

Erbitux® is an epidermal growth factor receptor (EGFR) antagonist indicated for treatment of:

Head and Neck Cancer

- Locally or regionally advanced squamous cell carcinoma of the head and neck in combination with radiation therapy. (1.1, 14.1)
- Recurrent locoregional disease or metastatic squamous cell carcinoma of the head and neck in combination with platinum-based therapy with 5-FU. (1.1, 14.1)
- Recurrent or metastatic squamous cell carcinoma of the head and neck progressing after platinum-based therapy. (1.1, 14.1)

Colorectal Cancer

K-Ras mutation-negative (wild-type), EGFR-expressing, metastatic colorectal cancer as determined by FDA-approved tests

- in combination with FOLFIPIR for first-line treatment,
- in combination with irinotecan in patients who are refractory to irinotecan-based chemotherapy,
- as a single agent in patients who have failed oxaliplatin- and irinotecan-based chemotherapy or who are intolerant to irinotecan. (1.2, 5.7, 12.1, 14.2)

Limitation of Use: Erbitux (cetuximab) is not indicated for treatment of K-Ras mutation-positive colorectal cancer. (5.7, 14.2)

DOSAGE AND ADMINISTRATION

- Premedicate with an H₁ antagonist. (2.3)
- Administer 400 mg/m² initial dose as a 120-minute intravenous infusion followed by 250 mg/m² weekly infused over 60 minutes. (2.1, 2.2)
- Initiate Erbitux one week prior to initiation of radiation therapy. Complete Erbitux administration 1 hour prior to platinum-based therapy with 5-FU (2.1) and FOLFIPIR (2.2).
- Reduce the infusion rate by 50% for NCI CTC Grade 1 or 2 infusion reactions and non-serious NCI CTC Grade 3 infusion reaction. (2.4)
- Permanently discontinue for serious infusion reactions. (2.4)
- Withhold infusion for severe, persistent acneiform rash. Reduce dose for recurrent, severe rash. (2.4)

DOSAGE FORMS AND STRENGTHS

- 100 mg/50 mL, single-use vial (3)
- 200 mg/100 mL, single-use vial (3)

CONTRAINDICATIONS

None. (4)

WARNINGS AND PRECAUTIONS

- **Infusion Reactions:** Immediately stop and permanently discontinue Erbitux for serious infusion reactions. Monitor patients following infusion. (5.1)
- **Cardiopulmonary Arrest:** Closely monitor serum electrolytes during and after Erbitux. (5.2, 5.6)
- **Pulmonary Toxicity:** Interrupt therapy for acute onset or worsening of pulmonary symptoms. (5.3)
- **Dermatologic Toxicity:** Limit sun exposure. Monitor for inflammatory or infectious sequelae. (2.4, 5.4)
- **Hypomagnesemia:** Periodically monitor during and for at least 8 weeks following the completion of Erbitux. Replete electrolytes as necessary. (5.6)

ADVERSE REACTIONS

The most common adverse reactions (incidence ≥25%) are: cutaneous adverse reactions (including rash, pruritus, and nail changes), headache, diarrhea, and infection. (6)

To report SUSPECTED ADVERSE REACTIONS, contact Bristol-Myers Squibb at 1-800-721-5072 or FDA at 1-800-FDA-1088 or www.fda.gov/medwatch

USE IN SPECIFIC POPULATIONS

- **Pregnancy:** Administer Erbitux to a pregnant woman only if the potential benefit justifies the potential risk to the fetus. (8.1)
- **Nursing Mothers:** Discontinue nursing during and for 60 days following treatment with Erbitux. (8.3)

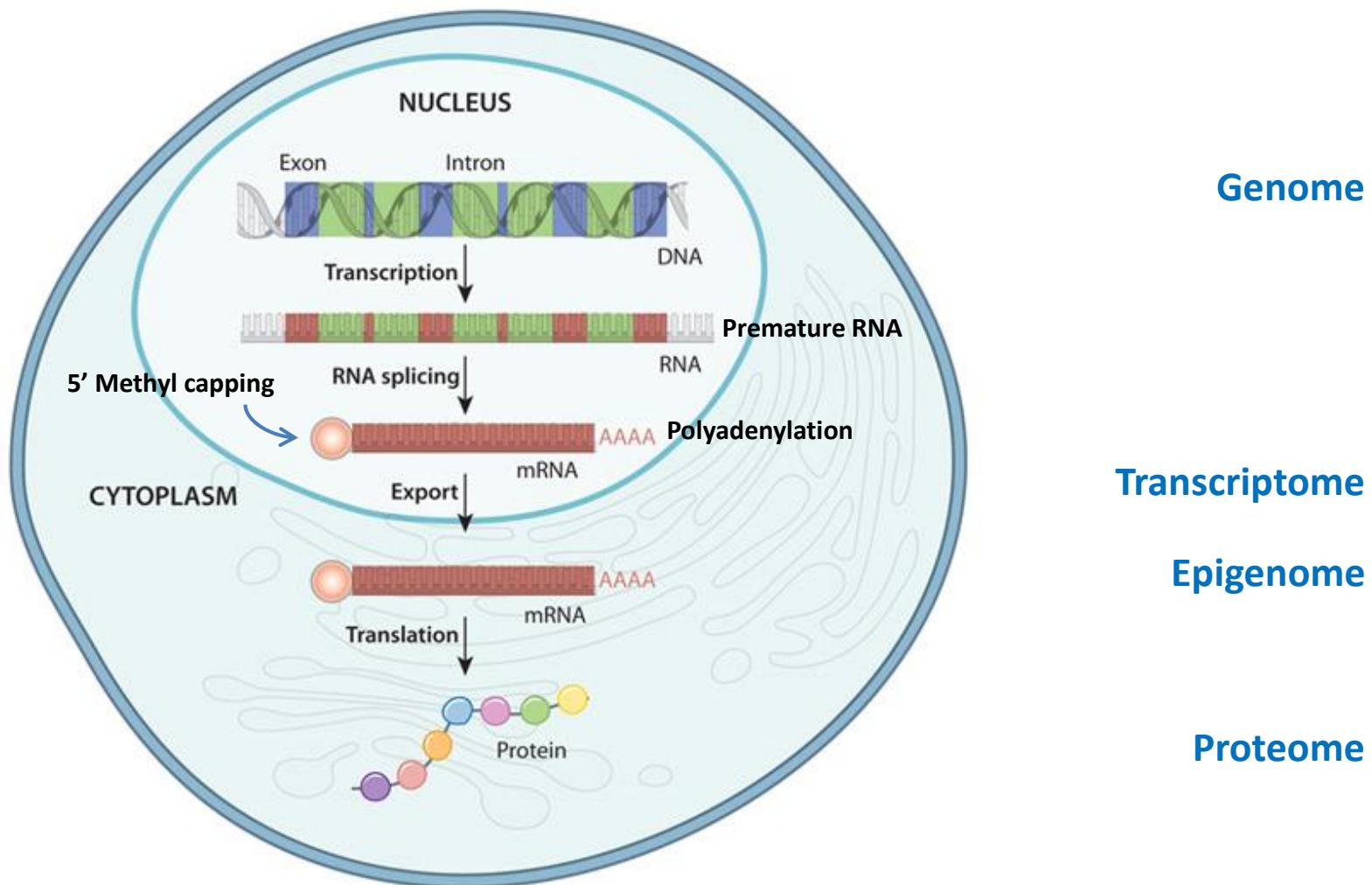
See 17 for PATIENT COUNSELING INFORMATION

BRCA1/2 mutations and Breast Cancer

In 2013...

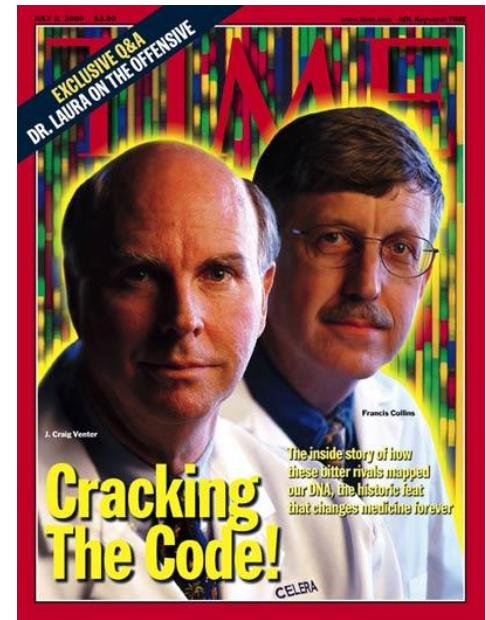


Linking the genetic components to the functional consequences in proteome level and the impact on the clinical decision point-personalized medicine.



Brief history of Human Genome Project

- Project goals: To identify all genes in human DNA, determine the sequences of the 3 billion chemical base pairs that make up human DNA, store this information in databases, improve tools for data analysis
- Estimated budget: \$ 3 billion
- Starting date: 1990
- Projected completion date for draft genome sequence: 2005
- Public venture coordinated by the U.S. Department of Energy and the National Institutes of Health and led by Francis Collins with major partners in UK, France, Japan, Germany, China
- Private venture led by Craig Venter at Celera.
- Draft sequence announced jointly May, 2000 at the White House.
- Draft sequences published in Nature (public) and Science (Celera) in Feb 2001, four years ahead of schedule.



As a result!

<http://www.ncbi.nlm.nih.gov/>

National Center for Biotechnology Information - Windows Internet Explorer

http://www.ncbi.nlm.nih.gov/ epigenomics databases

File Edit View Favorites Tools Help

Favorites Fonds de la Recherche Scie... Google Translate Citrix Access Gateway Matrix Science - Home Suggested Sites Web Slice Gallery

National Center for Biotechnology Information

NCBI Resources How To

All Databases Search

NCBI National Center for Biotechnology Information

Welcome to NCBI

The National Center for Biotechnology Information advances science and health by providing access to biomedical and genomic information.

[About the NCBI](#) | [Mission](#) | [Organization](#) | [Research](#) | [RSS Feeds](#)

Get Started

[Tools](#): Analyze data using NCBI software
[Downloads](#): Get NCBI data or software
[How-To's](#): Learn how to accomplish specific tasks at NCBI
[Submissions](#): Submit data to GenBank or other NCBI databases

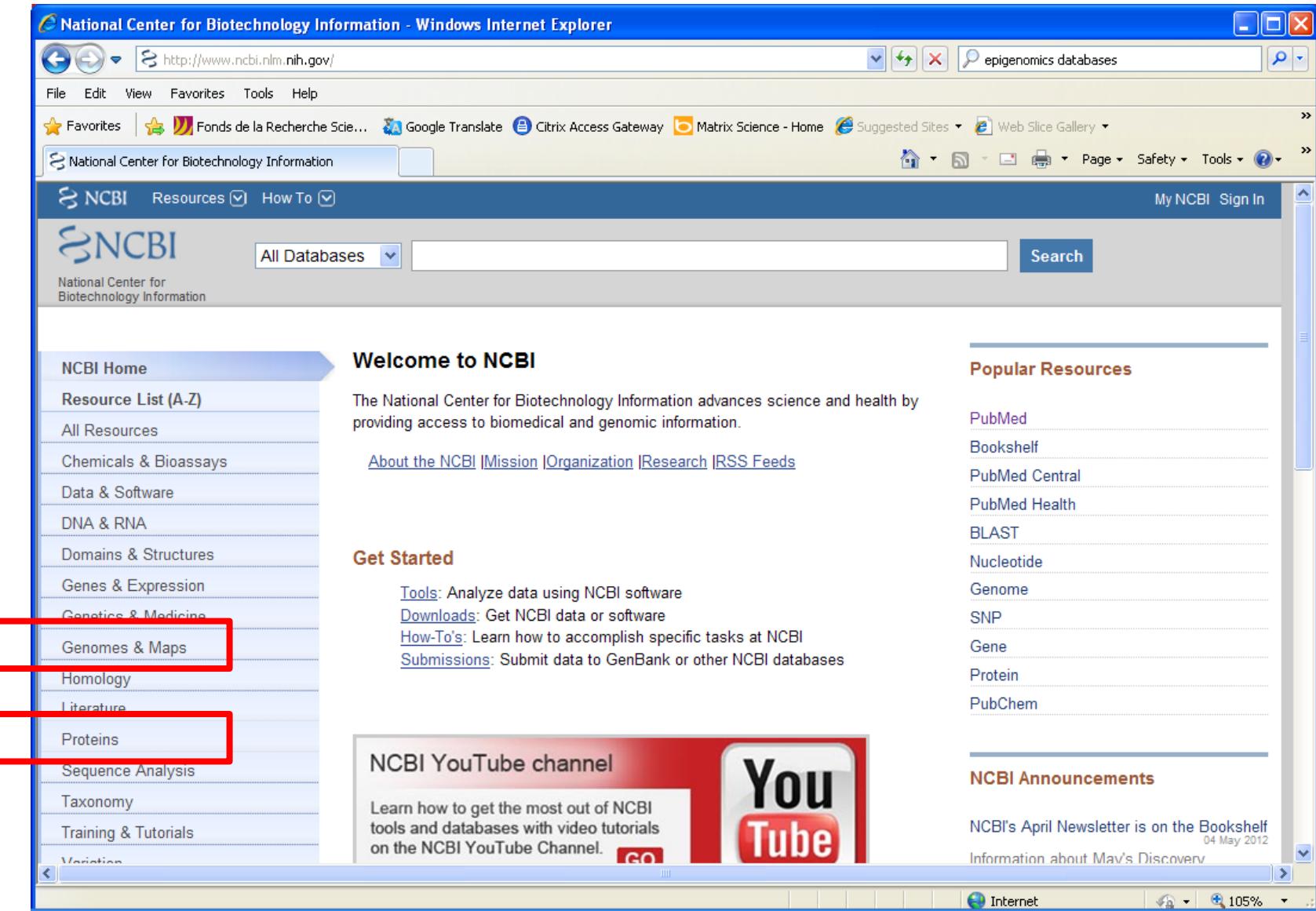
Popular Resources

PubMed
Bookshelf
PubMed Central
PubMed Health
BLAST
Nucleotide
Genome
SNP
Gene
Protein
PubChem

NCBI Announcements

NCBI's April Newsletter is on the Bookshelf
04 May 2012
Information about Mav's Discover

Internet 105% 13



As a result!

The screenshot shows the NCBI Genomes & Maps page. The left sidebar has a blue arrow pointing to the 'Genomes & Maps' link under the 'Resource List (A-Z)' section. The main content area has a red box highlighting the 'BioProject (formerly Genome Project)' section. A blue arrow points from the text 'ex) 1000 Genomes' to this red box. The 'BioProject (formerly Genome Project)' section contains a detailed description of the resource.

ex) 1000 Genomes

Databases

[BioProject \(formerly Genome Project\)](#)
A collection of genomics, functional genomics, and genetics studies and links to their resulting datasets. This resource describes project scope, material, and objectives and provides a mechanism to retrieve datasets that are often difficult to find due to inconsistent annotation, multiple independent submissions, and the varied nature of diverse data types which are often stored in different databases.

[CloneDB \(formerly Clone Registry\)](#)
A database that integrates information about clones and libraries, including sequence data, map positions and distributor information.

[Database of Genome Survey Sequences \(dbGSS\)](#)
A division of GenBank that contains short single-pass reads of genomic DNA. dbGSS can be searched directly through the Nucleotide GSS Database.

[Database of Genomic Structural Variation \(dbVar\)](#)
The dbVar database has been developed to archive information associated with large scale

Quick Links

- [BioProject \(formerly Genome Project\)](#)
- [Database of Genomic Structural Variation \(dbVar\)](#)
- [Genome](#)
- [Nucleotide Database](#)
- [PopSet](#)
- [Sequence Read Archive \(SRA\)](#)
- [Trace Archive](#)
- [UniSTS](#)
- [GenBank: tbl2asn](#)
- [Genome ProtMap](#)
- [Genome Workbench](#)
- [Map Viewer](#)
- [ProSplign](#)
- [Splign](#)

As a result!

1

Genome/Human Genome
-Gene Database/RefSeq
-dbSNP
rs# (reference snp id)

Genome

Contains sequence and map data from the whole genomes of over 1000 organisms. The genomes represent both completely sequenced organisms and those for which sequencing is in progress. All three main domains of life (bacteria, archaea, and eukaryota) are represented, as well as many viruses, phages, viroids, plasmids, and organelles.

2



"The goal of the GRC is to correct misassembled regions, to close remaining gaps, and to provide alternate assemblies of structurally variant positions (loci) in the genome."

Needs of Reference Genome?

3

Influenza Virus

A compilation of data from the NIAID Influenza Genome Sequencing Project and GenBank. It provides tools for flu sequence analysis, annotation and submission to GenBank. This resource also has links to other flu sequence resources, and publications and general information about flu viruses.

Nucleotide Database

A collection of nucleotide sequences from several sources, including GenBank, RefSeq, the Third Party Annotation (TPA) database, and PDB. Searching the Nucleotide Database will yield available results from each of its component databases.

Reference: Baseline

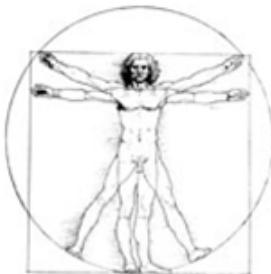
GRC Home Data Help Report an Issue Contact Us Credits Curators Only

Human Mouse Zebrafish Paper Supplemental Data

<http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/>

Genome Assemblies

The GRC has built tools to facilitate the curation of genome assemblies based on the sequence overlaps of long, high quality sequences (Clones and PCR products, not short sequence reads). The GRC currently supports production of assemblies for human, mouse or zebrafish. If your assembly data fits this model and you are interested in using these tools please contact us using the 'Contact Us' page.



Needs of Reference Genome?

- 1) Scientific communication/reference
- 2) Alignment of the sub-genome/genes

Human

The human genome assembly was produced as part of the [Human Genome Project \(HGP\)](#). The previous assembly ([NCBI36](#)) was the last one produced by the HGP and was described in 2004 ([PMID: 15496913](#)) ; this was the starting point for the GRC. The assembly is based largely on assembling overlapping clone sequences.

Human assembly information

Current Major Assembly	GRCh37
Regions with Alternate Loci	3
Assembly N50	46,395,641 bp
Remaining Gaps	357
Patch Release version	p8
Patches Released	Fix: 69; Novel: 71

Reference: Baseline

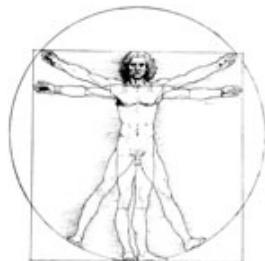
GRChome Data Help Report an Issue Contact Us Credits Curators Only

Human Mouse Zebrafish Paper Supplemental Data

<http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/>

Genome Assemblies

The GRC has built tools to facilitate the curation of genome assemblies based on the sequence overlaps of long, high quality sequences (Clones and PCR products, not short sequence reads). The GRC currently supports production of assemblies for human, mouse or zebrafish. If your assembly data fits this model and you are interested in using these tools please contact us using the 'Contact Us' page. [Subscribe](#) to the grc-announce email list to receive email notification for all GRC assembly updates.



Human

The human genome assembly was produced as part of the [Human Genome Project \(HGP\)](#). The previous assembly (NCBI36) was the last one produced by the HGP and was described in 2004 ([PMID: 15496913](#)); this was the starting point for the GRC. The assembly is based largely on assembling overlapping clone sequences.

Human assembly information

Current Major Assembly	GRCh38
Regions with Alternate Loci	178
Assembly N50	67,794,873 bp
Remaining Gaps	875

Visit 1000 genome webpage
FAQ site
<http://www.1000genomes.org/faq>

Q: "Which reference assembly do you use?"

GenBank

GenBank: The NIH genetic sequence database, an annotated collection of all publicly available DNA sequences*. There are 164 136 731 sequence records in the Genbank as of April 2013.
GenBank is part of the International Nucleotide Sequence Database Collaboration (INSDC).

GenBank at NCBI

DNA DataBank of Japan (DDBJ)

European Molecular Biology Laboratory (EMBL)

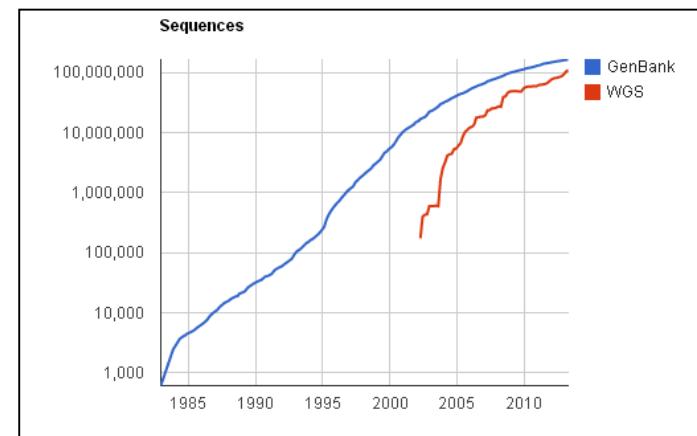
GI number ("gi") : a series of digits that are assigned consecutively to each sequence record processed by NCBI.

VERSION: the accession number of the database record followed by a dot and a version number.

The GI number has been used for many years by NCBI to track sequence histories in GenBank and the other sequence databases it maintains. The VERSION system of identifiers was adopted in February 1999 by the INSDC. The two systems of identifiers run in parallel to each other.

When any change is made to a sequence, it receives a new GI number AND an increase to its version number.

RefSeq : annotated information of DNA, RNA, and protein. constructed from sequence data submitted to the INSDC.



*Guidance and Introduction:
<http://nar.oxfordjournals.org/content/41/D1/D36.full.pdf+html>

Annotation for RefSeq

Accession prefix	Molecule type	Comment
AC_	Genomic	Complete genomic molecule, usually alternate assembly
NC_	Genomic	Complete genomic molecule, usually reference assembly
NG_	Genomic	Incomplete genomic region
NT_	Genomic	Intermediate genomic assemblies of BAC and/or WGS sequence data
NW_	Genomic	Contig or scaffold, primarily WGS ^a
NS_	Genomic	Environmental sequence
NZ ^b	Genomic	Unfinished WGS**
NM_	mRNA	Transcript products; mature messenger RNA (mRNA) transcripts.
NR_	RNA	Non-coding transcripts including structural RNAs, transcribed pseudogenes, and others.
XM ^c	mRNA	Predicted* model
XR ^c	RNA	Predicted model
AP_	Protein	Annotated on AC_ alternate assembly
NP_	Protein	Associated with an NM_ or NC_ accession
YP_	Protein	Protein products; no corresponding transcript record provided. Primarily used for bacterial, viral, and mitochondrial records.
XP_	Protein	Predicted model, associated with an XM_ accession
ZP_	Protein	Predicted model, annotated on NZ_ genomic records

*PREDICTED: Not curated. Automatically provided based on GenBank sequence data; limited or partial support for the transcript or protein. A portion of the transcript or protein may reflect an *ab initio* annotation prediction that was submitted to GenBank.

**WGS: Not curated. The RefSeq record represents a collection of whole genome shotgun (WGS) sequences. This status code is applied to genomic records

Example

Search: “homo sapiens SAA1” in nucleotide section.

RefSeq (13)

Transcript variant 1, mRNA

Check RefSeq accession, GI, version

Search the previous version

CDS for protein

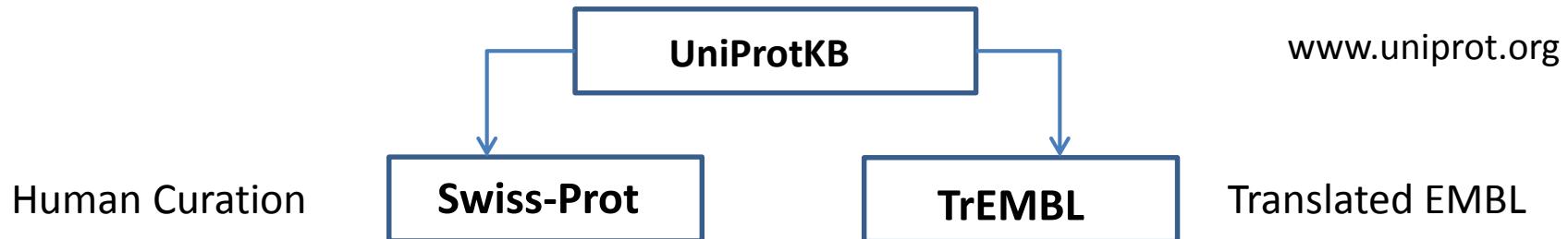
NP_accession

Summary

Entry	mRNA	Protein	mRNA version	GI	GENE	GENE
TR variant 1, mRNA	NM_000331	NP_000322.2	NM_000331.4	GI:295821191	GeneID:6288	HGNC:10513
TR variant 2, mRNA	NM_199161	NP_954630.1	NM_199161.3	GI:295821190	GeneID:6288	HGNC:10513
TR variant 3, mRNA	NM_001178006	NP_001171477.1	NM_001178006.1	GI:295821192	GeneID:6288	HGNC:10513

Protein DB

protein DB with evidence in protein level



The screenshot shows the UniProt homepage. At the top, there is a search bar with "Search in" set to "Protein Knowledgebase (UniProtKB)" and a query "homo sapiens SAA1". Below the search bar, the "WELCOME" section states: "The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information." The "What we provide" section includes a table:

UniProtKB	Protein knowledgebase, consists of two sections: <ul style="list-style-type: none">★ Swiss-Prot, which is manually annotated and reviewed.★ TrEMBL, which is automatically annotated and is not reviewed. Includes complete and reference proteome sets .
UniRef	Sequence clusters, used to speed up sequence similarity searches.
UniParc	Sequence archive, used to keep track of sequences and their identifiers.
Supporting data	Literature citations , taxonomy , keywords , subcellular locations , cross-referenced databases and more.

On the right side of the page, there is a "NEWS" section with a link to "UniProt release 2013_06 - May 29, 2013". Below the news section is a "SITE TOUR" section featuring a screenshot of the UniProt search interface.

Matrix Science - Mascot - Pep X

www.matrixscience.com/cgi/search_form.pl?FORMVER=2&SEARCH=PMF

**MATRIX
SCIENCE**

HOME | WHAT'S NEW | MASCOT | HELP | PRODUCTS | SUPPORT | TRAINING | CONTACT

Search Go

Mascot > Peptide Mass Fingerprint

MASCOT Peptide Mass Fingerprint

Your name	yjk	Email	yeounjin.kim@gmail.com
Search title			
Database(s)	SwissProt NCBInr contaminants cRAP	Enzyme	Trypsin
		Allow up to	1 missed cleavages
Taxonomy	All entries		
Fixed modifications	--- none selected ---	>	NIPCAM (C) Oxidation (HW) Phospho (ST) Phospho (Y) Propionamide (C) Pyridylethyl (C) Pyro-carbamidomethyl (N-term C) Sulfo (STY) TMT2plex (K) TMT2plex (N-term) TMT6plex (K)
	Display all modifications	<	
Variable modifications	Carbamidomethyl (C) Deamidated (NQ) Oxidation (M)	>	
	<		
Protein mass	[] kDa	Peptide tol. ±	1.2 Da
Mass values	<input checked="" type="radio"/> MH ⁺ <input type="radio"/> M _r <input type="radio"/> M-H ⁻	Monoisotopic	<input checked="" type="radio"/> Average <input type="radio"/>
Data file	Choose File No file chosen		
Query	822.43 947.50 1002.56 1227.67 1344.75 1518.67		
Decoy	<input type="checkbox"/>	Report top	AUTO hits
	Start Search ... Reset Form		

Thermo Proteome Discoverer 1.3.0.339 (Viewer)

File Search Report Quantification Tools Window Help



Multi Report from 3 Reports X

Proteins	Peptides	Search Input	Result Filters	Peptide Confidence	Search Summary							
		Sequence	# PSMs	# Proteins	# Protein Groups	Protein Group Accessions	Modifications	MH+ [Da]	B5: Area	C5: Area	A2	
+ 706	<input type="checkbox"/>	SCVGETTESTQCEDEELEHLR	3	8	1	45580688;P10643	C2(Carbamidomethyl); C12...	2509.04312	6.273e7	6.273e7	█	
+ 707	<input type="checkbox"/>	SDDKVTLER	5	7	1	115298678;P01024		1191.58476	1.816e8	6.679e7	█	
+ 708	<input type="checkbox"/>	SDFASNCCSINSPPPLYCDSEI...	4	8	2	139641;P02774;28373620	C7(Carbamidomethyl); C8(...	3332.47190	2.643e8	2.643e8	█	
+ 709	<input type="checkbox"/>	SDIAPVAR	3	7	2	224053;308153640;P01...		828.45690	1.469e8	1.469e8	█	
+ 710	<input type="checkbox"/>	SDLGAVISLLLWGR	3	3	1	1620396		1499.85857	4.382e6	4.382e6	█	
+ 711	<input type="checkbox"/>	SDNCEDTPEAGYFAVAVWK	2	11	3	110590597;119599572;...	C4(Carbamidomethyl)	2071.92558	0.000e0	2.723e7	█	
+ 712	<input type="checkbox"/>	SDNCEDTPEAGYFAVAVWK	2	11	3	110590597;119599572;...	C4(Carbamidomethyl)	2200.02042	0.000e0	7.722e7	█	
+ 713	<input type="checkbox"/>	SDVYTDWK	3	6	2	112877;P02763;229386		1112.52539	2.321e7	2.321e7	█	
+ 714	<input type="checkbox"/>	SDVYTDWKK	3	6	2	112877;P02763;229386		1240.62016	4.207e8	4.207e8	█	
+ 715	<input type="checkbox"/>	SEAEDASLLSFMQGYMK	2	1	1	186972736		1906.85417	0.000e0	9.417e6	█	
+ 716	<input type="checkbox"/>	SEGSSVNLSPPLEQCVPDRG...	3	13	1	4503635;P00734	C15(Carbamidomethyl)	2888.35398	1.731e8	1.731e8	█	
+ 717	<input type="checkbox"/>	SELTQQQLNALFQDK	3	5	2	178779;71773110		1634.83769	1.008e8	1.008e8	█	
+ 718	<input type="checkbox"/>	SELTQQQLNALFQDKLGEVNT ...	3	5	2	178779;71773110		3023.52884	7.477e7	7.477e7	█	
+ 719	<input type="checkbox"/>	SFFSFLGEAFDGar	8	19	5	225985;119588814;134...		1550.72624	3.192e8	3.192e8	█	
		Accession	Description				ΣCoverage	Σ# Proteins	Σ# Unique Peptides	Σ# Peptides	Σ# PSMs	B5: Area
+ 1	<input type="checkbox"/>	225985	amyloid related serum protein SAA				57.69 %	10	1	5	12	0.000e0
+ 2	<input type="checkbox"/>	119588814	serum amyloid A1, isoform CRA_a [Homo sapiens]				52.46 %	11	2	6	16	0.000e0
+ 3	<input type="checkbox"/>	134167;P02735	RecName: Full=Serum amyloid A protein; Short=SAA; Co...				52.46 %	11	2	6	30	7.630e8
+ 4	<input type="checkbox"/>	188497671	serum amyloid A2 isoform a [Homo sapiens]				48.36 %	9	2	6	16	0.000e0
+ 5	<input type="checkbox"/>	315075331	SAA2-SAA2 protein [Homo sapiens]				35.10 %	10	5	8	18	0.000e0
		Sequence	# PSMs	# Proteins	# Protein Groups	Protein Group Accessions	Modifications	MH+ [Da]	B5: Area	C5: Area	A2	
+ 720	<input type="checkbox"/>	SFQTGLFTAAR	2	6	1	192447438;P07225		1198.62102	1.449e7	0.000e0	█	
+ 721	<input type="checkbox"/>	SGAGTEDSACIPWAYYSTVD...	6	10	1	4557485;P00450	C10(Carbamidomethyl)	2505.12114	1.283e8	1.283e8	█	
+ 722	<input type="checkbox"/>	SGAGTEDSACIPWAYYSTVD...	4	10	1	4557485;P00450	C10(Carbamidomethyl); C3...	4061.96555	5.784e7	5.784e7		
+ 723	<input type="checkbox"/>	SGAQATWTELPWPHEK	6	7	1	11321561;P02790		1837.88706	5.655e8	5.655e8	█	
+ 724	<input type="checkbox"/>	SGBTFRPEVHLLPPPSEELAL...	11	1	1	70058	B3(D); B22(N); Z23(Q); C2...	3573.87403	0.000e0	5.328e9	█	

Ready

159/201 Protein Group(s), 9950/14372 Merged Protein(s), 1026/5068 Peptide(s), 3661/14422 PSM(s), 22931/22931 Search Input(s)

BLAST (The Basic Local Alignment Search Tool)

BLAST finds regions of local similarity between sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance of matches. BLAST can be used to infer functional and evolutionary relationships between sequences as well as help identify members of gene families.

The screenshot shows the NCBI BLAST interface. At the top, there's a navigation bar with links for Home, Recent Results, Saved Strategies, and Help. Below that, a breadcrumb trail shows 'NCBI/ BLAST/ blastn suite' and the title 'Standard Nucleotide BLAST'. A red box highlights the search mode dropdown which includes 'blastn', 'blastp', 'blastx', 'tblastn', and 'tblastx'. Below this, there's a section for 'Enter Query Sequence' with a text input field and a 'Clear' button. Further down, there are fields for 'Or, upload file' (with a 'Browse...' button) and 'Job Title' (with a placeholder 'Enter a descriptive title for your BLAST search'). There's also a checkbox for 'Align two or more sequences'. The 'Choose Search Set' section includes a 'Database' dropdown set to 'Human genomic + transcript' (selected with a radio button), and other options like 'Mouse genomic + transcript' and 'Others (nr)'. Below that, there are 'Exclude' and 'Optional' sections with checkboxes for 'Models (XM/XP)' and 'Uncultured/environmental sample sequences'. At the bottom, there's an 'Entrez Query' field with a placeholder 'Enter an Entrez query to limit search'.

blastn: nucleotide-nucleotide

blastp: protein-protein

blastx: nucleotide 6-frame translation-protein
(both strands)

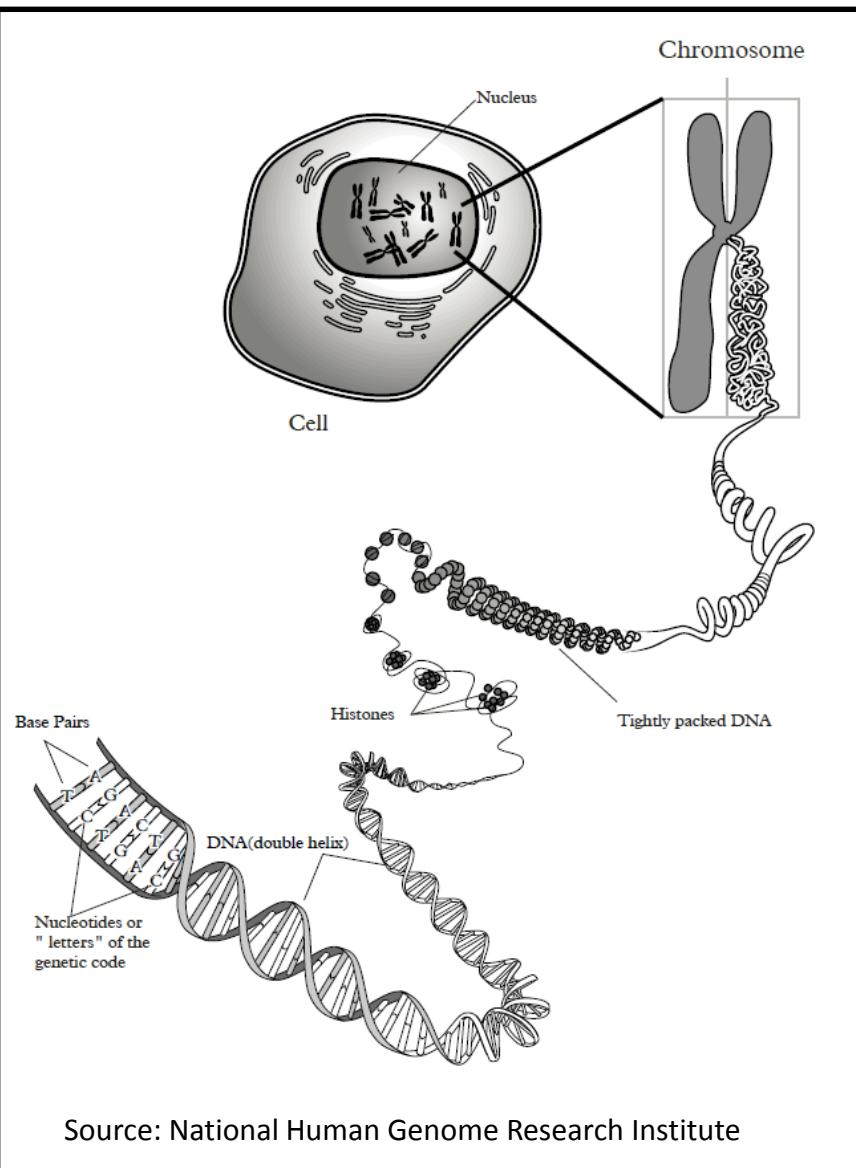
tblastn: Protein-nucleotide 6-frame translation

tblastx: nucleotide 6-frame translation-
nucleotide 6-frame translation

TCT AAC GCT AGT TGA AA
T CTA ACG CTA GTT GAA A
TC TAA CGC TAG TTG AAA

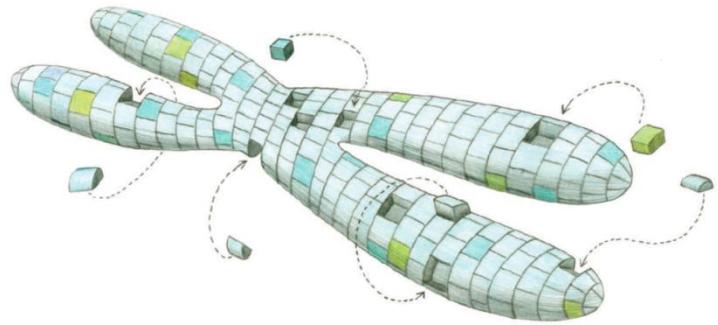
AAAGTTGATCGCAATCT

Human Genome Structure



- 3,000,000,000 base pairs (bp)
 - 24 chromosomes: 1-22, X, Y
 - 21,500 ~ 24,000 genes
 - Coding DNA: only 2% of the genome encodes genes
 - SNPs (single nucleotide polymorphisms)
 - Mutations (change of amino acids)
 - INDELs (insertion or deletions of 1~10Kb)
 - Structural variation (rearrangement or deletion of >10Kb)
 - Inversions (an entire fragment swaps ends)
 - Amplification (multiple repeated copies of a genome segment)
-
- Human genome sequenced in the large scale project serves scientists as a “reference” to which we can compare other human genomes...

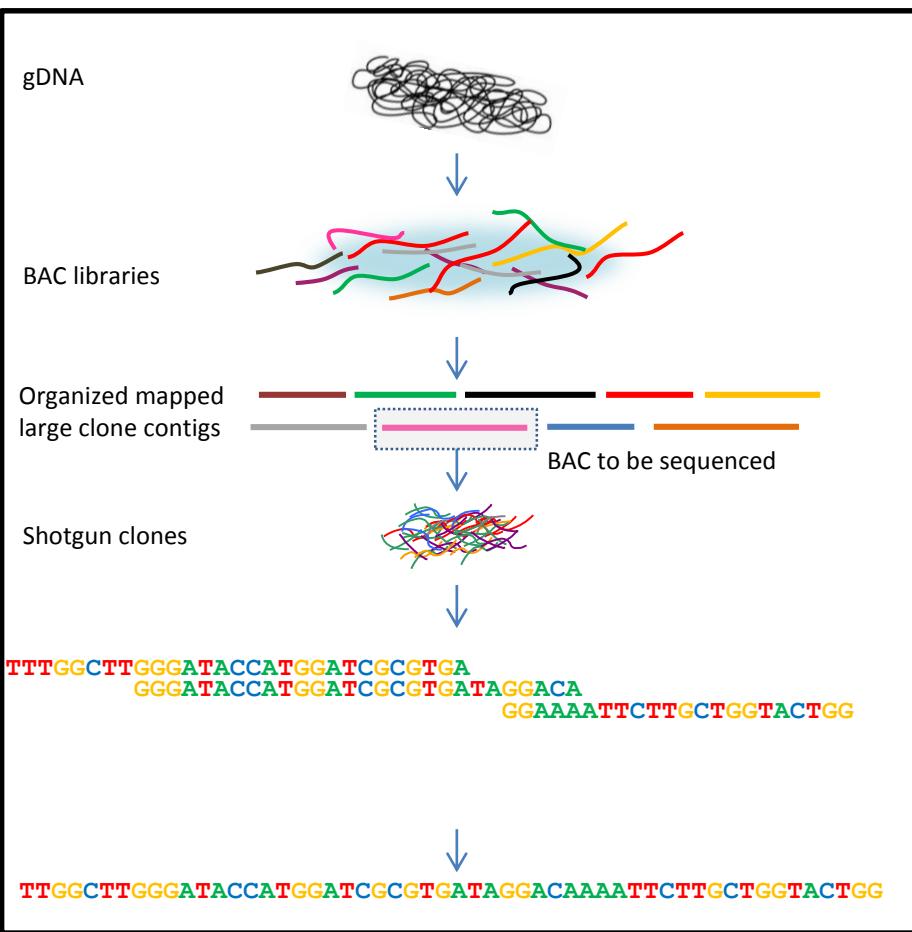
Still on going



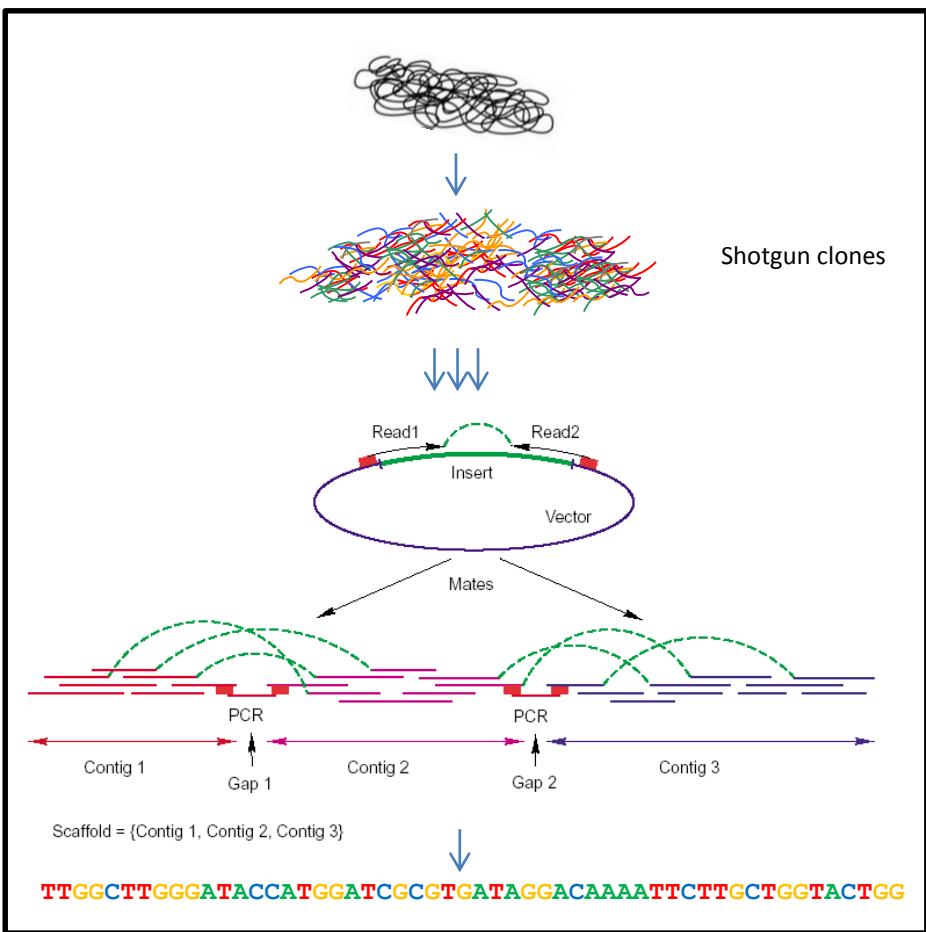
- Though the HGP is finished, analyses of the data will continue for many years.
- Re-sequencing: In general, human genomes are ~99.99% identical. However, we have little knowledge about the level of “normal” human variation beyond SNPs (amplification/deletion and inter-/intra-chromosomal rearrangements, etc.).
- Having a reference genome enables us to investigate variation in additional human genomes. A lot of current sequencing effort is spent on re-sequencing genomes of known species
 - Individual humans (1000 Genomes Project)
 - Experimental organisms – looking for genetic variation, copy number variation
- Challenge: to align millions of sequence reads to a reference genome with some % of mismatches to accurately call SNPs , indels, repeated sequences
- Disease related somatic mutations
 - TCGA consortium
 - Personalized medicine and companion diagnostics
- <http://genome.ucsc.edu/>
 - Exercise: Search for TMPRSS11E in GRCh36 vs GRCh37 **CASE I**

Sequencing strategies used in HGP

Public effort



Celera



Sequencing strategies used in HGP

Public effort

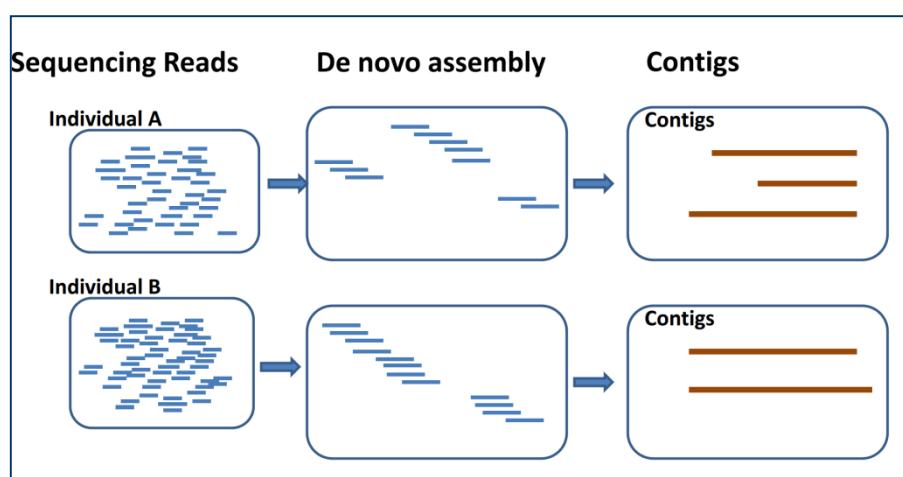


Celera

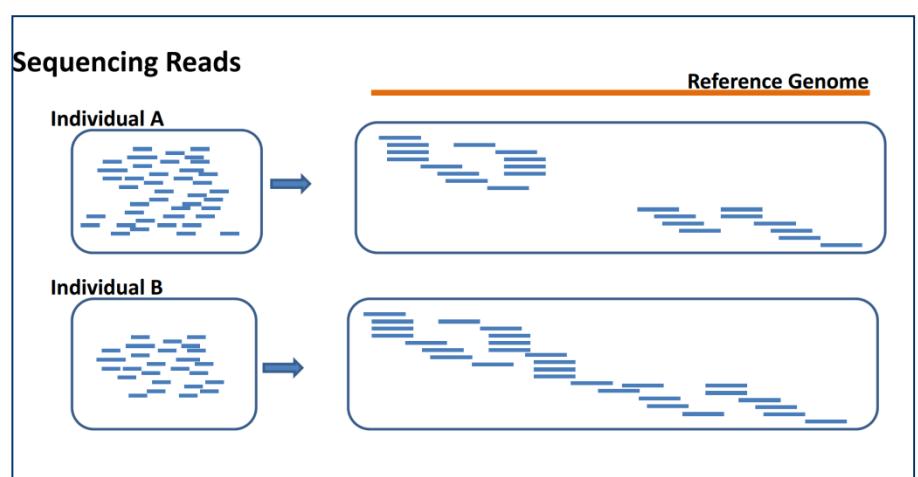


Assembly vs. Alignment

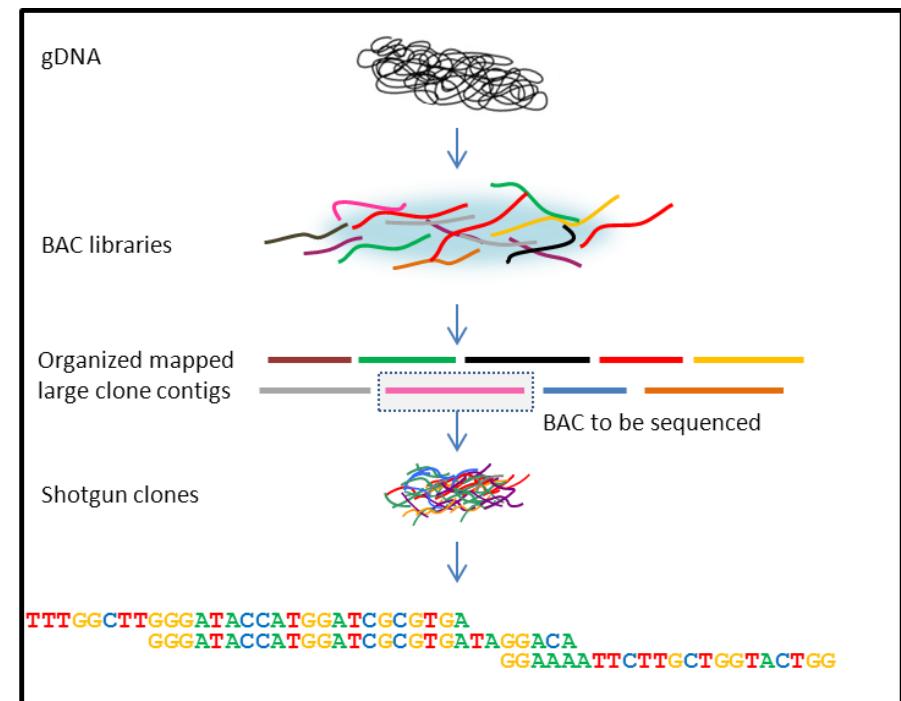
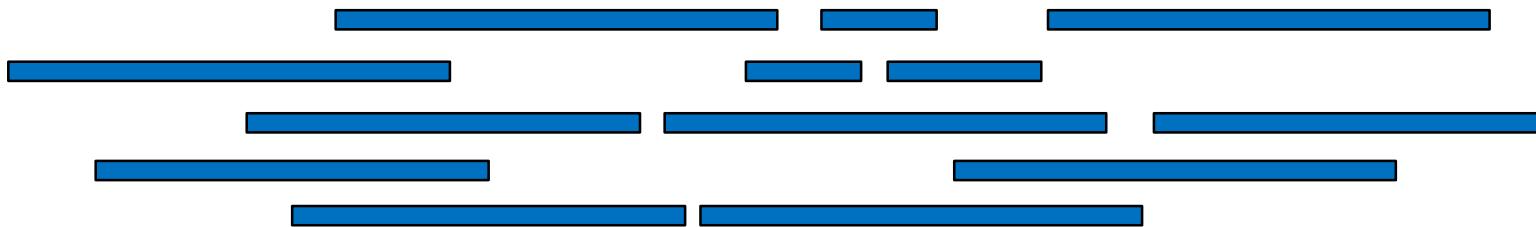
Assembly



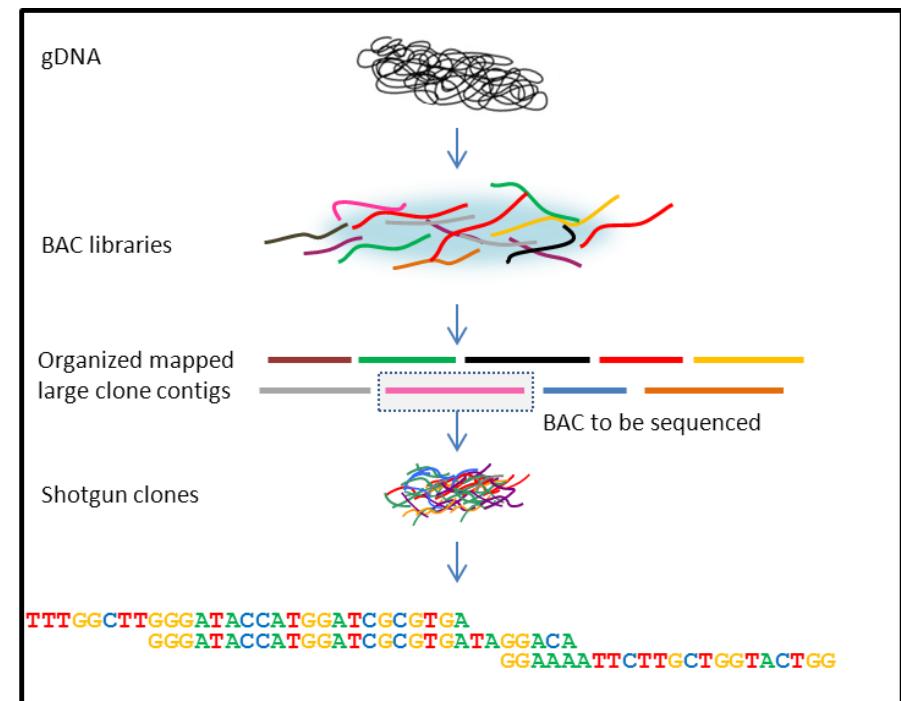
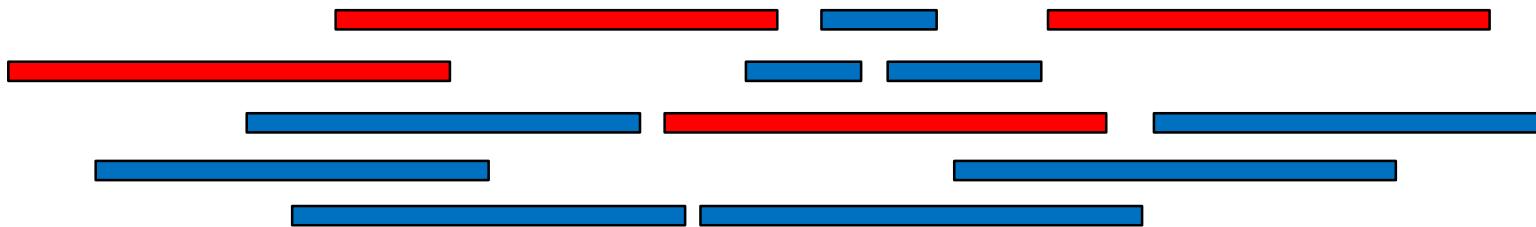
Alignment



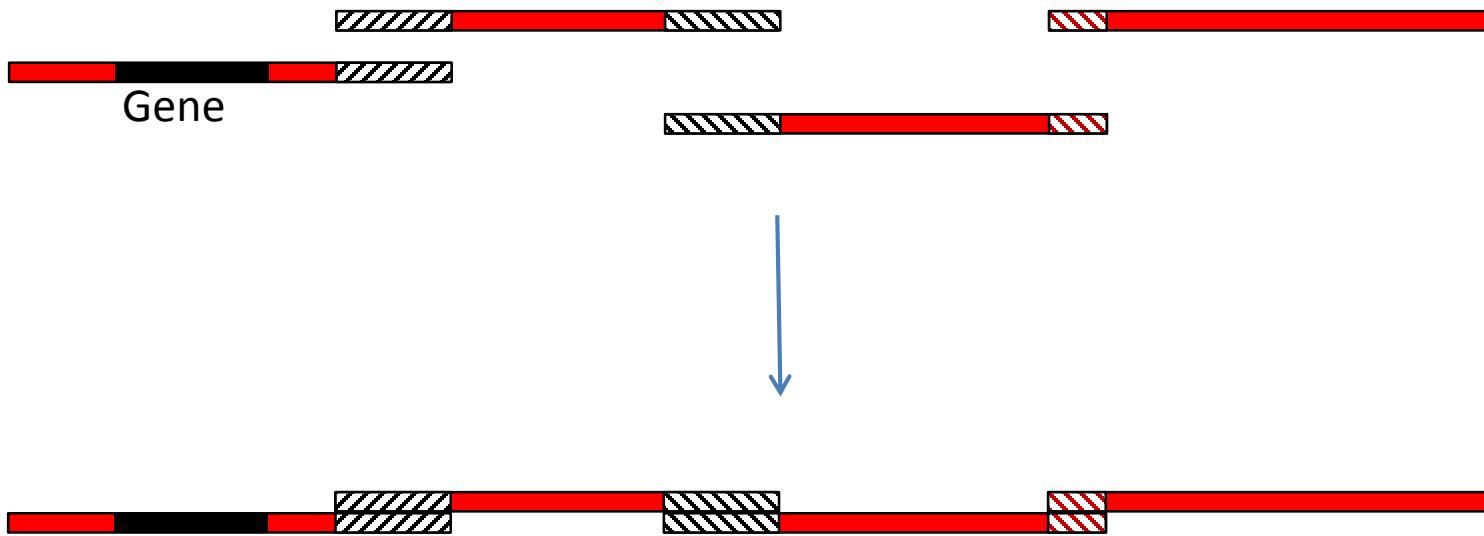
Minimal Tiling Path



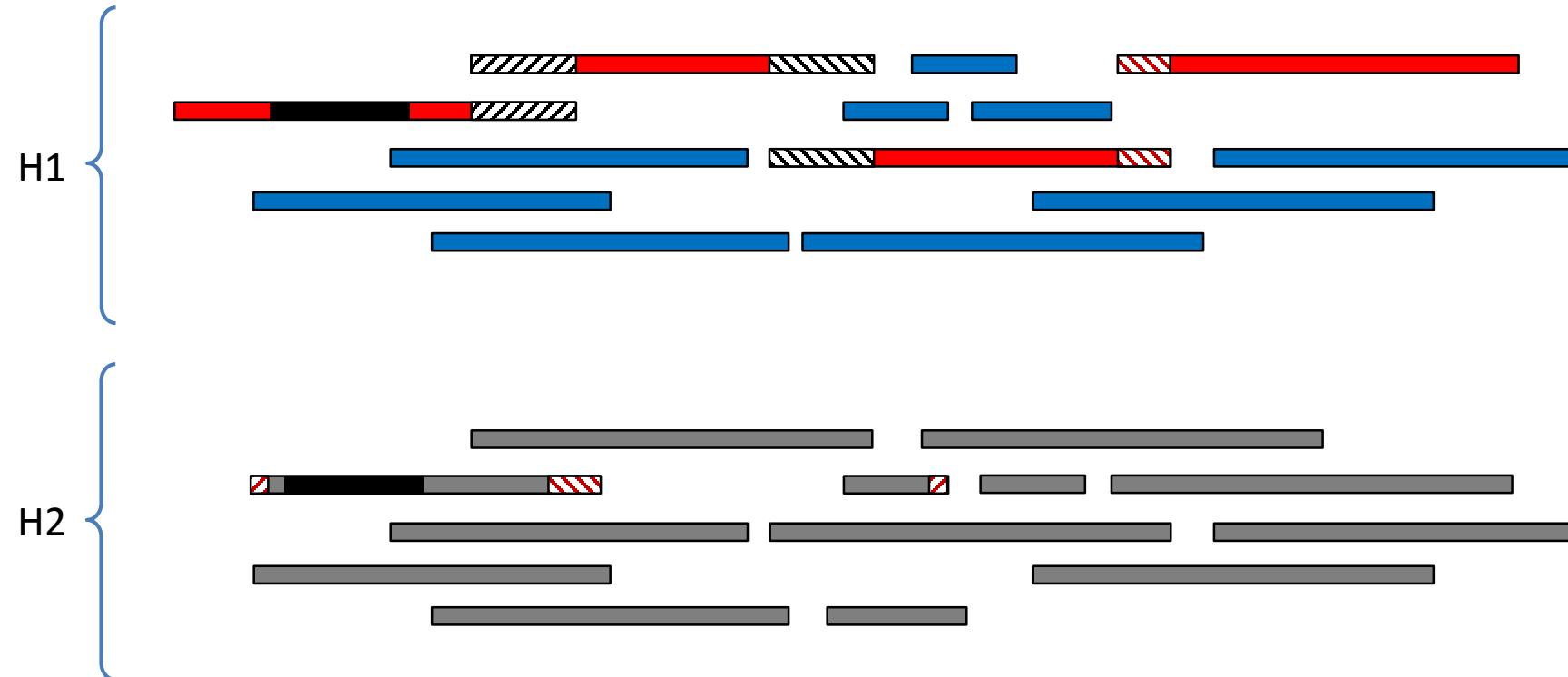
Minimal Tiling Path



Assembly



Misassembly

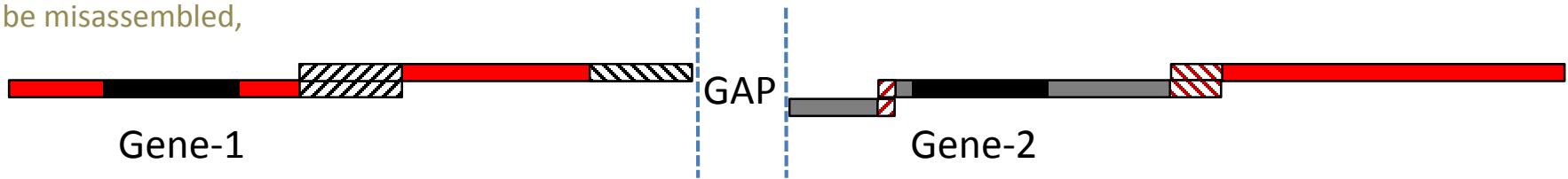


Misassembly

Instead of this,



Can be misassembled,



Misassembled genome in this case results in gene duplication and gap.

Closing the GAP

Home Genomes Blat Tables Gene Sorter PCR Session FAQ Help

Human (*Homo sapiens*) Genome Browser Gateway

The UCSC Genome Browser was created by the [Genome Bioinformatics Group of UC Santa Cruz](#).
Software Copyright (c) The Regents of the University of California. All rights reserved.

clade	genome	assembly	position or search term	gene
Mammal	Human	Feb. 2009 (GRCh37/hg19)	chr7:55086725-55275031	EGFR
Click here to reset the browser user interface settings to their defaults.				
track search	add custom tracks	track hubs	configure tracks and display	clear position

About the Human Feb. 2009 (GRCh37/hg19) assembly ([sequences](#))

The February 2009 human reference sequence (GRCh37) was produced by the [Genome Reference Consortium](#). For more information about this assembly see [GRCh37](#) in the NCBI Assembly database.

Sample position queries

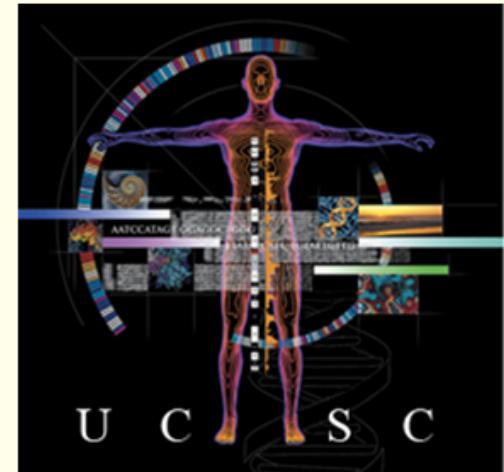
A genome position can be specified by the accession number of a sequenced genomic clone, an mRNA or EST or STS marker, a chromosomal coordinate range, or keywords from the GenBank description of an mRNA. The following list shows examples of valid position queries for the human genome. See the [User's Guide](#) for more information.

Request: **Genome Browser Response:**

chr7 Displays all of chromosome 7

chrUn_g1000212 Displays all of the unplaced contig g1000212

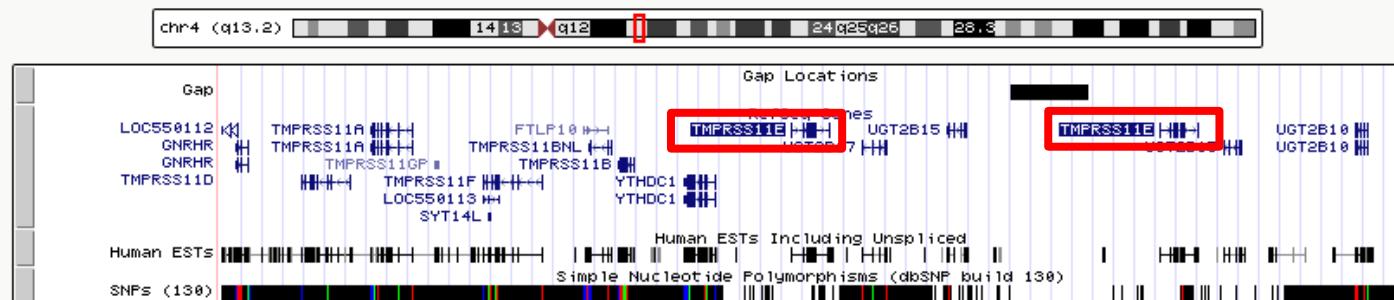
20p13 Displays region for band p13 on chr 20



UCSC Genome Browser on Human Mar. 2006 (NCBI36/hg18) Assembly

move <<< << < > >> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x

chr4:68,268,500-69,773,179 1,504,680 bp. enter position, gene symbol or search terms

"30 x zoom out"


move start < 2.0 > move end < 2.0 >

Click on a feature for details. Click or drag in the base position track to zoom in. Click side bars for track options. Drag side bars or labels up or down to reorder tracks. Drag tracks left or right to new position.

track search default tracks default order hide all add custom tracks track hubs configure reverse resize refresh

collapse all expand all

Use drop-down controls below and press refresh to alter tracks displayed.
Tracks with lots of items will automatically be displayed in more compact modes.

Mapping and Sequencing Tracks

Base Position	Chromosome Band	STS Markers	FISH Clones	Recomb Rate	deCODE Recomb
hide	hide	hide	hide	hide	hide

Map Contigs	Assembly	Gap	Coverage	BAC End Pairs	Fosmid End Pairs
hide	hide	dense	hide	hide	hide

GC Percent	Hg19 Diff	Hi Seq Depth	Wiki Track	BU ORCHID	Mapability
hide	hide	hide	hide	hide	hide

Short Match	Restr Enzymes
hide	hide

Phenotype and Disease Associations

refresh

Genes and Gene Prediction Tracks

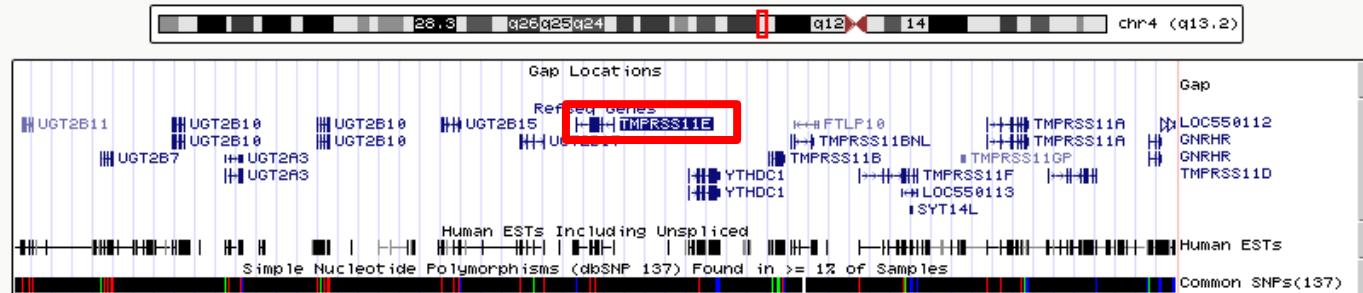
refresh

UCSC Genes	Old UCSC Genes	UCSC Alt Events	Gencode Genes	CCDS	RefSeq Genes
hide	hide	hide	hide	hide	pack

UCSC Genome Browser on Human Feb. 2009 (GRCh37/hg19) Assembly

move <<< << < > >> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x

chr4:68,585,905-70,090,584 1,504,680 bp. enter position, gene symbol or search terms go



track search default tracks default order hide all add custom tracks track hubs configure reverse resize refresh

collapse all

expand all

Use drop-down controls below and press refresh to alter tracks displayed.
Tracks with lots of items will automatically be displayed in more compact modes.



Mapping and Sequencing Tracks

refresh

Base Position

hide

Chromosome Band

hide

STS Markers

hide

18 FISH Clones

hide

Recomb Rate

hide

18 deCODE Recomb

hide

ENCODE Pilot

hide

Map Contigs

hide

Assembly

hide

GRC Map Contigs

hide

Gap

dense

BAC End Pairs

hide

18 Fosmid End Pairs

hide

GC Percent

hide

GRC Patch Release

hide

Hg18 Diff

hide

GRC Incident

hide

Hi Seq Depth

hide

Wiki Track

hide

BU ORCHID

hide

Mapability

hide

Short Match

hide

Restr Enzymes

hide



Phenotype and Disease Associations

refresh



Genes and Gene Prediction Tracks

refresh

UCSC Genes

hide

GENCODE...

hide

Old UCSC Genes

hide

UCSC Alt Events

hide

CCDS

hide

RefSeq Genes

pack

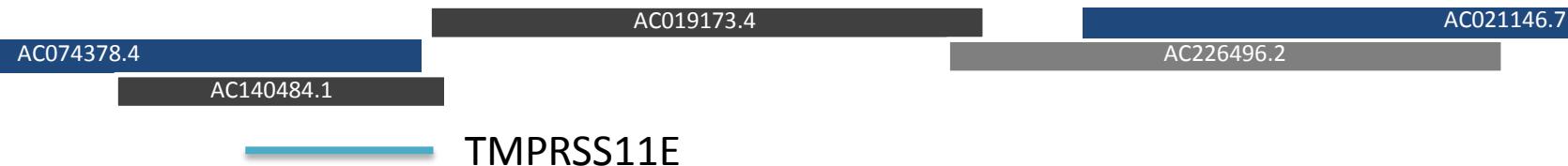
NCBI36 NC_000004.10 (chr4) Tiling Path



GRCh37 NC_000004.11 (chr4) Tiling Path



GRCh37: NT_167250.1 (UGT2B17 alternate locus)

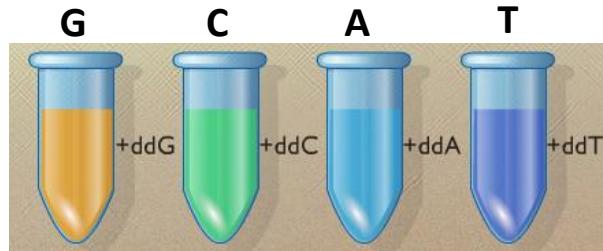


Sequencing

- Basic molecular biology
- Sanger method
- HT sequencing
- Human genome sequencing strategies
 - Directed (Top-down) vs. Shotgun (Bottom-up)
- Next-Generation (Massive Parallel) Sequencing
 - Pyrosequencing (Roche/454)
 - Sequencing by synthesis (Illumina)
 - Sequencing by ligation (SOLiD, AB)

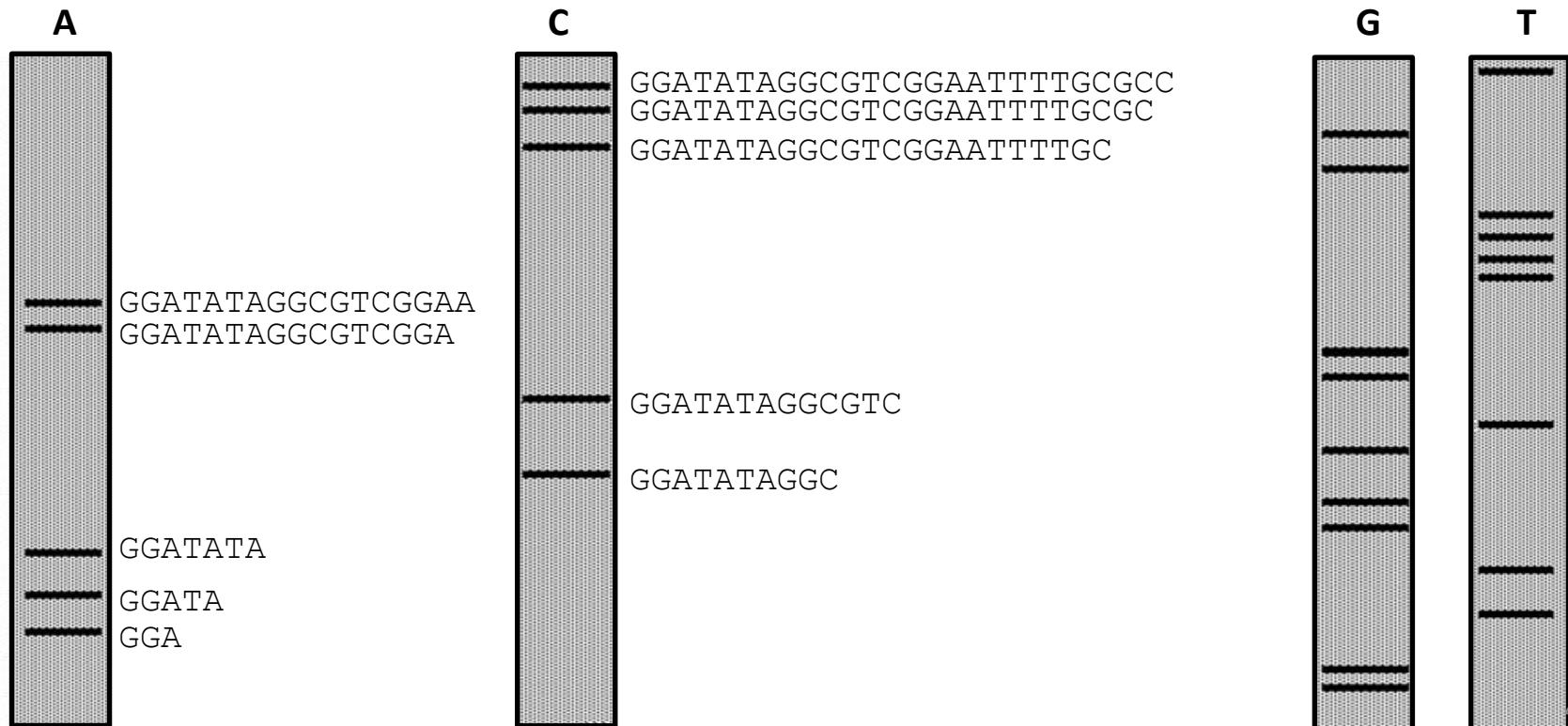
Generation I: Manual Sequencing

Sanger Method: Chain Termination Method



- A: all four dNTP's, **ddATP** and DNA polymerase
- C : all four dNTP's, **ddCTP** and DNA polymerase
- G: all four dNTP's, **ddGTP** and DNA polymerase
- T : all four dNTP's, **ddTTP** and DNA polymerase

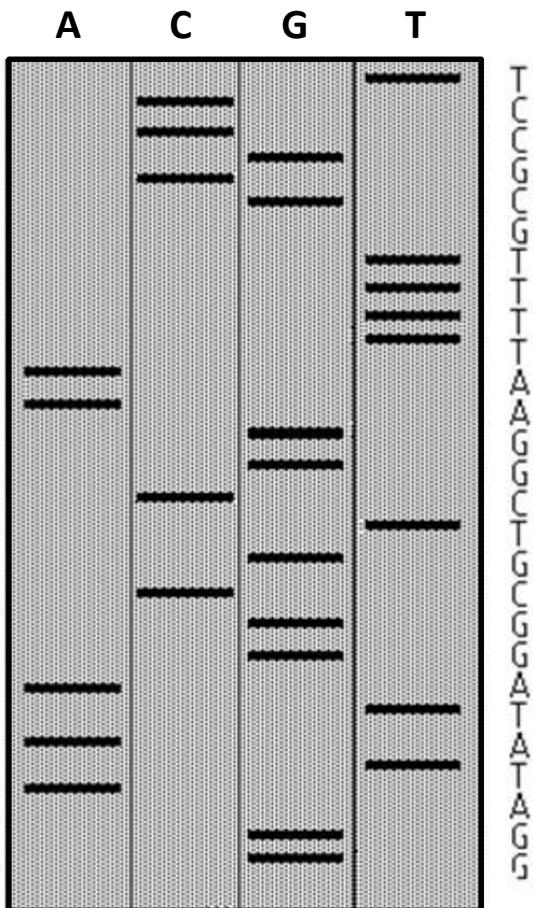
GGATATAAGCGTCGGAATTTCGCCT



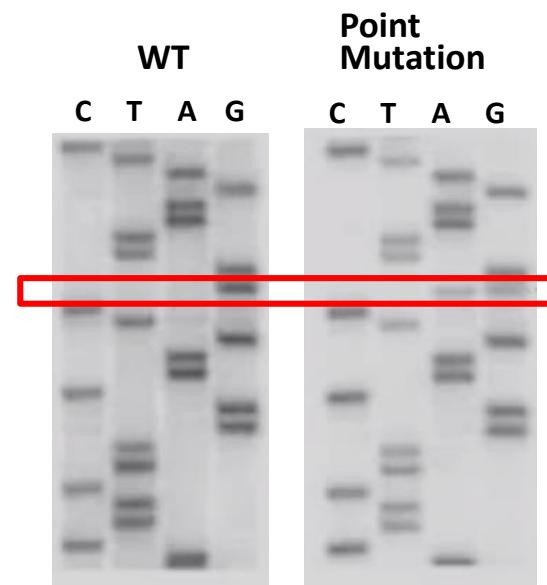
Generation I: Manual Sequencing

Sanger Method: Chain Termination Method

Composition

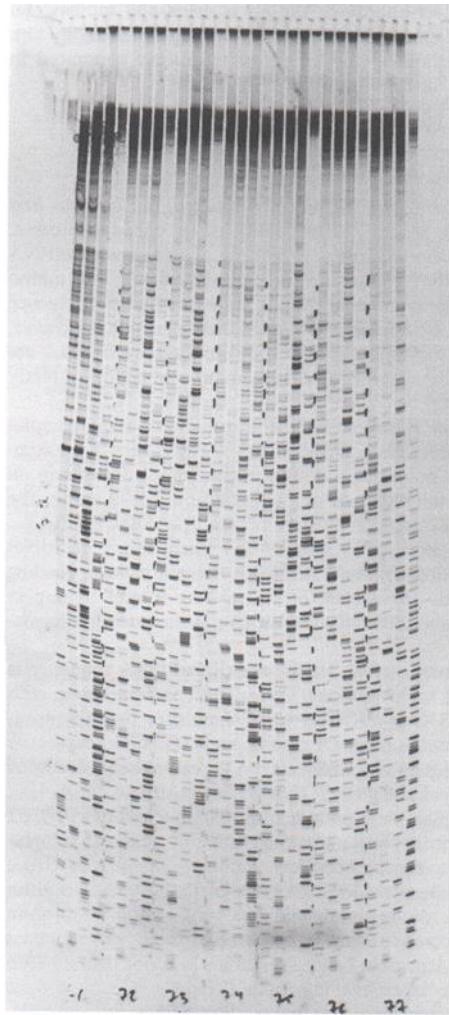


Radioisotope labeled ddNTP
Loading each ddNTP reaction in a different lane



Generation I: Manual Sequencing

Sanger Method: Chain Termination Method



Generation II: Automated Sequencing

<http://www.genomenewsnetwork.org/resources/timeline/>

Introduction | Overview

2004 Rat

2002 Mouse

2001 30,000 Genes

2000 The Human Genome

1999 Fruit Fly

1998 Worm

1996 An Extremophile

1996 Yeast

1995 *Haemophilus*

1991 Venter

1986 Human Genome

1986 Hood

1983 Mullis

1978 Botstein

1977 Gilbert & Sanger

1973 Boyer & Cohen

1972 Berg

1970 Smith

1970 Temin & Baltimore

1969 Beckwith

1967 Weiss & Green

1961 Jacob & Monod

1961 Nirenberg

1960 mRNA

1957 Crick

1956 Kornberg

1953 Crick & Watson

1950 Chargaff

1944 Avery

1943 Delbrück & Luria

1941 Beadle & Tatum

1934 Bernal

1927 Muller

1913 Sturtevant

Genetics and Genomics Timeline

1986

Leroy Hood (1938-) develops the automated sequencer

The techniques for sequencing DNA, developed independently by Walter Gilbert and Frederick Sanger in the late 1970s, represented dramatic advances for research in molecular biology. But they were labor intensive and costly. Automation was always a goal, and just about the time that scientists began to confront the prospect of sequencing the human genome, it became a reality. In 1986, Applied Biosystems Incorporated (ABI) announced the first automatic sequencer, invented by Leroy Hood.

A biologist at the California Institute of Technology and a founder of ABI, Hood improved the existing Sanger' method of enzymatic sequencing, which was becoming the laboratory standard. In this method, DNA to be sequenced is cut apart, and a single strand serves as a template for the synthesis of complementary strands. The nucleotides used to build these strands are randomly mixed with a radioactively labeled and modified nucleotide that terminates the synthesis. Fragments of all different lengths result. The resulting array, sent through a separation gel, reveals the order of the bases. Transferred to film, an "autoradiograph" provides a readable sequence from raw data. This data could be transferred to a computer by a human reader.

In automating the process, Hood modified both the chemistry and the data-gathering processes. In the sequencing reaction itself, he sought to replace the use of radioactive labels, which were unstable, posed a health hazard, and required separate gels for each of the four DNA bases.

- In place of radioisotopes, Hood developed chemistry that used fluorescent dyes of different colors—one for each of the four DNA bases. This system of "color-coding" eliminated the need to run several reactions in overlapping gels.

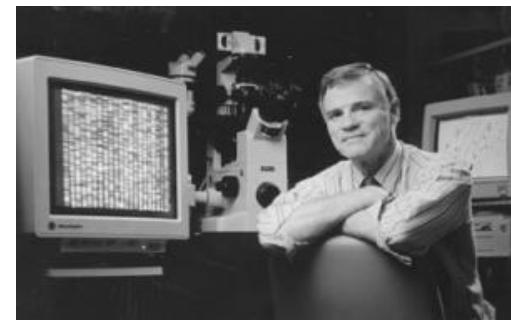
The fluorescent labels were also aspects of the larger system that revolutionized the end stage of the process—the way in which sequence data was gathered. Hood integrated laser and computer technology, eliminating the tedious process of information-gathering by hand.

- As the fragments of DNA percolated through the gel, a laser beam stimulated the fluorescent labels, causing them to glow. The light they emitted was picked up by a lens and photomultiplier, and transmitted as digital information directly into a computer.

 ShareThis

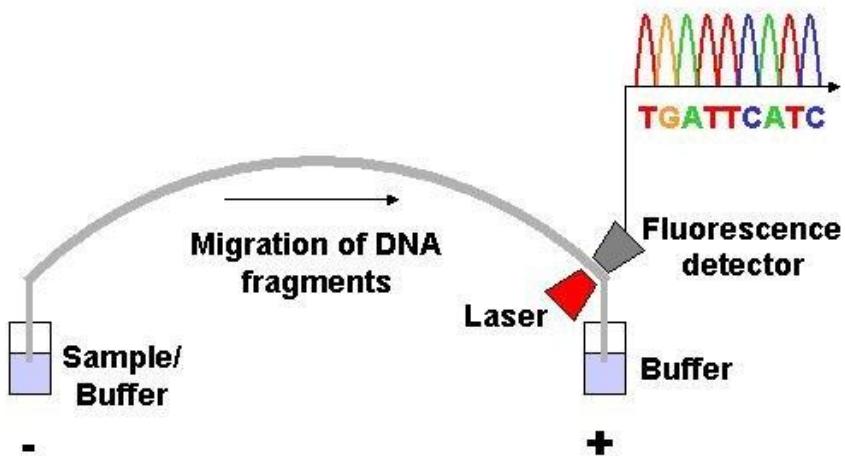


Printer
Friendly

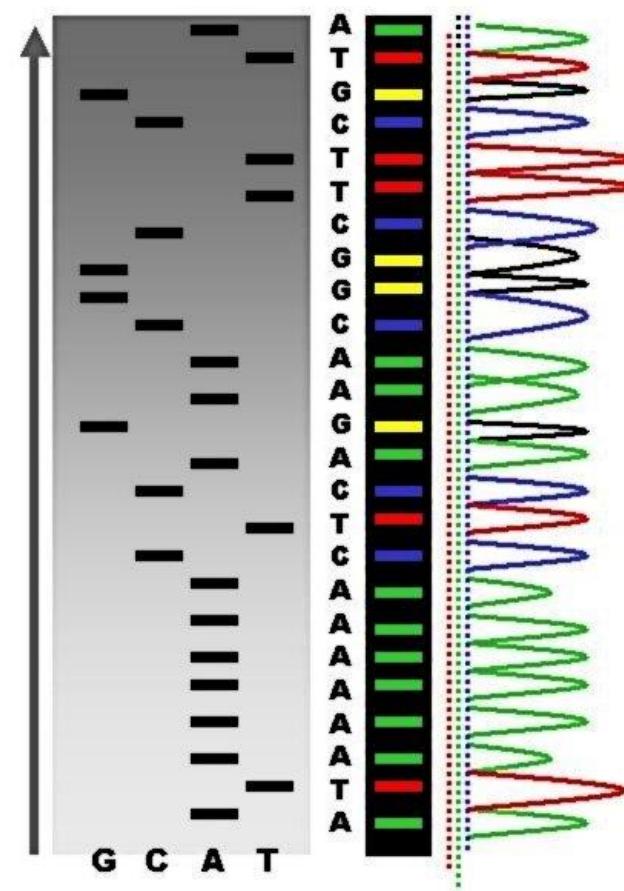


- **Fluorescent dye labeled ddNTP**
- **HT Data handling by integrated laser/computer technology**

Generation II: Automated Sequencing



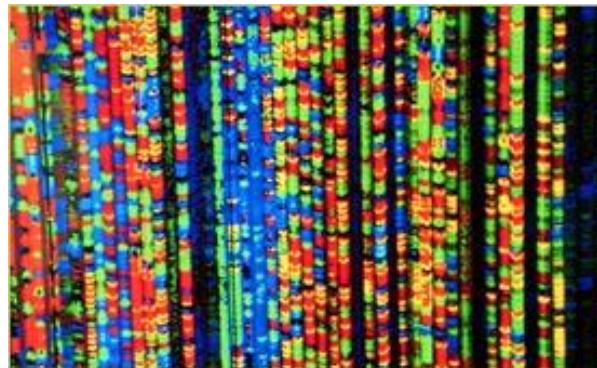
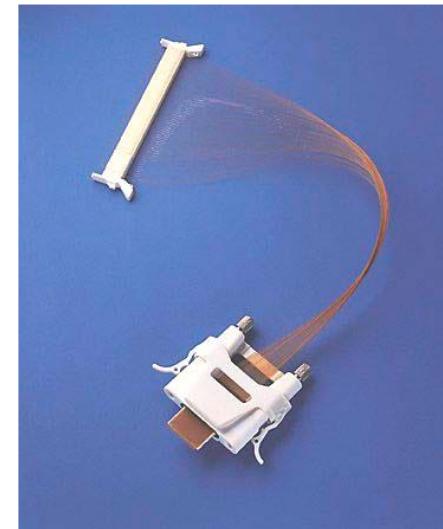
Automated sequencing machines use four colors, so they can read all four bases at once.



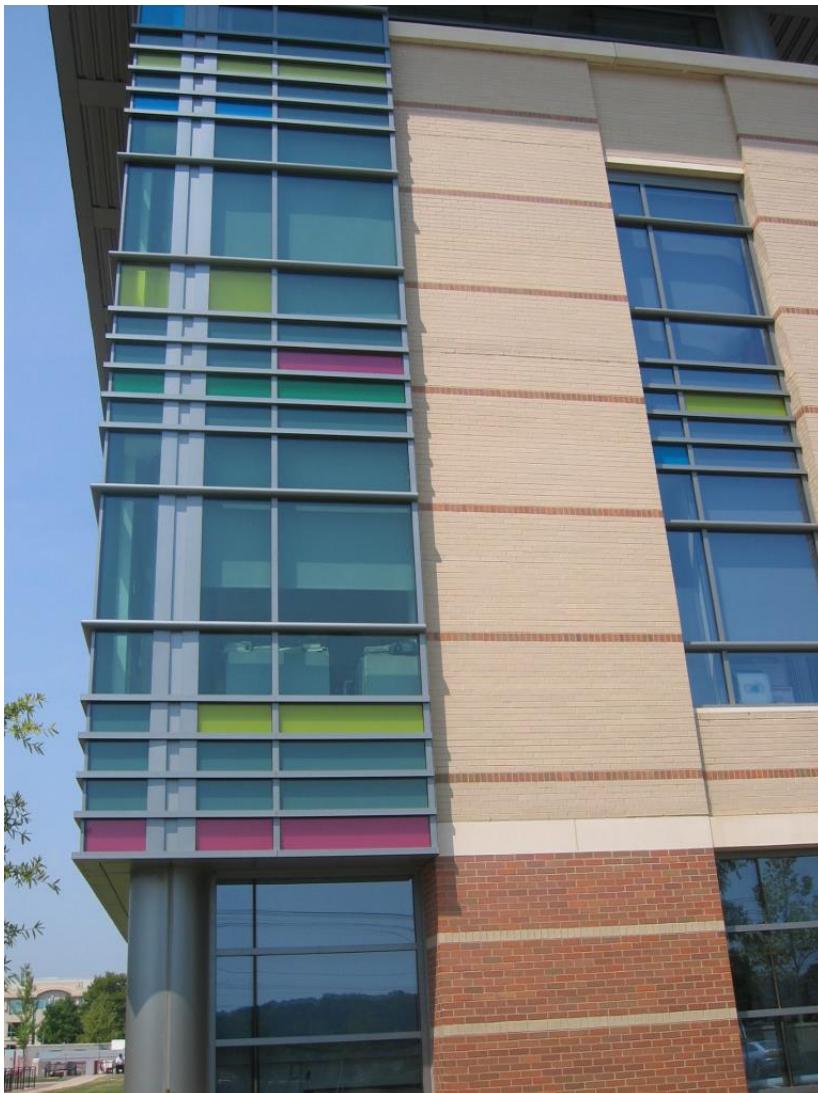
Sequence ladder by radioactive sequencing compared to fluorescent peaks

Generation II: Automated Sequencing

ABI 3730xl



Generation II: Automated Sequencing

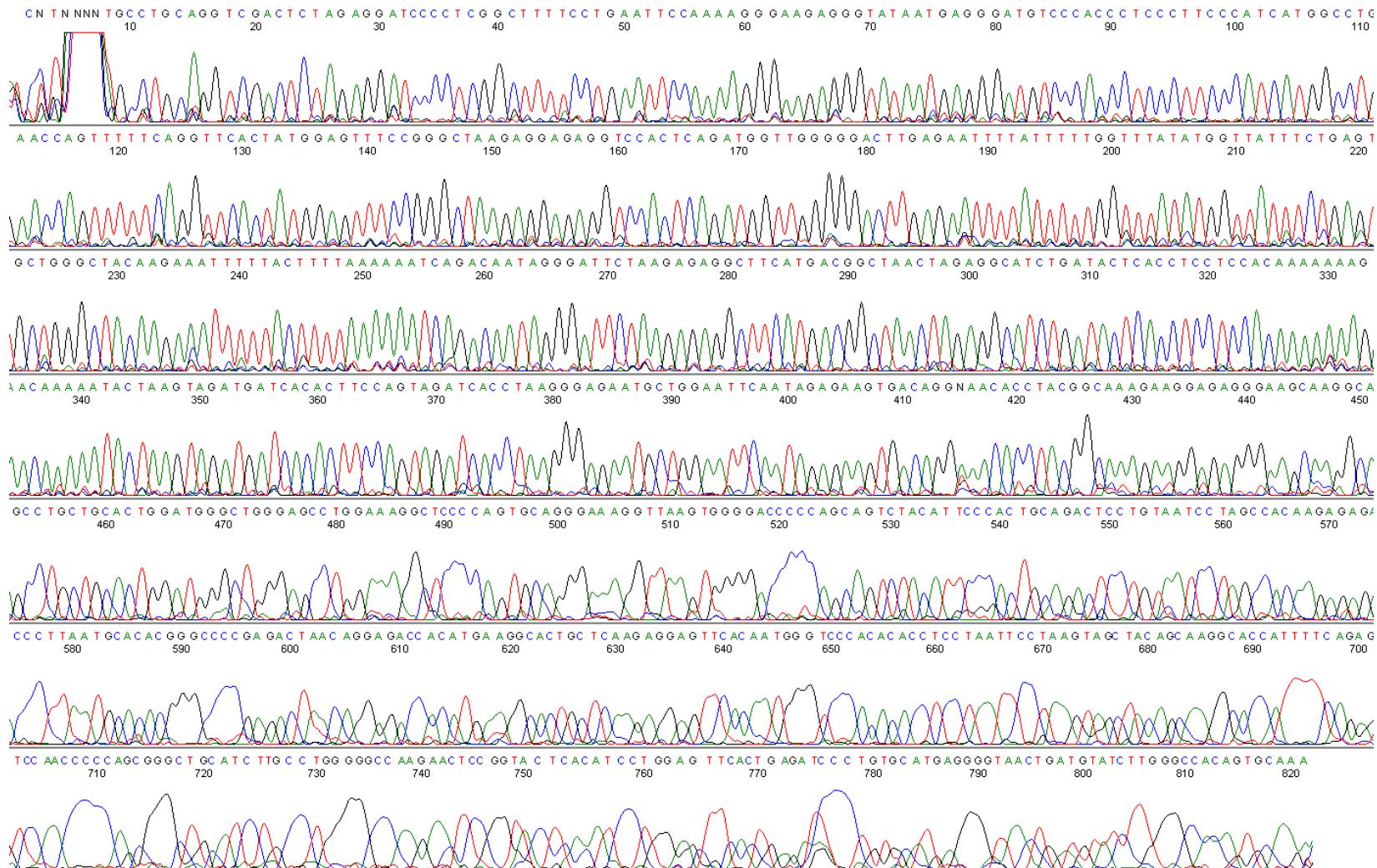


AGCT

JC Venter Institute, Rockville MD

Generation II: Automated Sequencing

Trace File



Exercise

- <http://www.clcbio.com/index.php?id=92>
- Download free trial
- Transfer sequencing data

CASE II. The provided two data sets contain CE based sequencing information from ABI platform. The first set contains different fragments of DNAs from one patient, the second set contains PCR products of the same region of specific gene from multiple patients. Import the data set to your computer, and analyze the data accordingly (1) to identify the region of sequences in reference genome, and (2) to find patients who have variants and the meaning of the variants.

Exercise

<Part1>

1. Assemble the gene fragments derived from a patient.
2. Identify the gene using BLASTn (RefSeq mRNA).
3. Download the sequence file of identified gene from NCBI.
4. Map the contigs to the reference sequence.

<Part2>

5. Align the PCR products from multiple patients.
6. Find the conflict sequences.
7. BLAT at the genome browser to map the variant in the reference and look for annotations
8. Find an associated SNP track
9. SNP DB-> rs number
10. NCBI -> Gene information
11. UniProt/SwissProt -> Protein information, look for variant annotation
12. Disease association

Why trimming?

Trim Sequences

1. Select nucleotide sequences
2. Set trim parameters

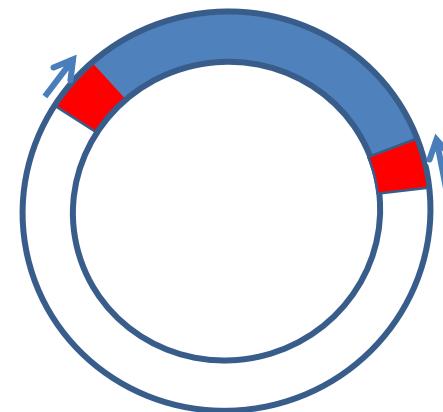
Set trim parameters

Sequence trimming

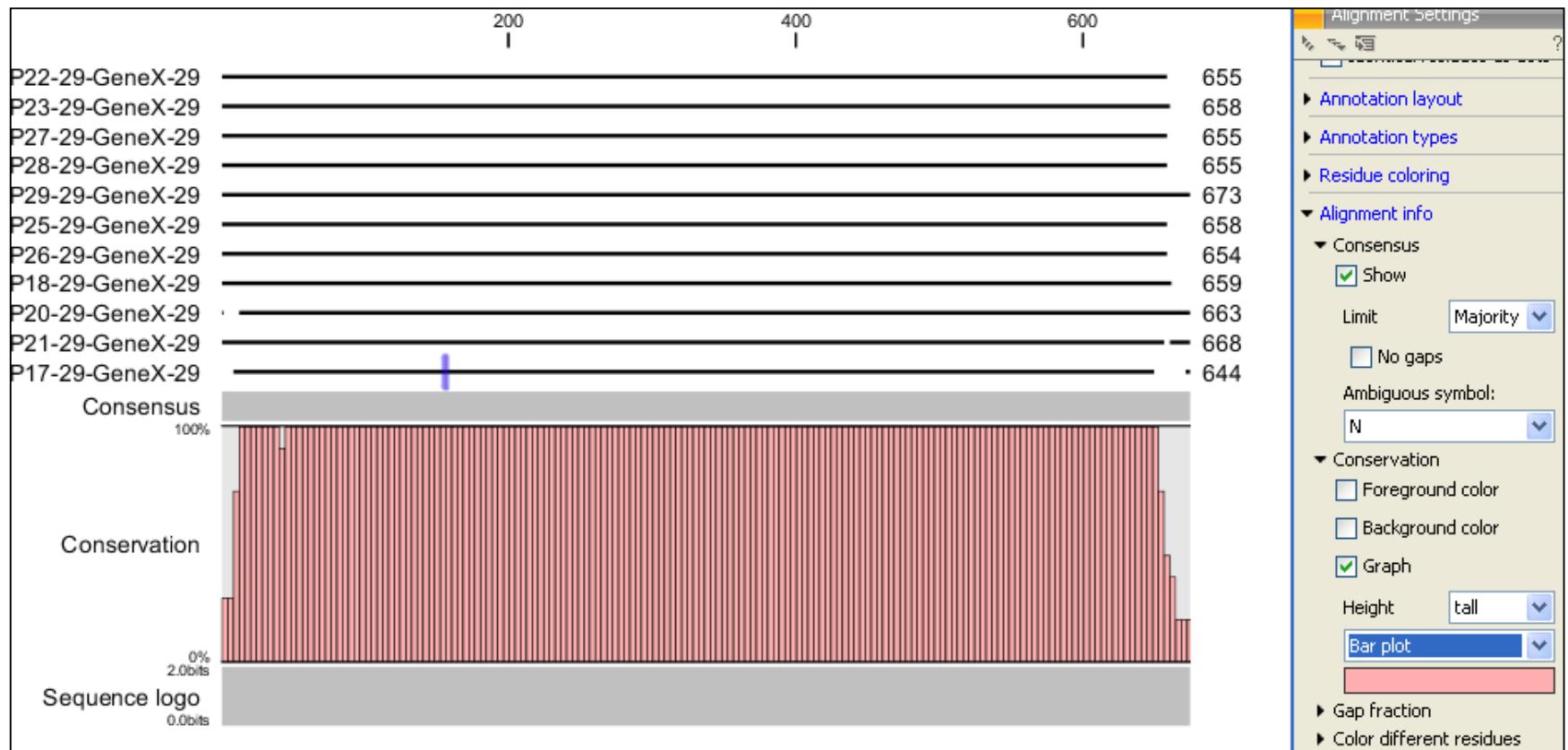
Ignore existing trim information
 Trim using quality scores
Limit: 0,01 **Sequencing error rate**
 Trim using ambiguous nucleotides
Residues: 2

Vector trimming

Trim contamination from vectors in UniVec database **Vector DB**
 Trim contamination from saved sequences (to be chosen in the next step)
Hit limit: moderate



Find SNP



BLAT search and linkage of DBs

BLAT in Genome Browser using GRCh37

Gene Name:

SNP ID:

Synonymous or Non-synonymous:

Disease Association:

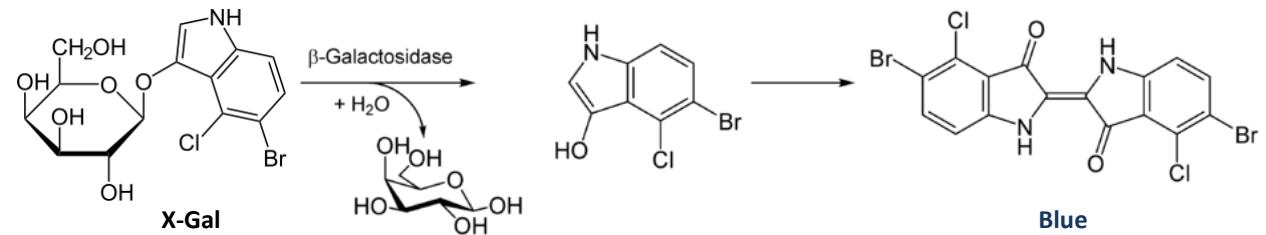
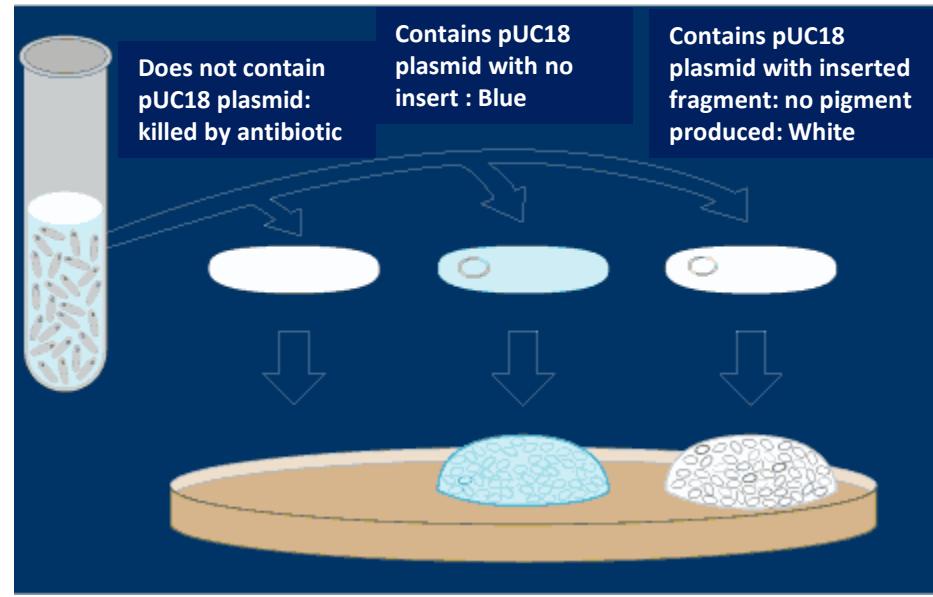
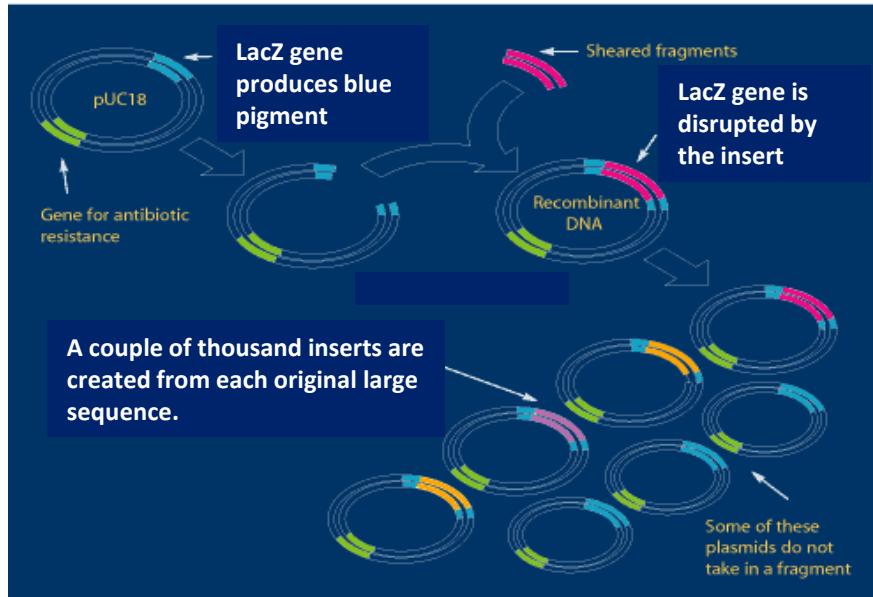
Find RefSeq Protein accession number:

Search the AN in UniProt:

Review and find the SNP ID:

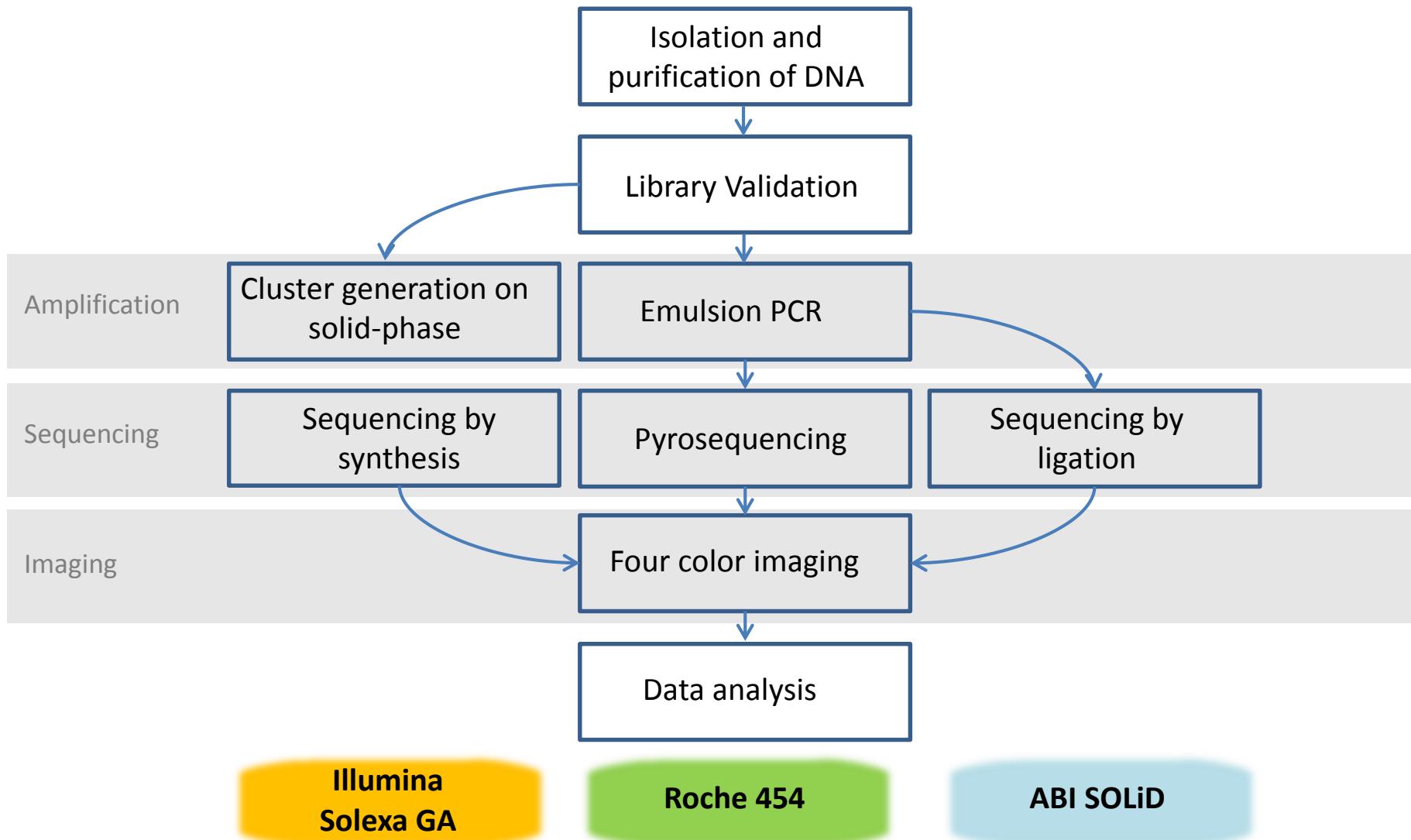
Check reference

Cloning was necessary for prior to G-III



Generation III: Solid Phase Clusters

(Massively Parallel Sequencing)



Generation III: Solid Phase Clusters

(Massively Parallel Sequencing)

Illumina
Solexa GA

Pyrosequencing

Roche 454

Seq by Synthesis

ABI SOLiD

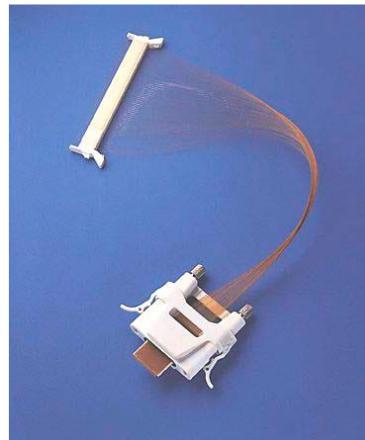
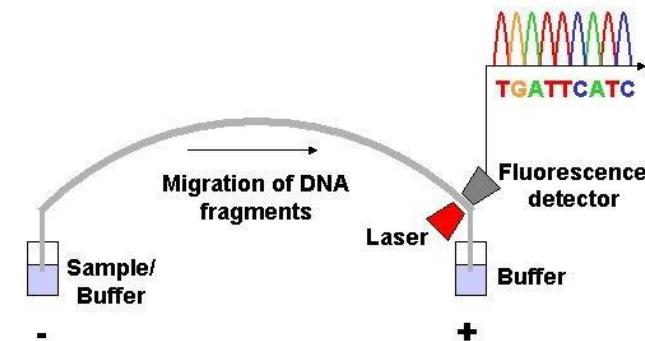
Seq by Ligation

Questions:

1. Major differences of the 2nd Generation Sequencing (eg ABI3730) from classic Sanger Methods.
2. Major differences of NGS from the 2nd GS.

Throughput

Speed is matter



If 500b read per capillary

$$500 \times 100 \text{ b} = 50,000 \text{ b per 3 hr}$$

$$50,000 \text{ b} \times 8 = 400,000 \text{ b per 1 day}$$

If 1000b read per capillary

$$800,000 \text{ b per 1 day}$$

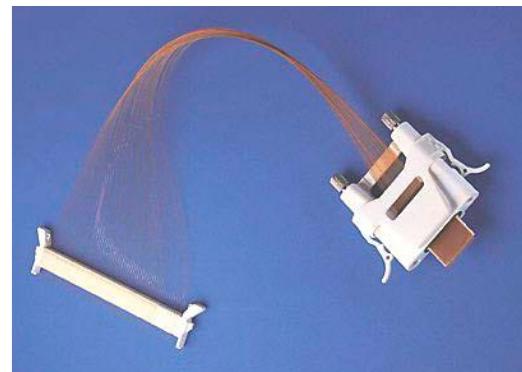
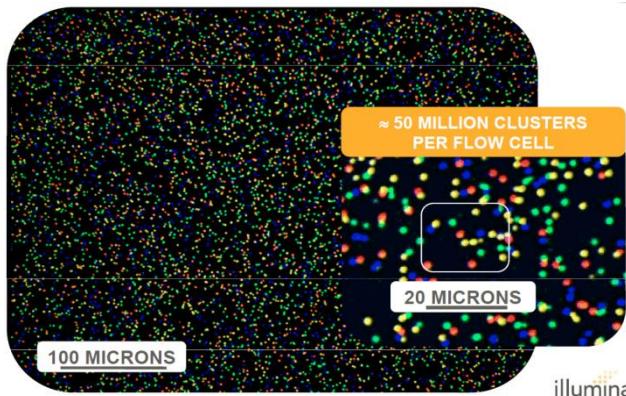
For whole genome

$$3 \times 10^9 \text{ b} \rightarrow 3,750 \text{ d} \rightarrow 10.3 \text{ y}$$

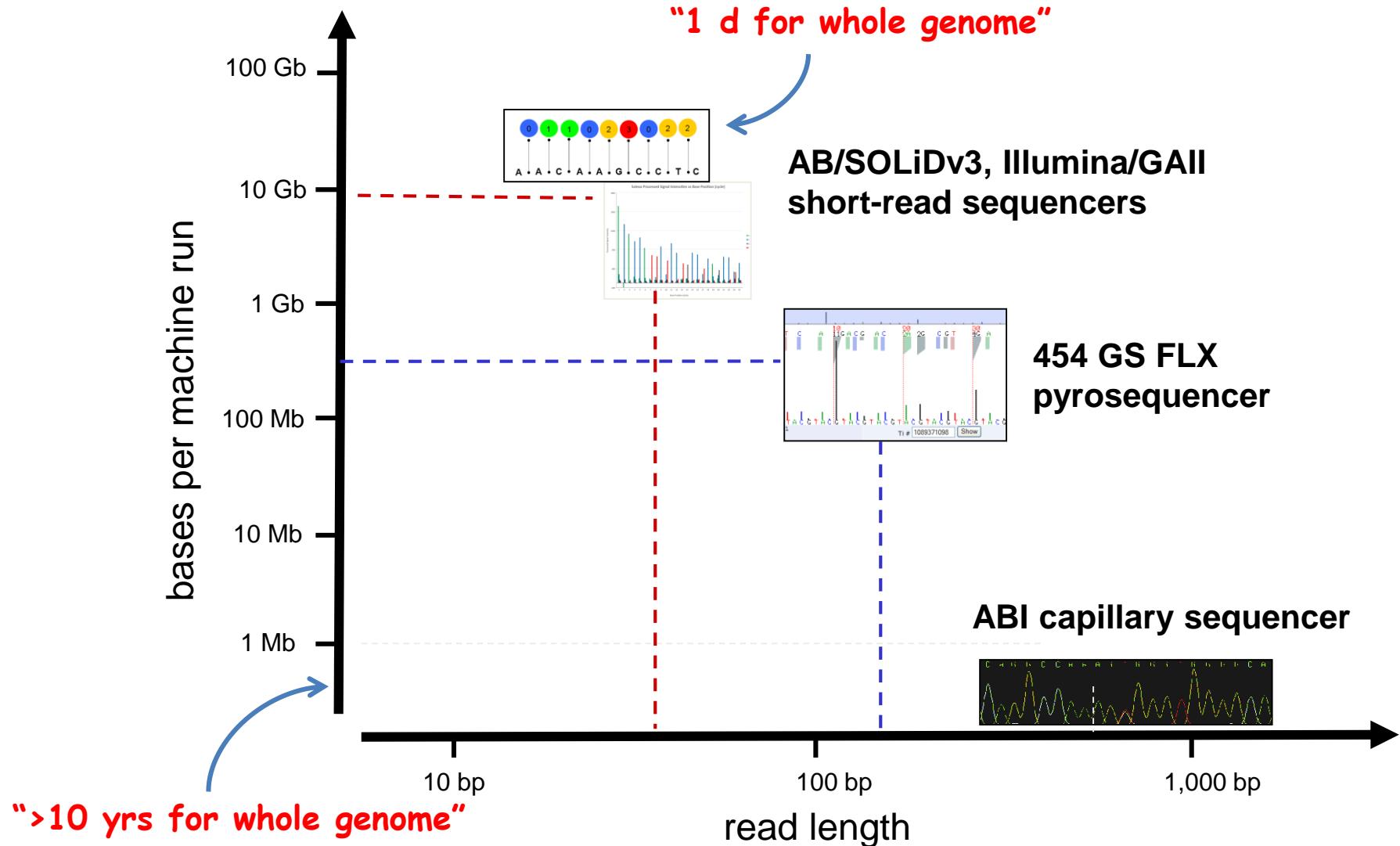
Comparison of sequencers

Vendor	Roche			Illumina			ABI			ABI
Technology	454			Solexa GA			SOLiD			3730
Platform	GS20	FLX	Ti	I	II	IIx	1	2	3	
Reads (M)	0.5	0.5	1.25	28	100	250	40	115	320	9.60E-05
Fragment										
Read length	100	200	400	35	50	100	25	35	50	800
Run time (d)	0.25	0.3	0.4	3.0	3.0	5.0	6.0	5.0	8.0	0.125
Yield (Gb)	0.05	0.1	0.5	1.00	5.00	25.00	1.00	4.00	16.00	7.68E-05
Rate (Gb/d)	0.2	0.33	1.25	0.33	1.67	5.00	1.67	0.8	2.00	6.14E-04
For 3Gb	15.0	9.1	2.4	9.09	1.8	0.6	1.8	3.8	1.5	4882.8

<http://www.politigenomics.com/next-generation-sequencing-informatics>



Comparison of sequencers



The 1000 Genomes project

Nature method, 2012

The 1000 Genomes Project was launched as one of the largest distributed data collection and analysis projects ever undertaken in biology. In addition to the primary scientific goals of creating both a deep catalog of human genetic variation and extensive methods to accurately discover and characterize variation using new sequencing technologies, the project makes all of its data publicly available. Members of the project data coordination center have developed and deployed several tools to enable widespread data access.

High-throughput sequencing technologies, including those from Illumina, Roche Diagnostics (454) and Life Technologies (SOLiD), enable whole-genome sequencing at an unprecedented scale and at dramatically reduced costs over the gel capillary technology used in the human genome project. These technologies were at the heart of the decision in 2007 to launch the 1000 Genomes Project, an effort to comprehensively characterize human variation in multiple populations. In the pilot phase of the project, the data helped create an extensive population-scale view of human genetic variation¹.

The larger data volumes and shorter read lengths of high-throughput sequencing technologies created substantial new requirements for bioinformatics, analysis and data-distribution methods. The initial plan for the 1000 Genomes Project was to collect 2× whole genome coverage for 1,000 individuals, representing ~6 giga-base pairs of sequence per individual and ~6 tera-base pairs (Tbp) of sequence in total. Increasing

~40× whole-genome sequence for 500 individuals in total (~25-fold increase in sequence generation over original estimates). In fact, the 1000 Genomes Pilot Project collected 5 Tbp of sequence data, resulting in 38,000 files and over 12 terabytes of data being available to the community¹. In March 2012 the still-growing project resources include more than 260 terabytes of data in more than 250,000 publicly accessible files.

As in previous efforts^{2–4}, the 1000 Genomes Project members recognized that data coordination would be critical to move forward productively and to ensure the data were available to the community in a reasonable time frame. Therefore, the Data Coordination Center (DCC) was set up jointly between the European Bioinformatics Institute (EBI) and the National Center for Biotechnology (NCBI) to manage project-specific data flow, to ensure archival sequence data deposition and to manage community access through the FTP site and genome browser.

Here we describe the methods used by the 1000 Genomes Project members to provide data resources to the community from raw sequence data to project results that can be browsed. We provide examples drawn from the project's data-processing methods to demonstrate the key components of complex workflows.

Data flow

Managing data flow in the 1000 Genomes Project such that the data are available within the project and to the wider community is the fundamental bioinformatics

The Cancer Genome Atlas

National Cancer Institute National Human Genome Research Institute

The Cancer Genome Atlas  Understanding genomics to improve cancer care

Launch Data Portal | Contact Us | For the Media

Search  Search

Home About Cancer Genomics Cancers Selected for Study Research Highlights Publications News and Events About TCGA

Home > About TCGA > Program Overview > How It Works > Genome Sequencing Centers

Genome Sequencing Centers

The Cancer Genome Atlas (TCGA) Genome Sequencing Centers (GSCs) perform large-scale DNA sequencing using the latest sequencing technologies. Supported by the National Human Genome Research Institute (NHGRI) large-scale sequencing program, the GSCs generate the enormous volume of data required by TCGA, while continually improving existing technologies and methods to expand the frontier of what can be achieved in cancer genome sequencing. All sequencing data are available in the TCGA Data Portal or from the TCGA page at NIH's database of Genotype and Phenotype (dbGaP).

Throughout the TCGA program, the GSCs have continued to evolve their approaches, as seen in this brief timeline:

- October 2008: TCGA publication on the glioblastoma multiforme genome includes polymerase chain reaction/Sanger dideoxy method for sequencing of 601 target genes. At the same time, GSCs are validating protocols using new second-generation sequencing instruments.
- March 2009: GSCs introduce hybrid-capture procedure and second-generation sequencing instruments (Illumina and ABI SOLiD) to enable analysis of more than 6,000 known cancer-associated target genes and at production scale. **This item is highlighted with a red box.**
- July 2009: GSCs submit first of 24 whole genome sequence (i.e., entire 6 billion nucleotides from both tumor and blood specimens from a cancer case) datasets from the glioblastoma multiforme and ovarian tumor projects.
- January 2010: GSCs validate whole exome capture methods, thereby expanding analysis of each tumor sample from 6,000 genes to all protein-coding and RNA genes.

Whole Exome vs. Whole Genome

Two DNA samples from every TCGA cancer case – one from the tumor specimen and the second from either blood or non-malignant tissue – are sent from a TCGA Biospecimen Core Resource site to a GSC. The non-tumor DNA serves as a control to confirm that mutations discovered in the tumor DNA are unique to the tumor and not normal genetic variations within the individual. All samples are analyzed by whole exome sequencing using second-generation sequencing instruments. Such instruments can generate the exome data from 8 to 16 samples in a single run in 8 to 14 days.

Next, more than 10 percent of the samples from each TCGA tumor project undergo whole genome sequencing to reveal mutations that lie outside of the exome regions. **This item is highlighted with a red box.**

Launch Data Portal

The Cancer Genome Atlas (TCGA) Data Portal provides a platform for researchers to search, download, and analyze data sets generated by TCGA.

Questions About Cancer

Visit www.cancer.gov

Call 1-800-4-CANCER

Use [LiveHelp Online Chat](#)

Multimedia Library

 Images

 Videos and Animations

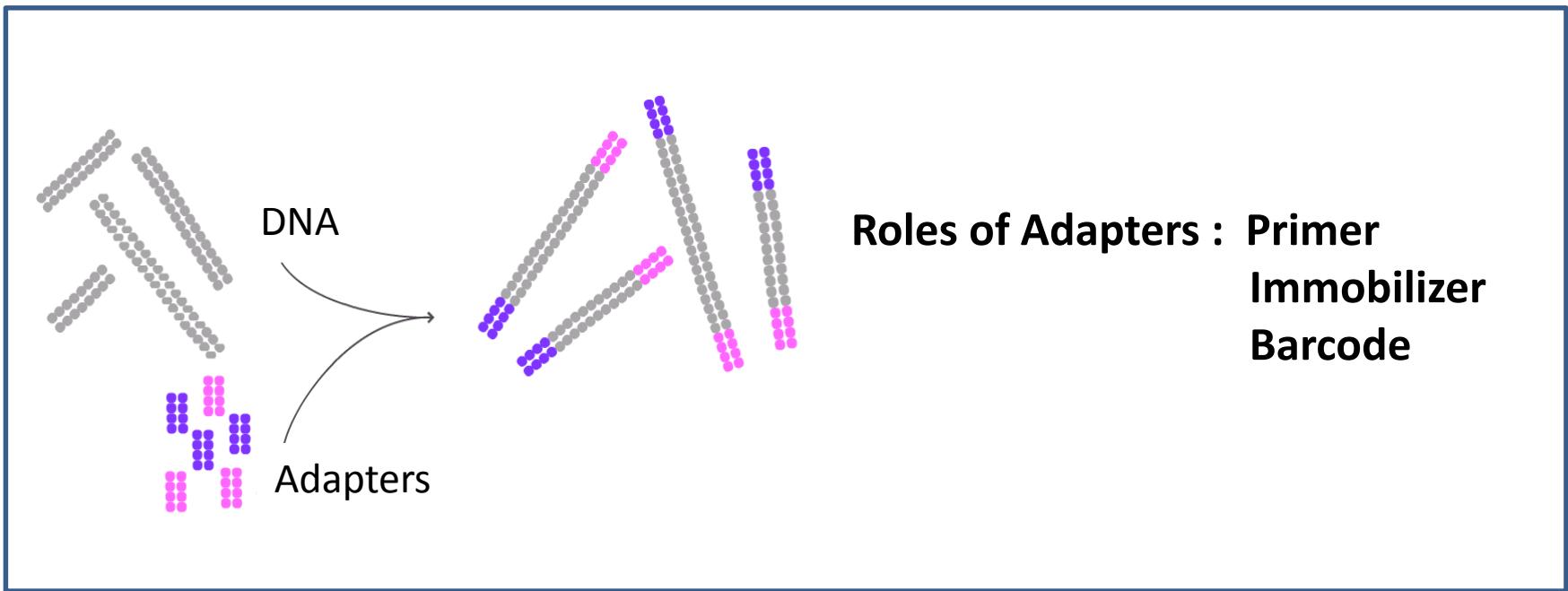
 Podcasts

 Interactive

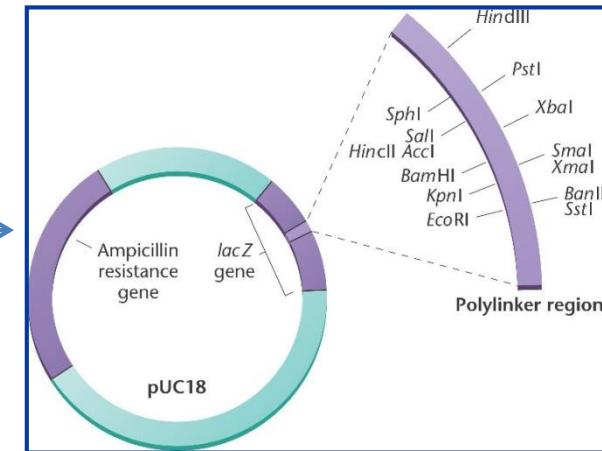
Stay Connected

Sign up for email updates

In NGS, adapters for DNA preparation

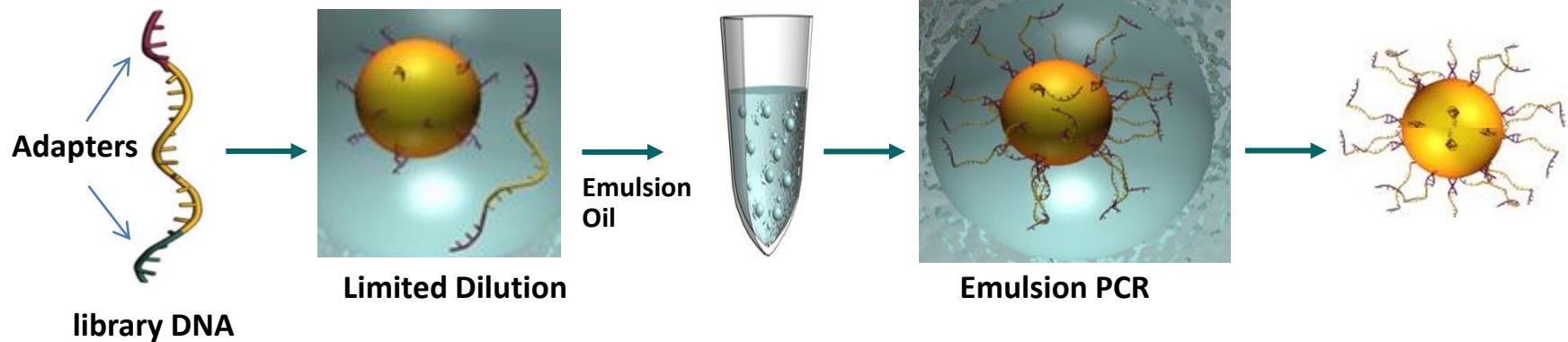


instead of this →

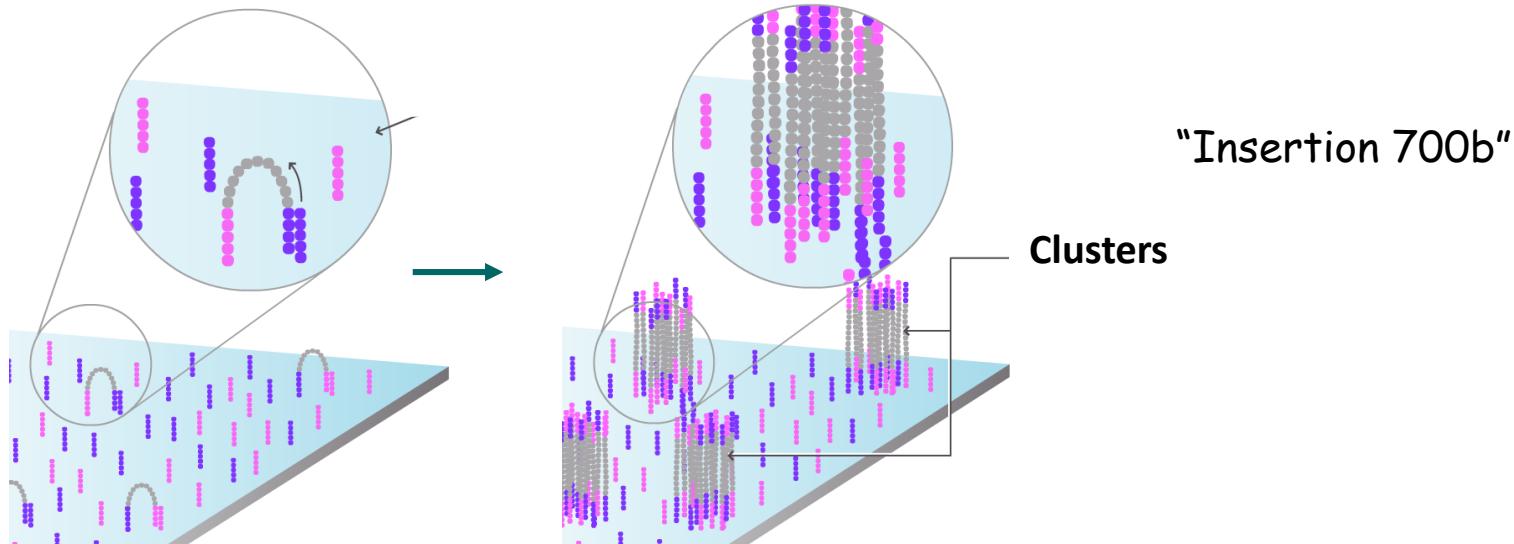


In NGS, instead of CLONING!

1. Emulsion PCR (Roche, AB)



2. Bridge Amplification (Illumina)

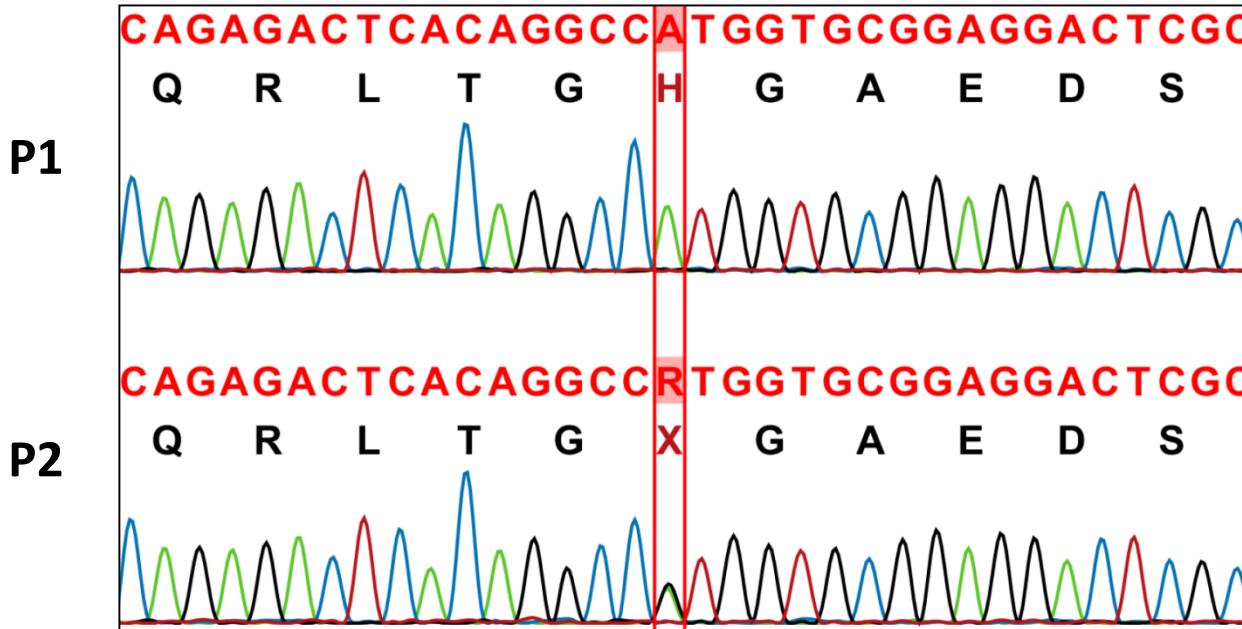


How to handle these cases in G-II vs G-III?

(1)

ATGGATCAGAGACTCACAGGCCA**A**TGGTGC^GGGAGGGACTCGCATGGAT
ATGGATCAGAGACTCACAGGCCA**G**TGGTGC^GGGAGGGACTCGCATGGAT

ATGGATCAGAGACTCACAGGCCA**A**TGGTGC^GGGAGGGACTCGCATGGAT
ATGGATCAGAGACTCACAGGCCA**G**TGGTGC^GGGAGGGACTCGCATGGAT

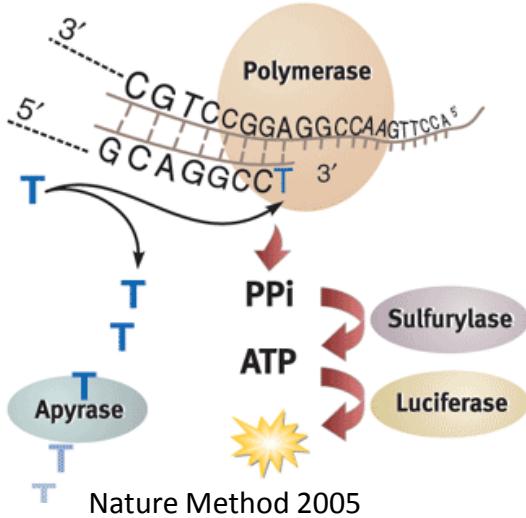


(2)

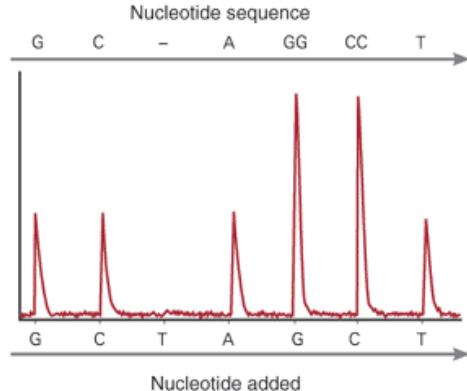
ATGGATTTGGCTTGGGATACCA**T**GGATCGCGTGAAATTCTATGGAT
ATGGATGGGATACCA**T**GGATCGCGTGAAATTCTTAGGACAATGGAT
ATGGAT**G**AAAATTCTTGTGGTACTGGAAATTCTTAATTCTATGGAT
ATGGATTTGGCTTGGGATACCA**T**GGATCGCGTGAAATTCTATGGAT
ATGGATGGGATACCA**T**GGATCGCGTGAAATTCTTAGGACAATGGAT
ATGGAT**G**AAAATTCTTGTGGTACTGGAAATTCTTAATTCTATGGAT

In NGS, for Base Calling

1. Pyrosequencing (Roche)



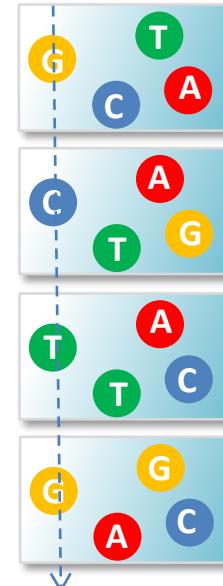
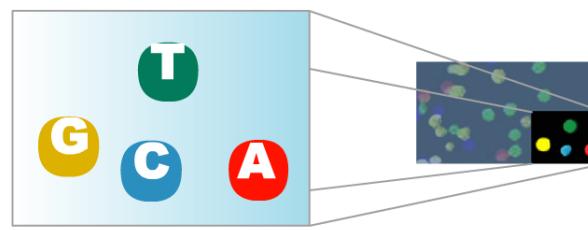
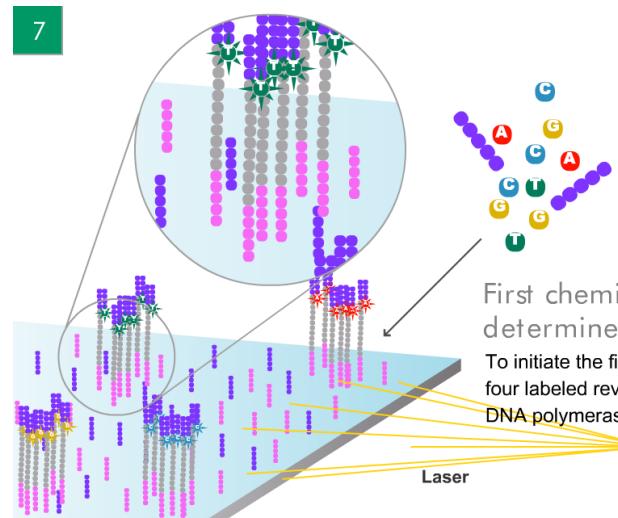
Nature Method 2005



Excessive use of dNTPs, Enzymes (Sulfurylase, Luciferase, Polymerase, Apyrase)

Weak: Homopolymeric repeats

2. Seq by synthesis (Illumina)

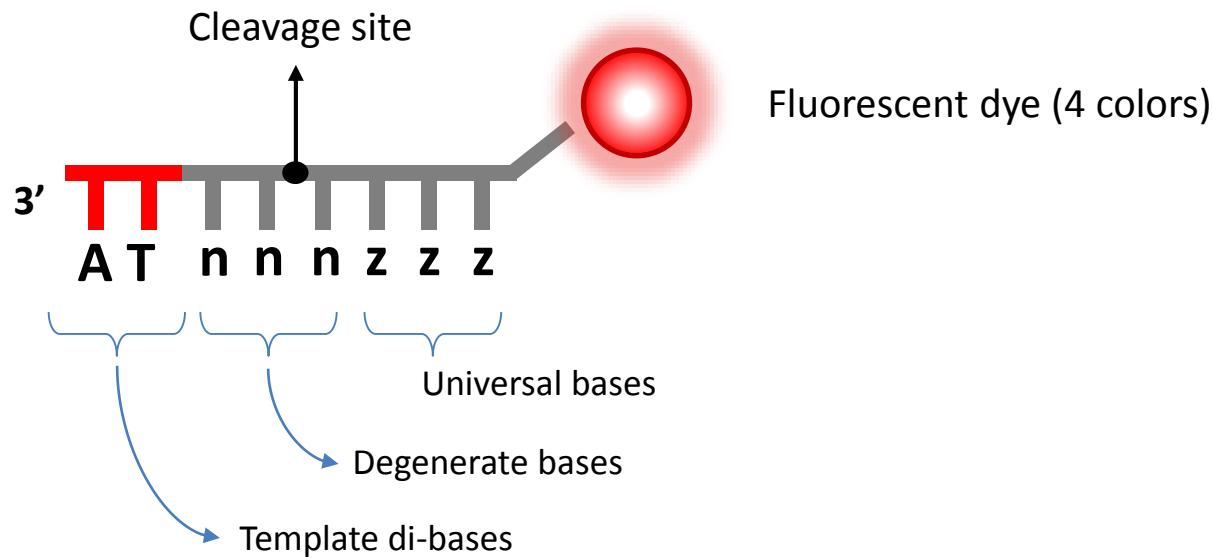


← PREV

NEXT →

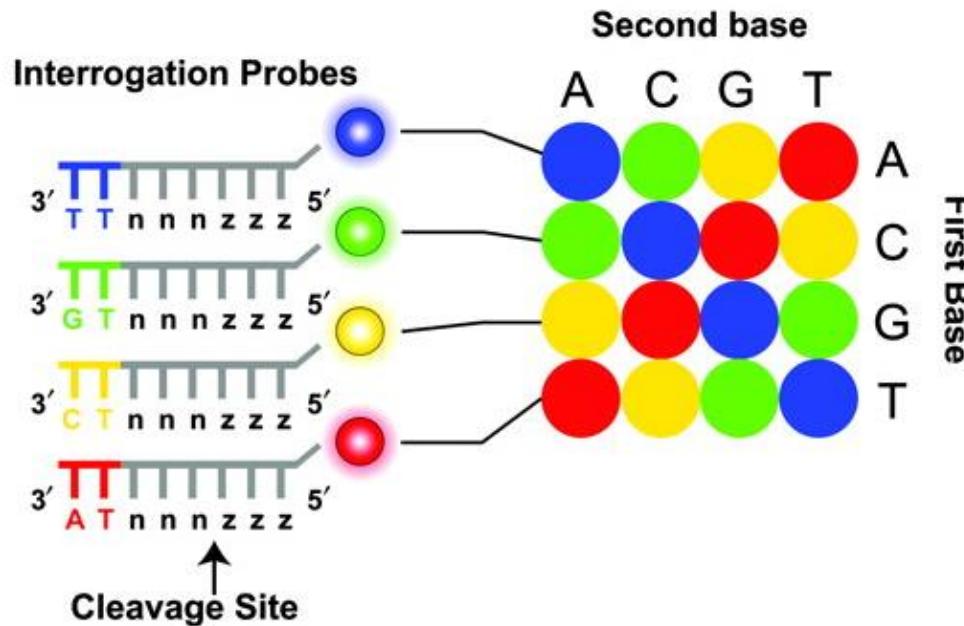
In NGS, for Base Calling

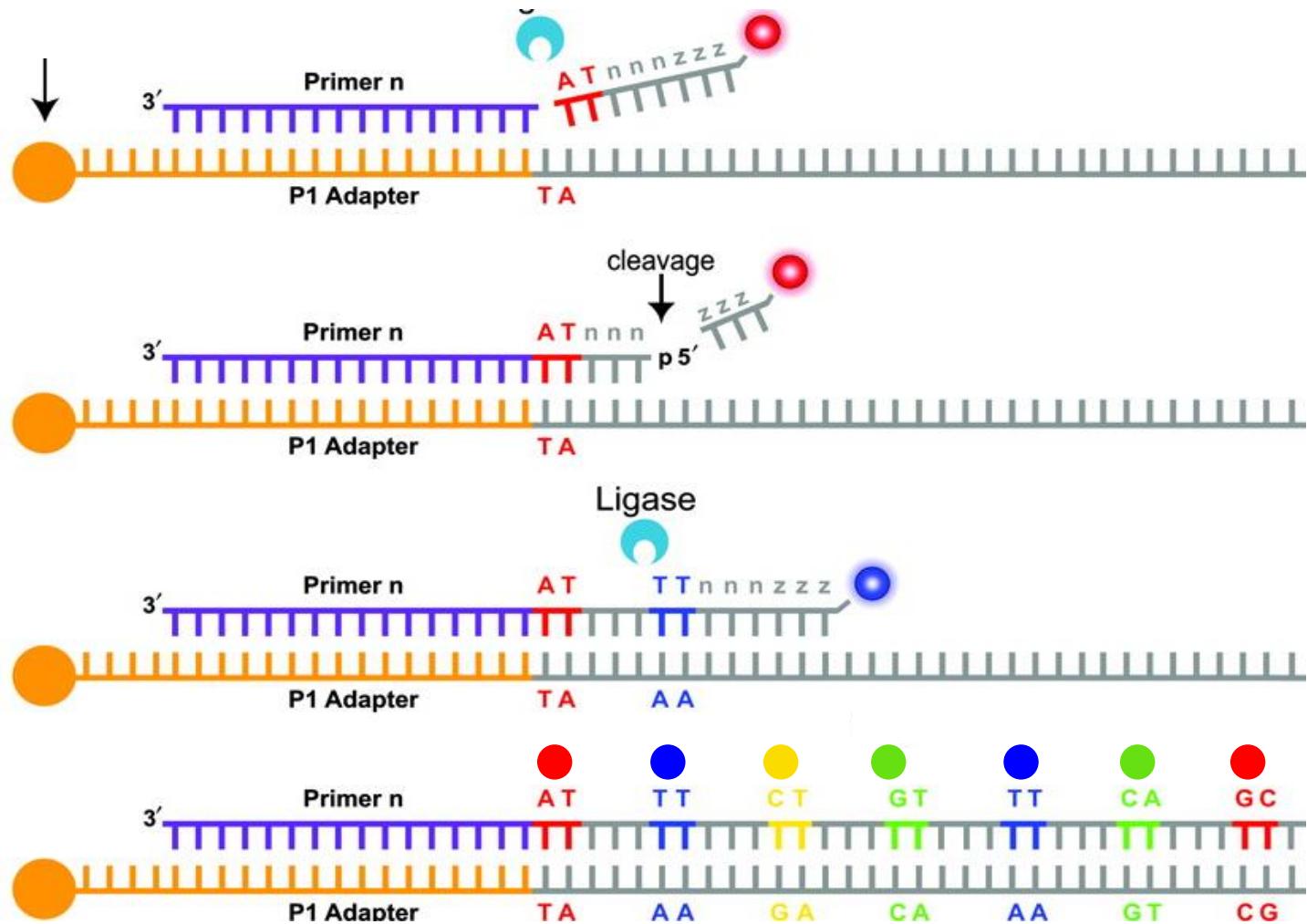
(Special attention to SOLiD Platform)

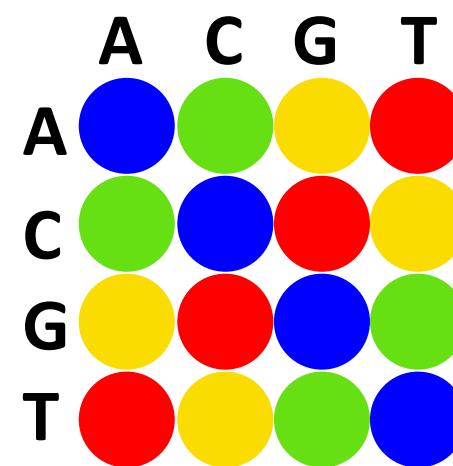
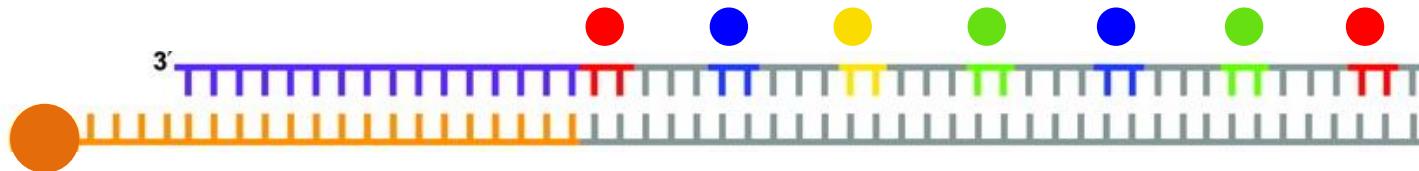


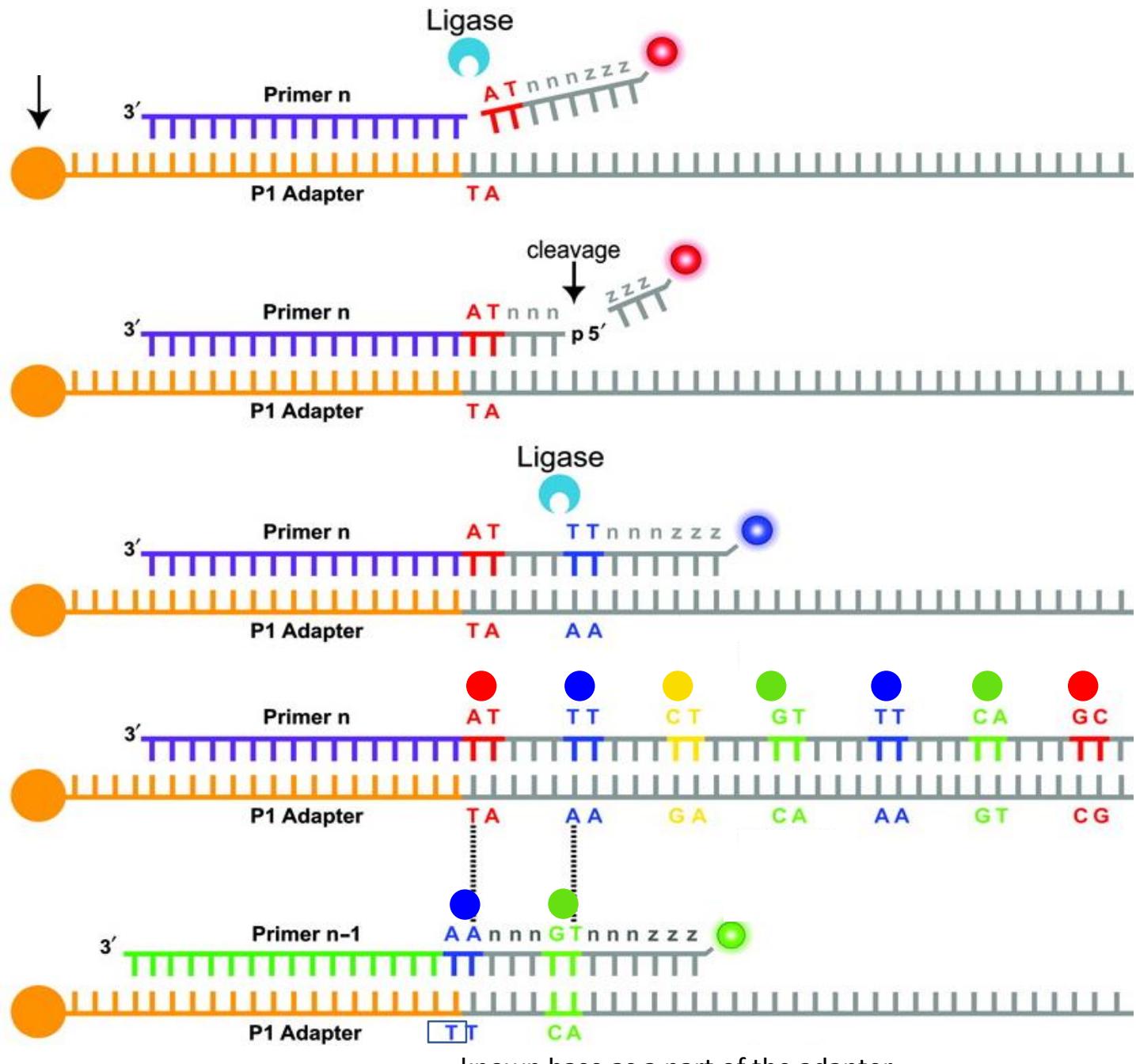
In NGS, for Base Calling

(Special attention to SOLiD Platform)



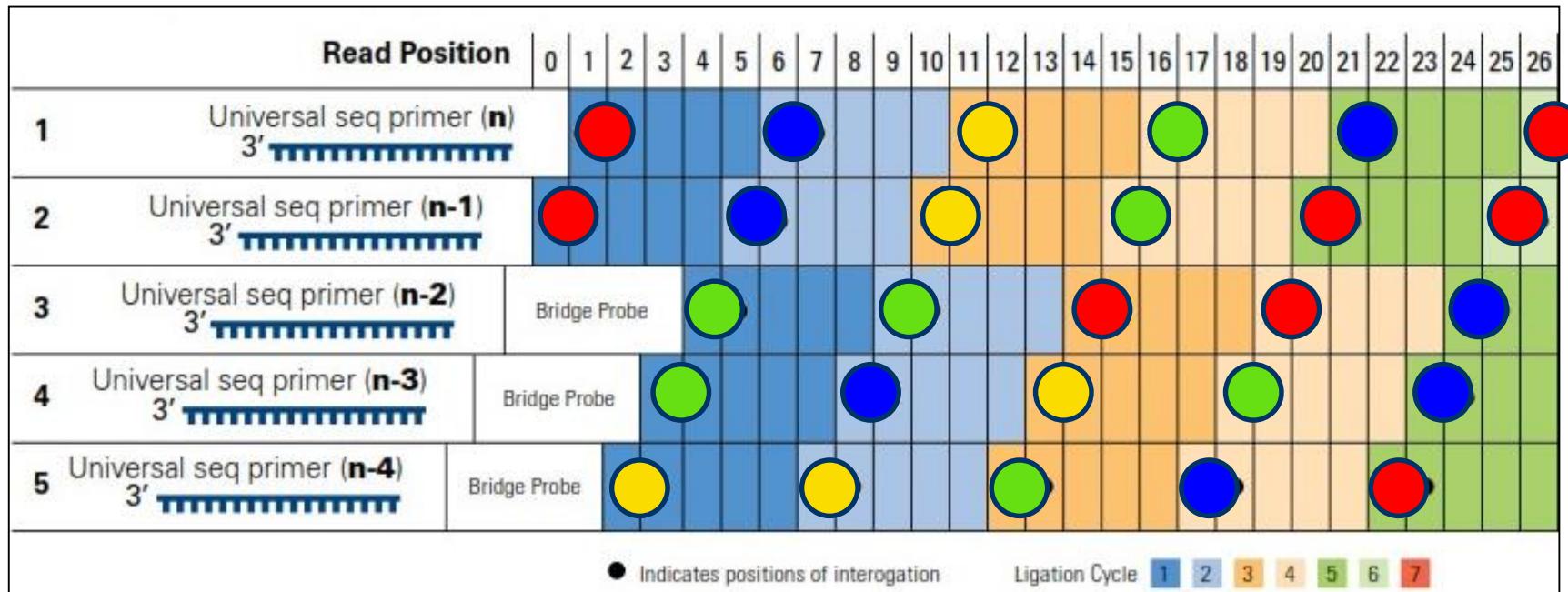






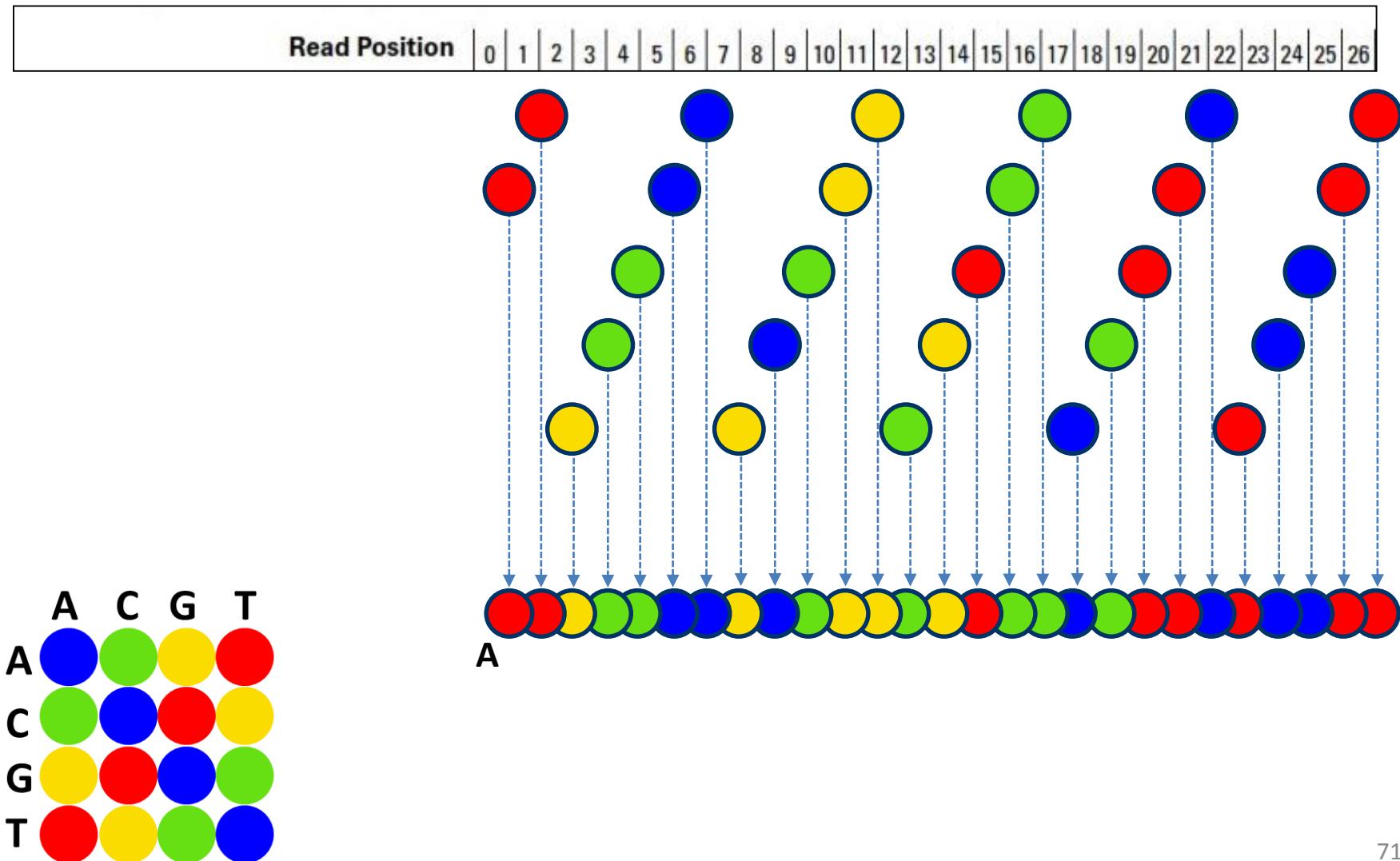
In NGS, for Base Calling

(Special attention to SOLiD Platform)



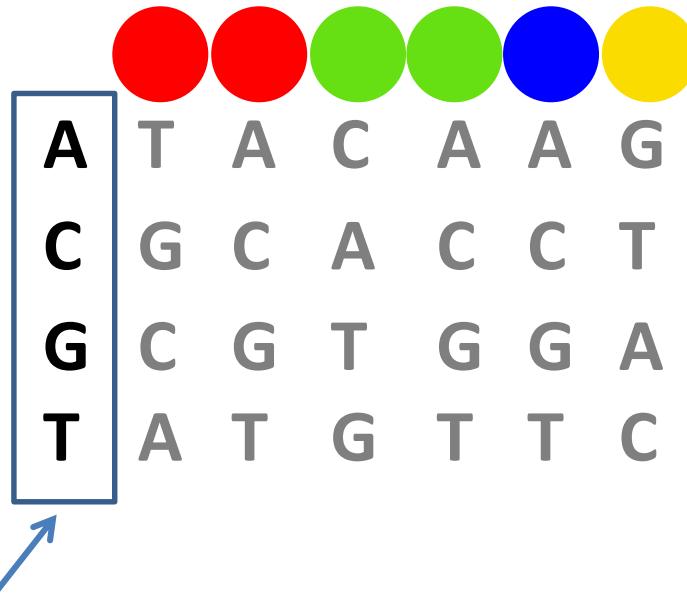
In NGS, for Base Calling

(Special attention to SOLiD Platform)



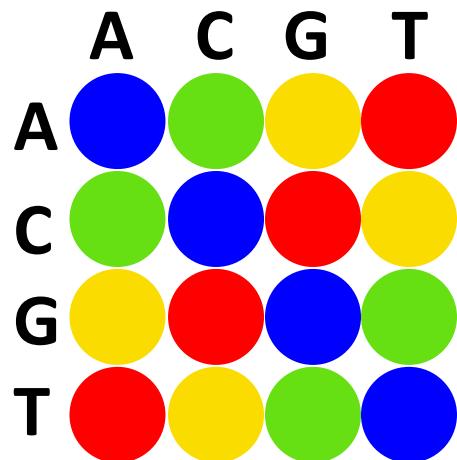
In NGS, for Base Calling

(Special attention to SOLiD Platform)

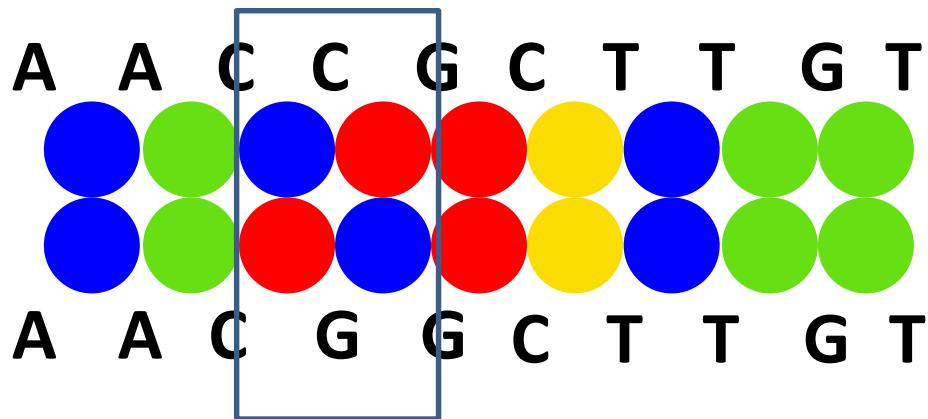


Must know identity of
first base to decode color
space

Strength in SNP



	A	C	G	T
A	0	1	2	3
C	1	0	3	2
G	2	3	0	1
T	3	2	1	0



Reassurance of the mutation result with the adjacent color

Example

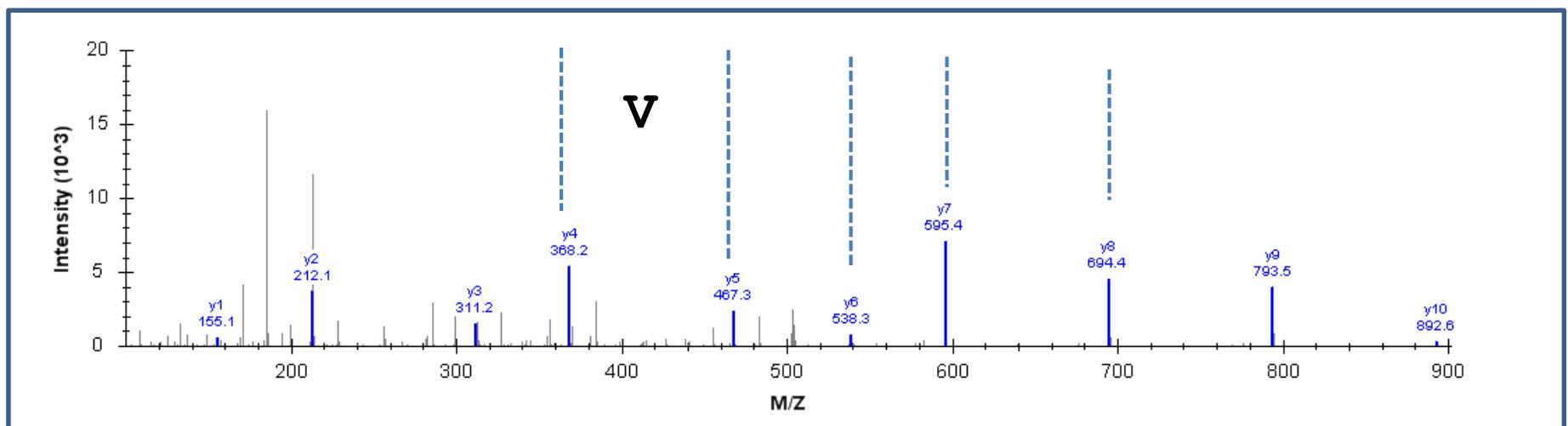
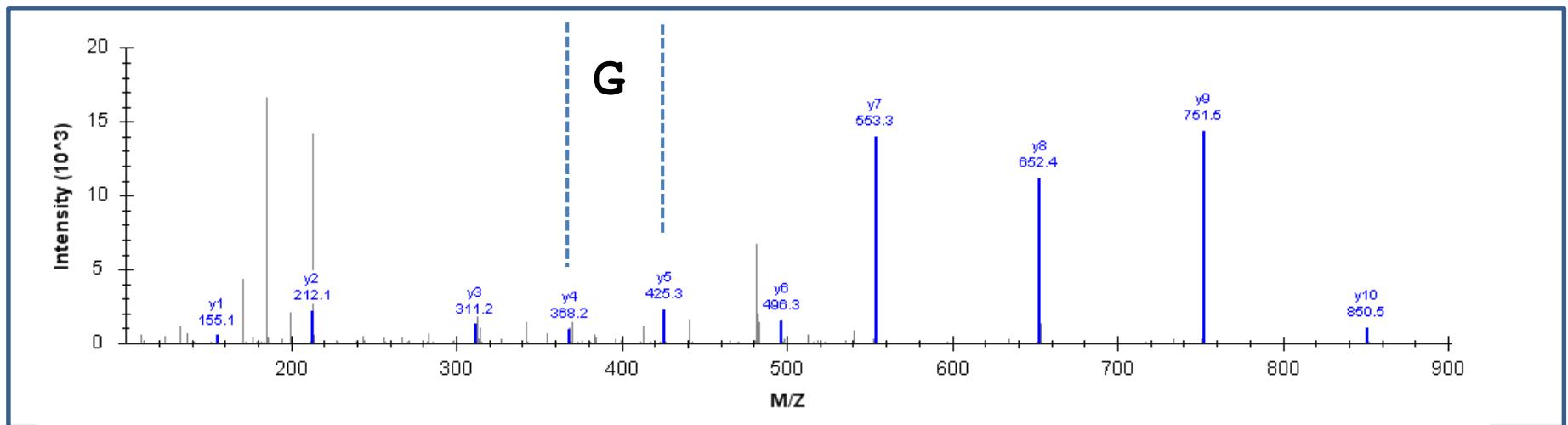
	A	C	G	T
A	0	1	2	3
C	1	0	3	2
G	2	3	0	1
T	3	2	1	0

T2011011313101223102311201000322003

Summary of G-III

- Library construction: No needs of run e-coli cloning, isolation and picking.
 - Random fragmentation of starting DNA (sonication, nebulization or shearing) followed by DNA repair and end-polishing
 - Custom made, platform-specific adaptor ligation
 - Straightforward processing: quicker, cheaper
- Library amplification: Beads or glass (solid support)
- Sequencing and detection: No separation involved, direct step-by-step detection of each nucleotide base incorporated during the sequencing reaction.
- Throughput: Upto 100M reactions imaged per instrument run (Massive parallel sequencing)
- Shorter read lengths than capillary based sequencers
- Digital read type that enables direct quantitative comparison

LVVVGAGGVGK
LVVVGA**V**GVGK



A proteomics approach for the identification and cloning of monoclonal antibodies from serum

Wan Cheung Cheung^{1,2}, Sean A Beausoleil^{1,2}, Xiaowu Zhang¹, Shuji Sato¹, Sandra M Schieferl¹, James S Wieler¹, Jason G Beaudet¹, Ravi K Ramenani¹, Lana Popova¹, Michael J Comb¹, John Rush¹ & Roberto D Polakiewicz¹

We describe a proteomics approach that identifies antigen-specific antibody sequences directly from circulating polyclonal antibodies in the serum of an immunized animal. The approach involves affinity purification of antibodies with high specific activity and then analyzing digested antibody fractions by nano-flow liquid chromatography coupled to tandem mass spectrometry. High-confidence peptide spectral matches of antibody variable regions are obtained by searching a reference database created by next-generation DNA sequencing of the B-cell immunoglobulin repertoire of the immunized animal. Finally, heavy and light chain sequences are paired and expressed as recombinant monoclonal antibodies. Using this technology, we isolated monoclonal antibodies for five antigens from the sera of immunized rabbits and mice. The antigen-specific activities of the monoclonal antibodies recapitulate or surpass those of the original affinity-purified polyclonal antibodies. This technology may aid the discovery and development of vaccines and antibody therapeutics, and help us gain a deeper understanding of the humoral response.

W. Cheung, Nature Biotech. **30**, 447 (2012)

Discussion Point

1. Challenges of applying proteomics-based identification for Abs.
2. Reference database used in this study.
3. Why the use of 454 NGS was critical to this study (compared to CE-based sequencing?)