

joanne@msl.ubc.ca

Bioinformatics

Common tools, useful databases, and tricks of the trade.



bioteach.ubc.ca/bioinfo2008

Workshop Schedule

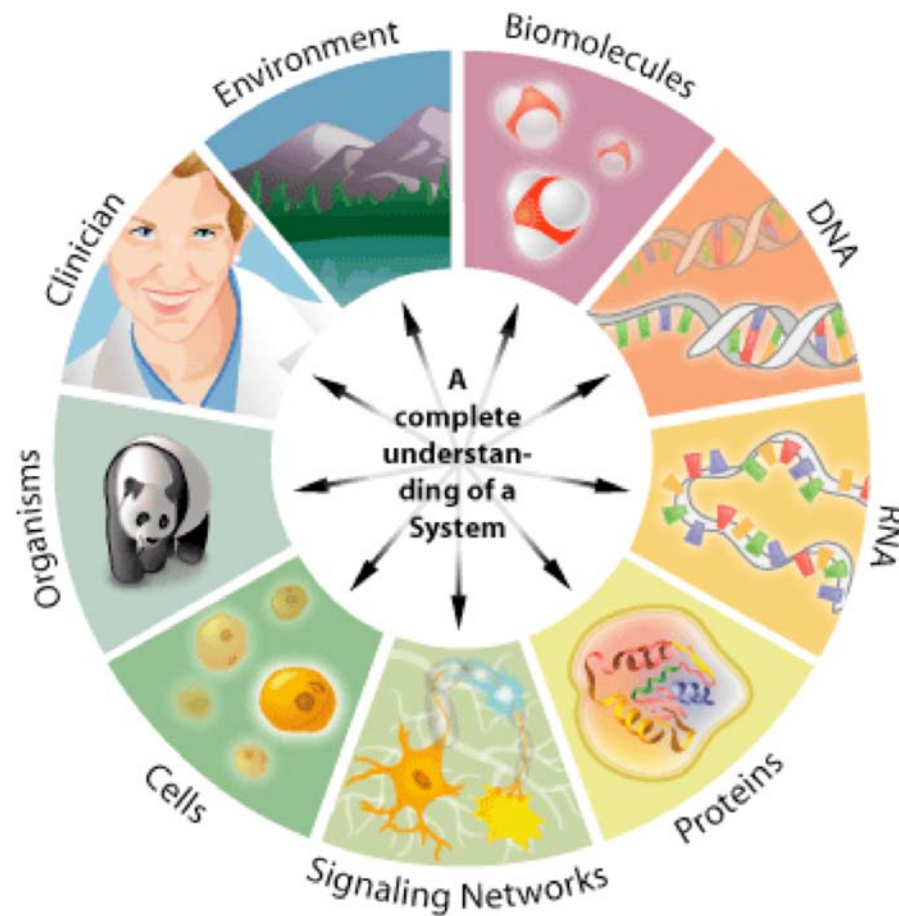
- Laptops, available here for your use 9am - 4:30pm
- wireless login
mslguest
4myguest
- Vancouver guide books available



Today's Topics

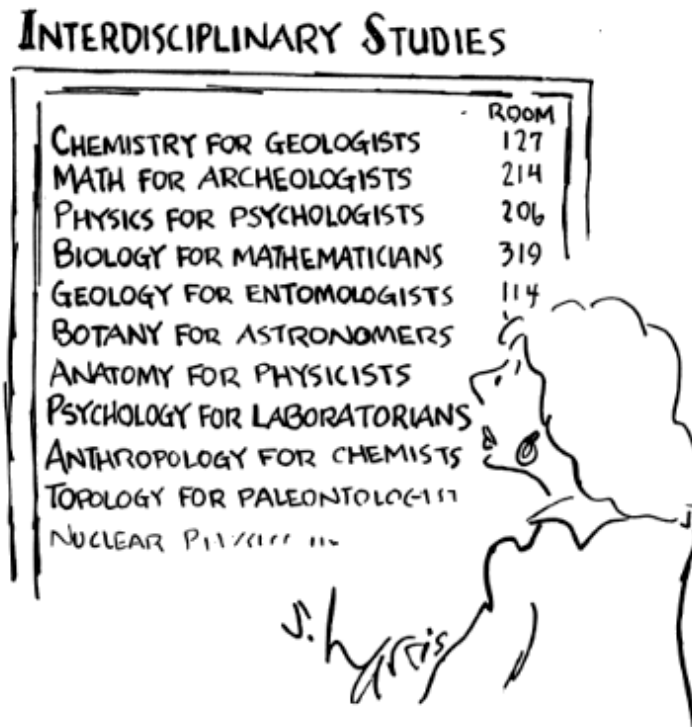
- **DNA Sequencing** - Generating Data & Emerging Technologies
- **Sequence Databases** - Public Resources at the NCBI
- **GUIDED TOUR** - Advanced Tips & Tricks for Searching Entrez
- **PRACTICAL EXERCISES** - Navigating Links, Retrieving Data with Entrez, and Searching PubMed

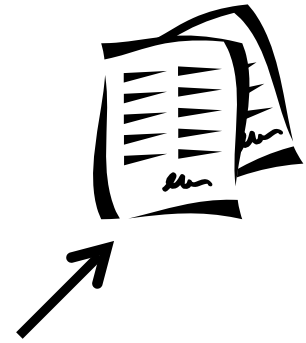
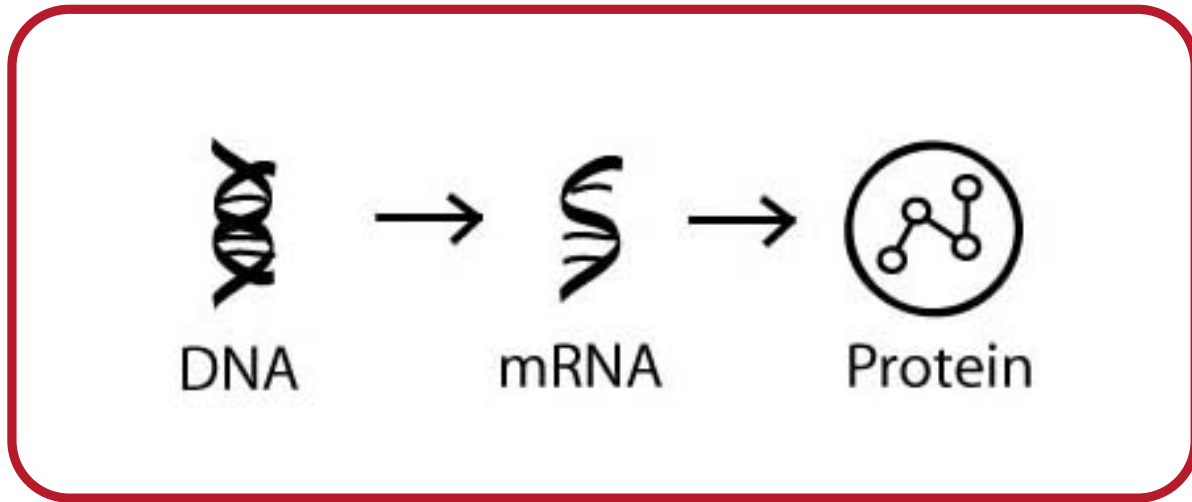
What is Bioinformatics?

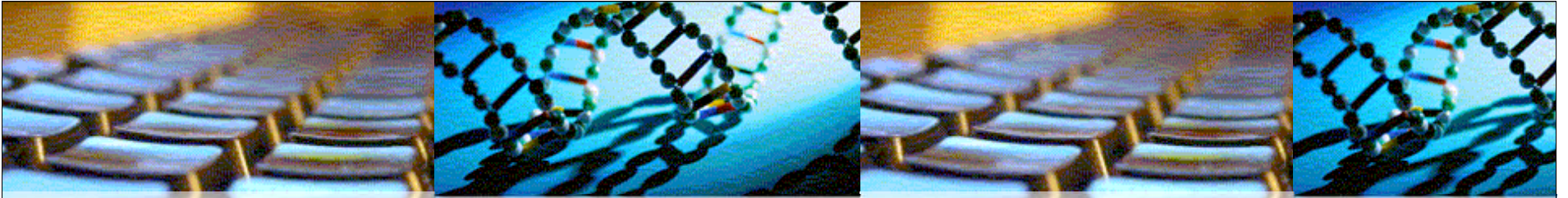


Bioinformatics for ~~Biologists~~

Government
Employees

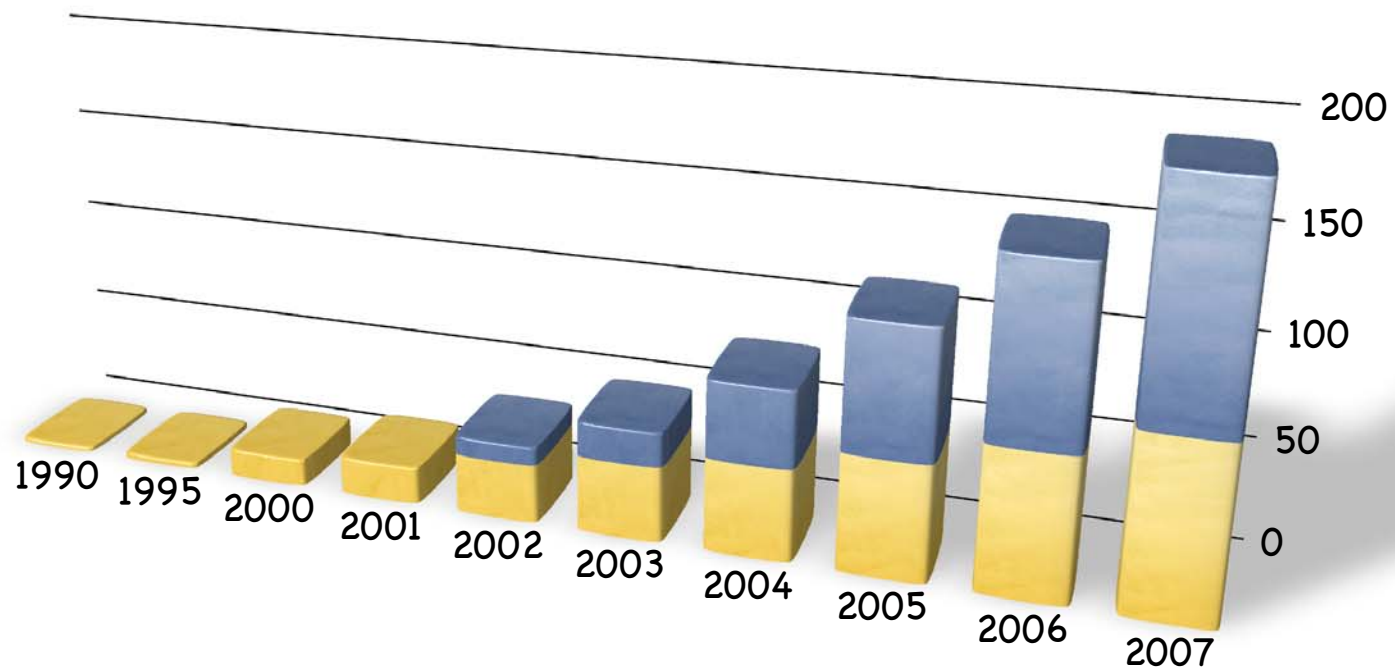






By the end of the
morning session,
you will define
bioinformatics
in your own words

Growth of GenBank

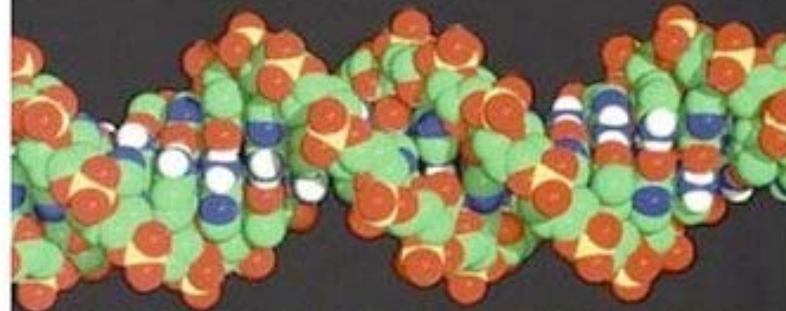


In 2005, International
sequence databases
exceed 100 gigabases

NATIONAL BESTSELLER

"A fascinating tour of the human genome. . . . If you want to catch a glimpse of the biotech century that is now dawning . . . *Genome* is an excellent place to start." —*Wall Street Journal*

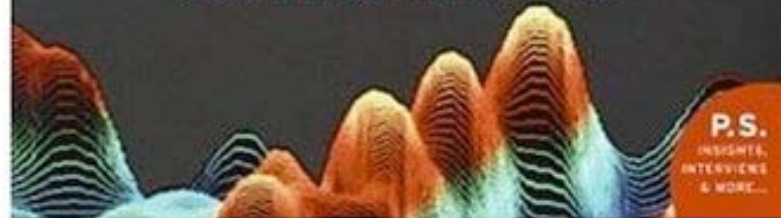
GENOME



THE AUTOBIOGRAPHY OF A
SPECIES IN 23 CHAPTERS

MATT RIDLEY

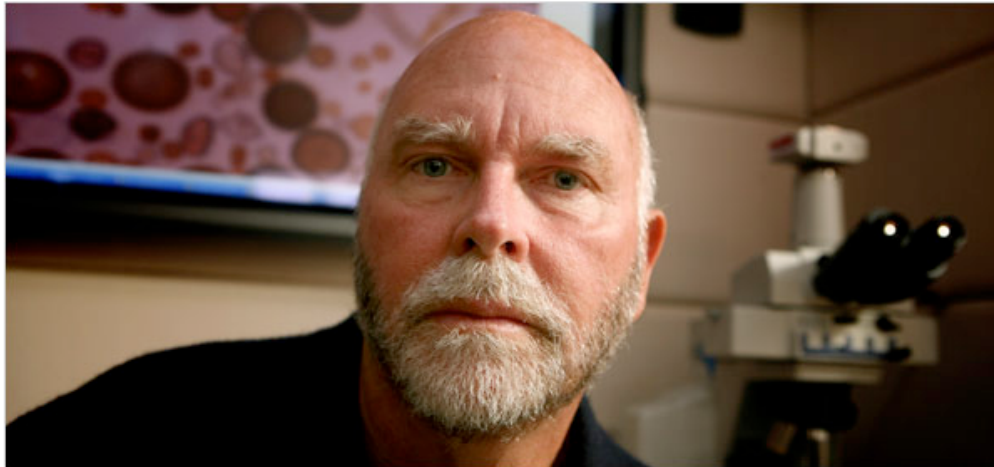
AUTHOR OF *THE AGILE GENE*
AND *FRANCIS CRICK*



P.S.
INSIGHTS,
INTERVIEWS
& MORE...

Personalized Medicine?

In the Genome Race, the Sequel Is Personal



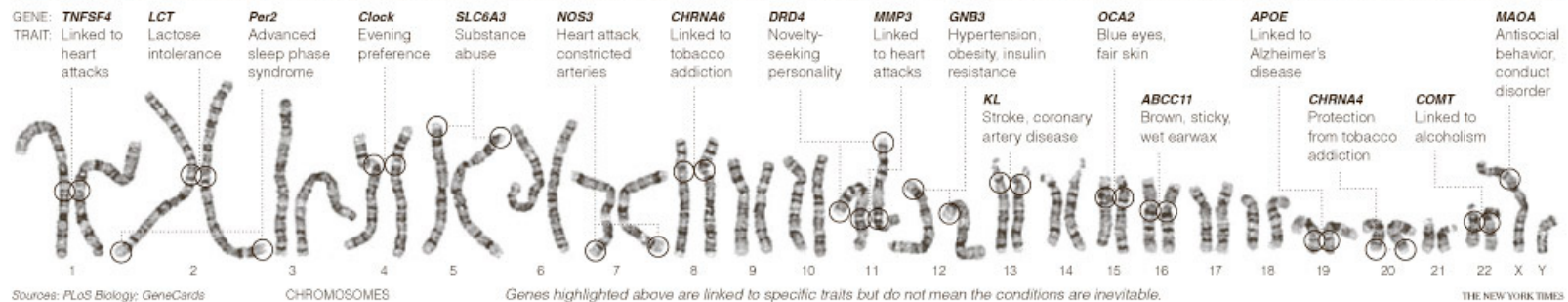
Thor Swift for The New York Times

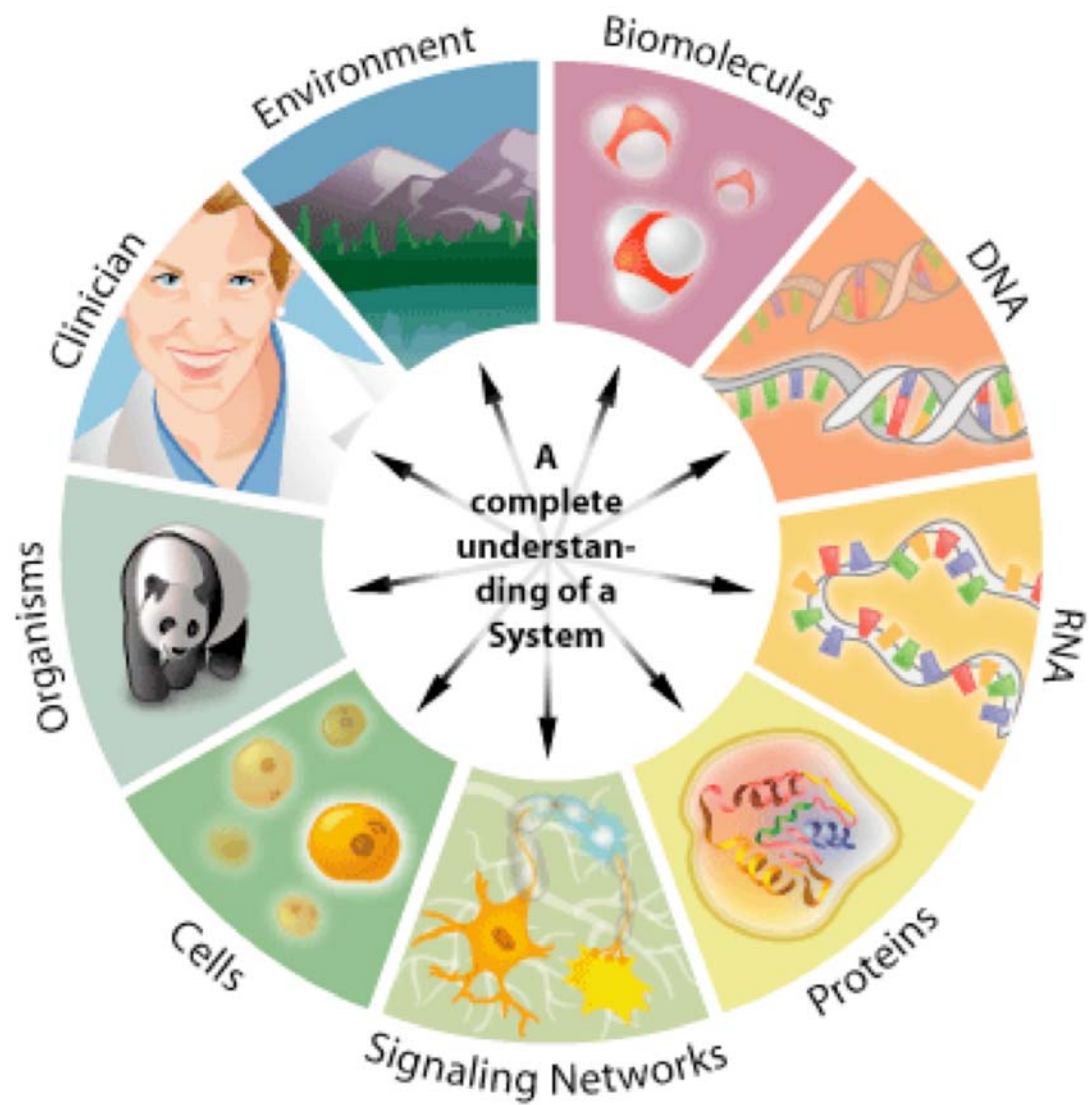
A team led by J. Craig Venter, above, has finished the first mapping of a full, or diploid, genome, made up of DNA inherited from both parents. The genome is Dr. Venter's own.

The New York Times

September 3, 2007

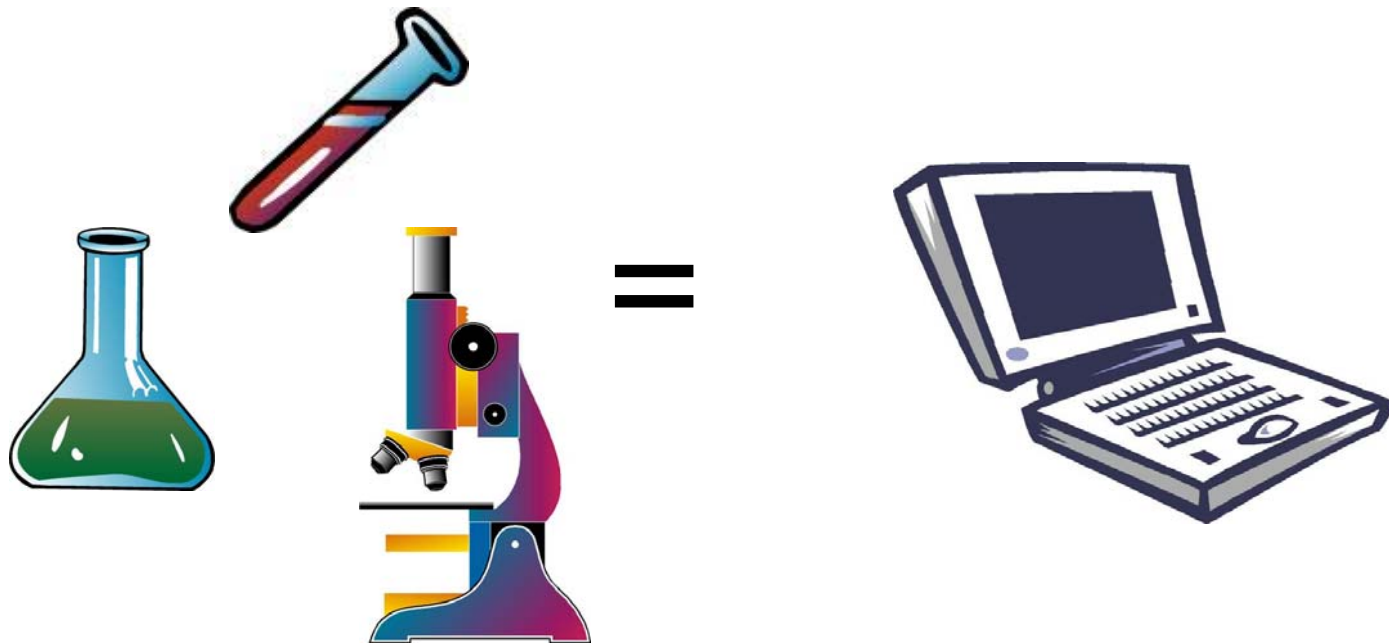
DECODING HIMSELF A team led by J. Craig Venter, above, has finished the first mapping of a full, or diploid, genome, made up of DNA inherited from both parents. The genome is Dr. Venter's own.



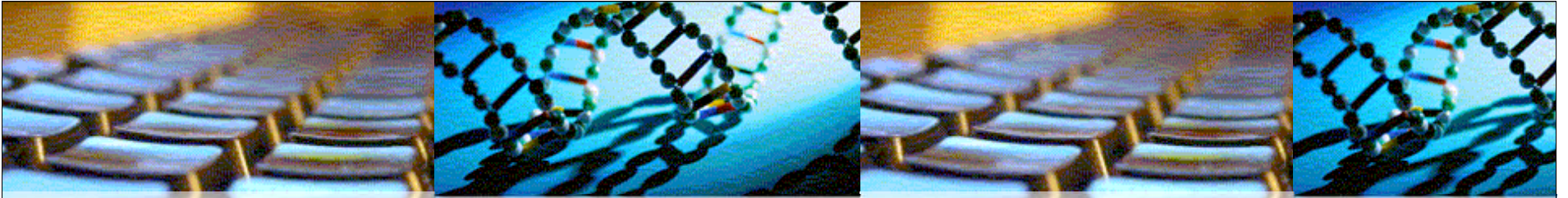


Linking Biological Information

- Nucleic acid & Protein Sequence data
- Sequence similarity implies homology or similar function
- Literature / Books
- Papers on related topics
- Macromolecular 3D Structure Data
- Similar protein fold implies homology
- Regulatory Pathways, Expression Data
- Two genes / enzymes in the same pathway
- Taxonomy
- Linkage of organisms by evolution



What is Bioinformatics?



Bioinformatics is an fast-paced interdisciplinary research field that involves the integration of computers, software tools, and databases in an effort to address biological questions.



Genomics refers to the analysis of all of the genes and transcripts included within the genome. **Proteomics**, on the other hand, refers to the analysis of the complete set of proteins or proteome.

Bioinformatics Questions

- What is encoded by the genome?
 - Genes, regulatory, and functional regions
- How is genome information expressed?
 - Function of genes and gene products (proteins)
 - Structure of proteins
- How can we interpret the information encoded in the genome?
 - Linking knowledge to the biological entities.
 - Systems biology approach
 - drugs, metabolites, ...
- How does the genome interact with its environment?

Summary

- An article called, “What is Bioinformatics?” is available from the Science Creative Quarterly.

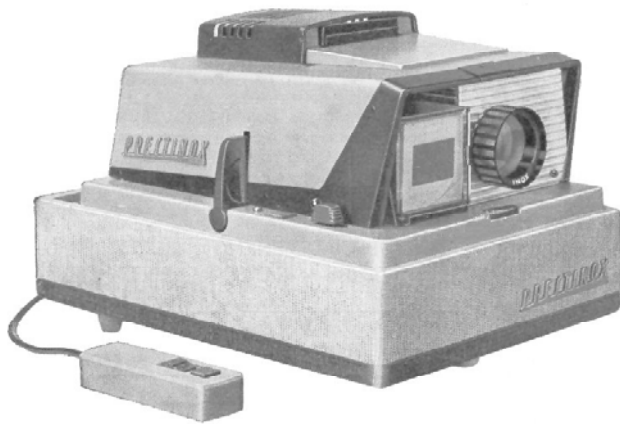
<http://www.scq.ubc.ca/what-is-bioinformatics/>

DNA Sequencing

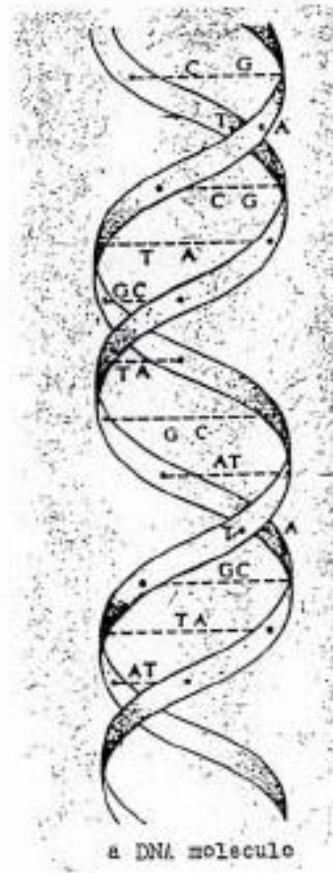
Generating Data & Emerging Technologies



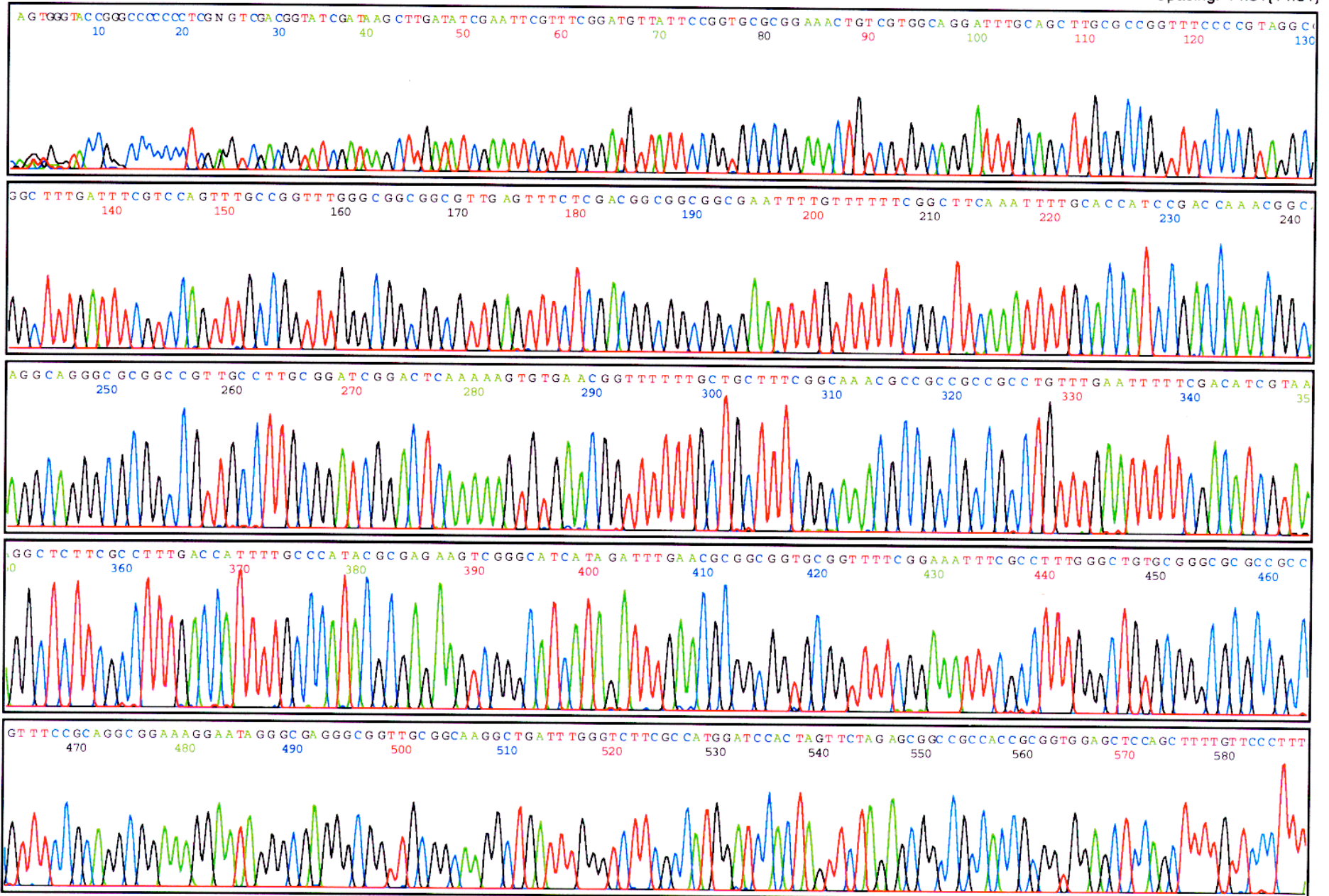
Technology



ggagcctcgg gaggtggtgg agtgacctgg cccagtgct gcgtccttat cagccgagcc
ggtcccagct cttgctcctg cctgtttgcc tggaaatggc cacgcttctc cttctccttg
gggtgctggt ggtaagccca gacgctctgg ggagcacaac agcagtgcag acaccacact
ccggagagcc tttggtctct actagcgagc ccctgagctc aaagatgtac accacttcaa
taacaagtga ccctaaggcc gacagcactg gggaccagac ctcagcccta cctccctcaa
cttccatcaa tgagggatcc cctctttgga cttccattgg tgccagcact ggttcccctt
tacctgagcc aacaacctac caggaagttt ccatcaagat gtcatacagt cccagggaaa
ccctcatgc aaccagtcac cctgctgttc ccataacage aaactctcta ggatcccaca
ccgtgacagg tggaaccata acaacgaact ctccagaaac ctccagtagg accagtggag
ccctgttac cacggcagct agctctctgg agacctccag aggcacctct ggaccccctc
ttaccatggc aactgtctct ctggagactt ccaaaggcac ctctggacct cctgttacca
tggcaactga ctctctggag acctccactg ggaccactgg accccctgtt accatgacaa
ctggctctct ggagccctcc agcggggcca gtggaccca ggtctctagc gtaaaactat
ctacaatgat gtctccaacg acctccacca acgcaagcac tgtgcccttc cggaaaccag
atgagaactc acgaggcatg ctgccagtgg ctgtgcttgt ggccctgctg gcggtcatag
tcctcgtggc tctgctcctg ctgtggcgcc ggcggcagaa gcggcggact ggggccctcg
tgctgagcag aggcggcaag cgtaacgggg tggtaggcgc ctgggctggg ccagcccagg
tcctgagga gggggccgtg acagtgaccg tgggagggtc cgggggcgac aagggtctctg
ggttccccga tggggagggg tctagccgtc ggcccacgct caccacttctc tttggcagac
ctggctctct ggagccctcc agcggggcca gtggaccca ggtctctagc gtaaaactat
ctacaatgat gtctccaacg acctccacca acgcaagcac tgtgcccttc cggaaaccag
atgagaactc acgaggcatg ctgccagtgg ctgtgcttgt ggccctgctg gcggtcatag



= SNOT



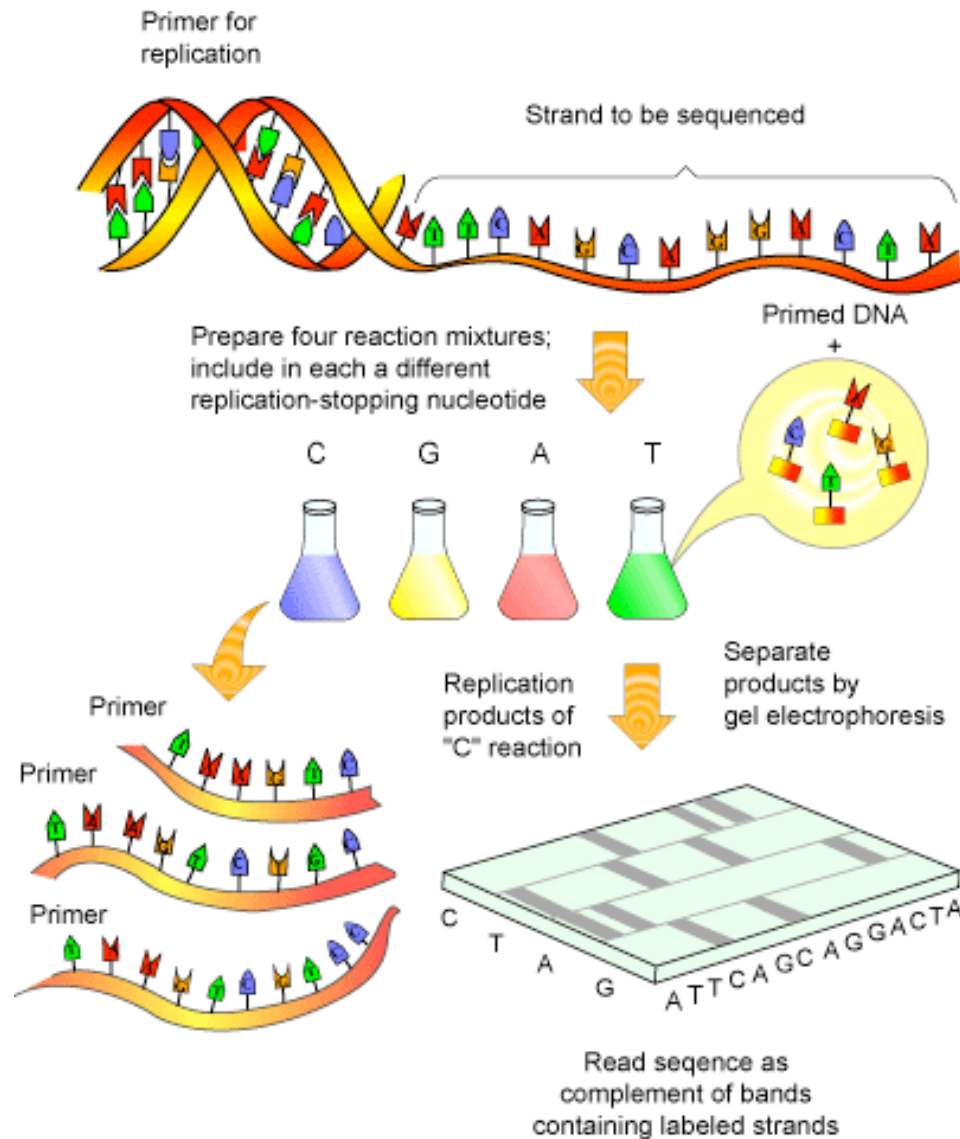


Figure 1. The Sanger sequencing reaction. Single stranded DNA is amplified in the presence of fluorescently labelled ddNTPs that serve to terminate the reaction and label all the fragments of DNA produced. The fragments of DNA are then separated via polyacrylamide gel electrophoresis and the sequence read using a laser beam and computer.

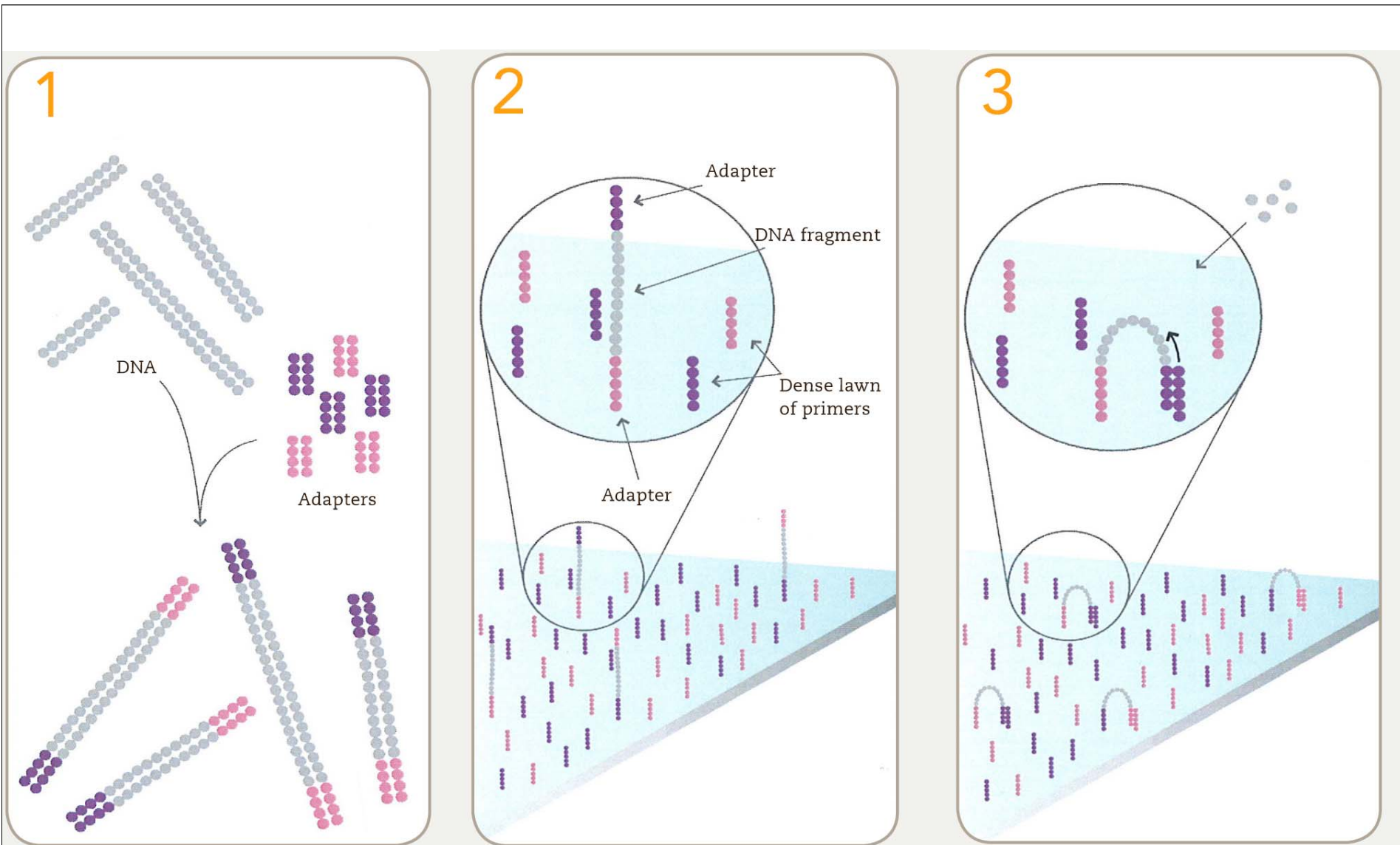
source: <http://www.scq.ubc.ca/genome-projects-uncovering-the-blueprints-of-biology/>

Every few years, a new technology comes along that dramatically changes how fundamental questions in biology are addressed. The impact of the technology is not always appreciated at first ...

- Stanley Fields

Solexa Technology

- DNA sequencing by synthesis
- approach built around very large number of short sequence reads
- key points:
 - ✓ solid phase amplification = no cloning necessary
 - ✓ reversible chemistry
 - ✓ data generated by imaging
 - ✓ read lengths 30-50 bp
 - ✓ < 1% cost, ultra high-throughput



PREPARE GENOMIC DNA SAMPLE
Randomly fragmented genomic DNA and ligate adaptors to both ends of the fragments

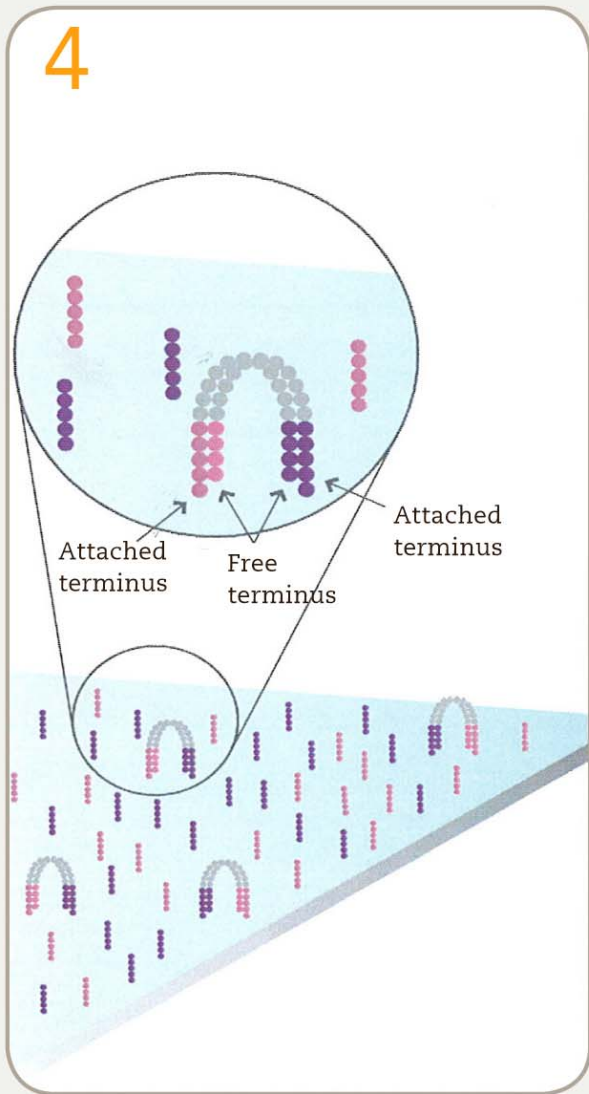
© 2007, Illumina Inc. All rights reserved.

ATTACH DNA TO SURFACE
Bind single stranded fragments randomly to the inside surface of the flow cell channels.

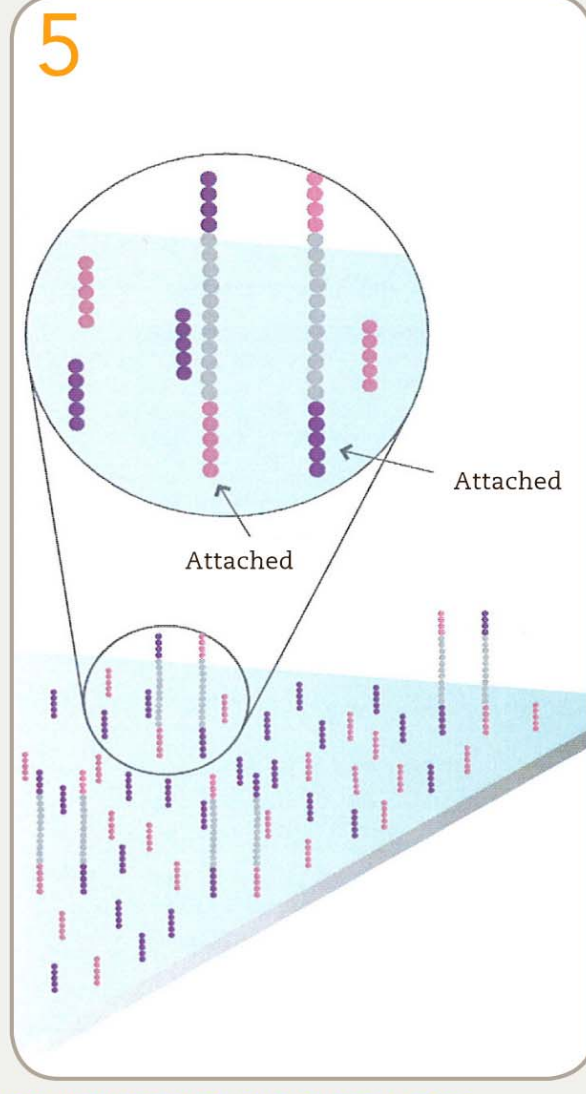
© 2007, Illumina Inc. All rights reserved.

BRIDGE AMPLIFICATION
Add unlabeled nucleotides and enzyme to initiate solid-phase bridge amplification.

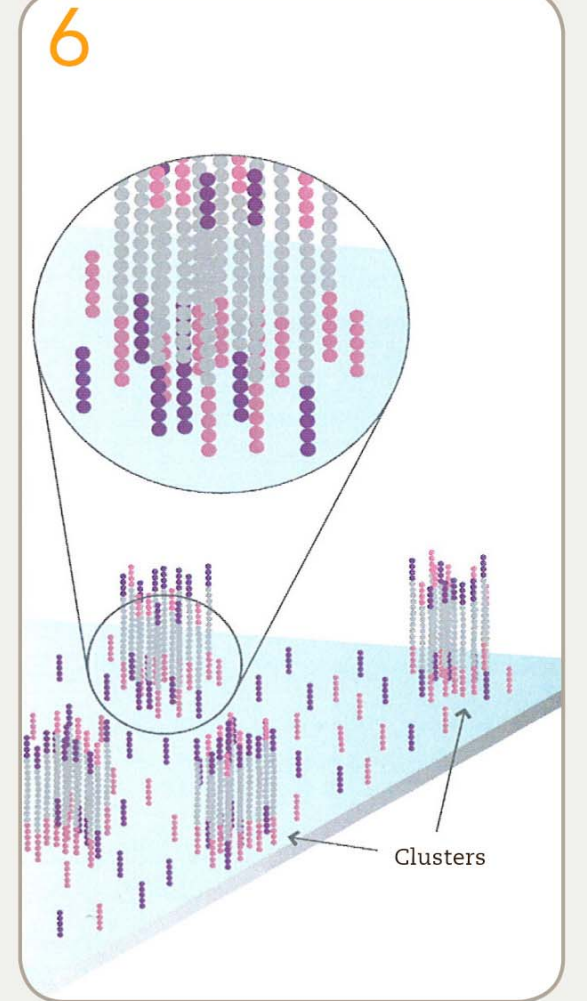
© 2007, Illumina Inc. All rights reserved.



© 2007, Illumina Inc. All rights reserved.

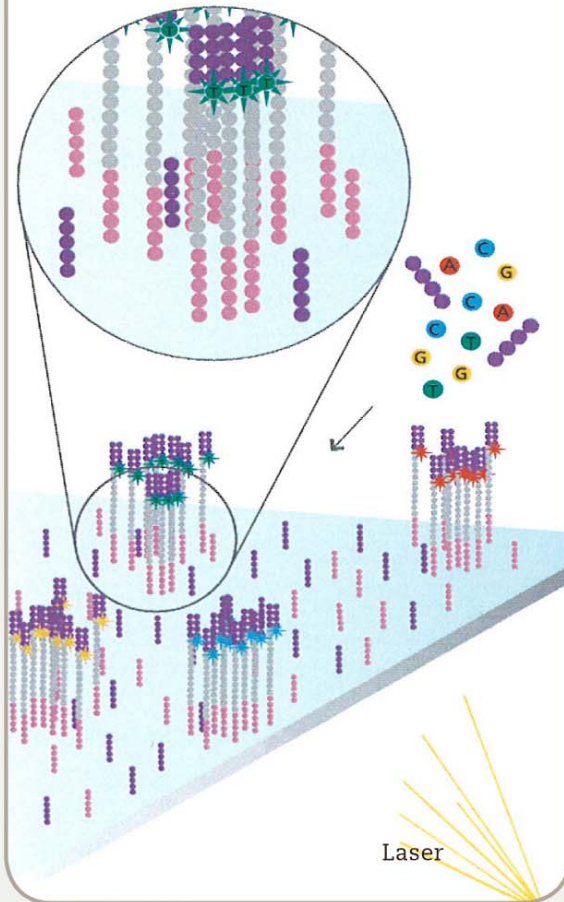


© 2007, Illumina Inc. All rights reserved.



© 2007, Illumina Inc. All rights reserved.

7



**FIRST CHEMISTRY CYCLE:
DETERMINE FIRST BASE**

To initiate the first sequencing cycle, add all four labeled reversible terminators, primers and DNA polymerase enzyme to the flow cell.

© 2007, Illumina Inc. All rights reserved.

8

After laser excitation, capture the image of emitted fluorescence from each cluster on the flow cell. Record the identity of the first base for each cluster.

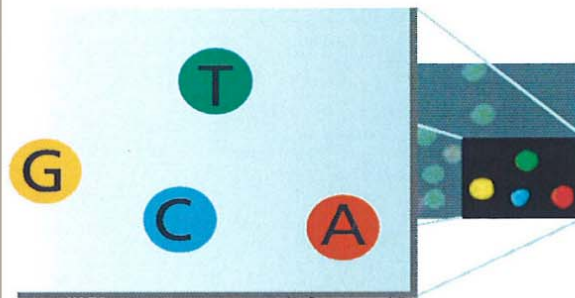
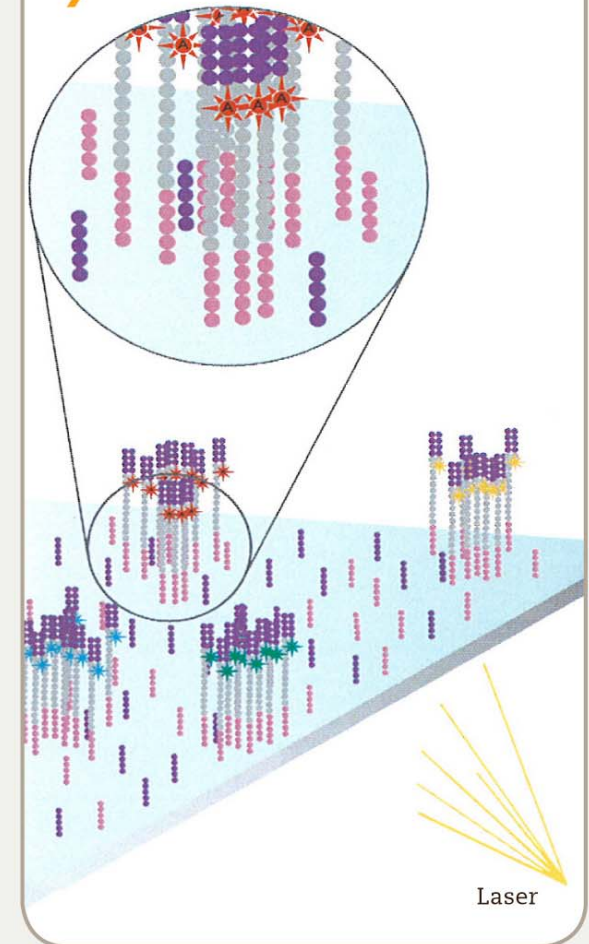


IMAGE OF FIRST CHEMISTRY CYCLE

© 2007, Illumina Inc. All rights reserved.

9



**SECOND CHEMISTRY CYCLE: DETERMINE
SECOND BASE**

To initiate the next sequencing cycle, add all four labeled reversible terminators and enzyme to the flow cell.

© 2007, Illumina Inc. All rights reserved.

10

After laser excitation, collect the image data as before. Record the identity of the second base for each cluster.

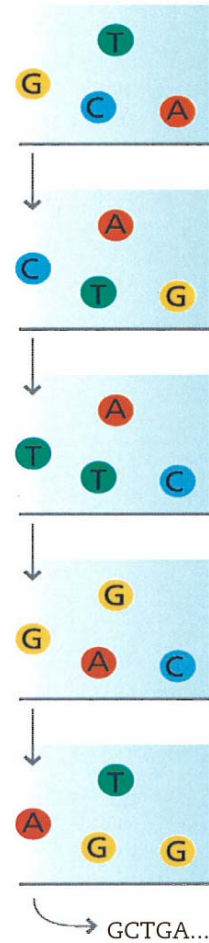


IMAGE OF SECOND CHEMISTRY CYCLE IS CAPTURED BY THE INSTRUMENT.

© 2007, Illumina Inc. All rights reserved.

11

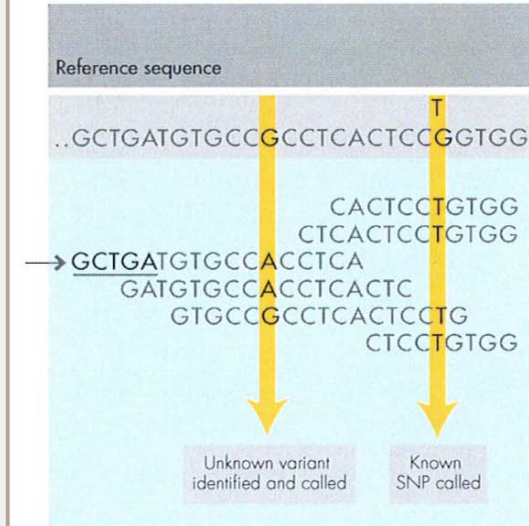
Repeat cycles of sequencing to determine the sequence of bases in a given fragment a single base at a time.



SEQUENCE READS OVER MULTIPLE CHEMISTRY CYCLES

© 2007, Illumina Inc. All rights reserved.

12

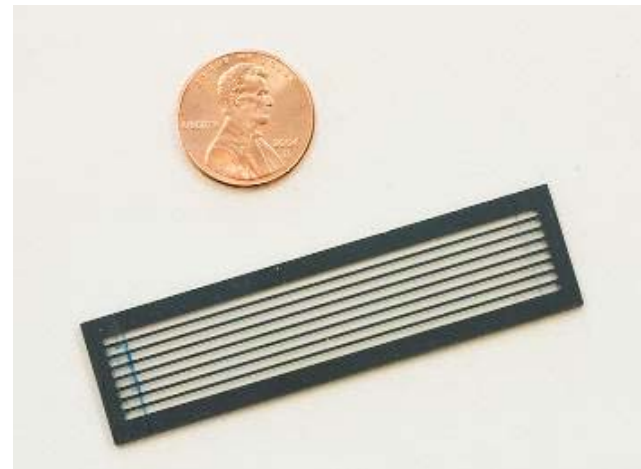


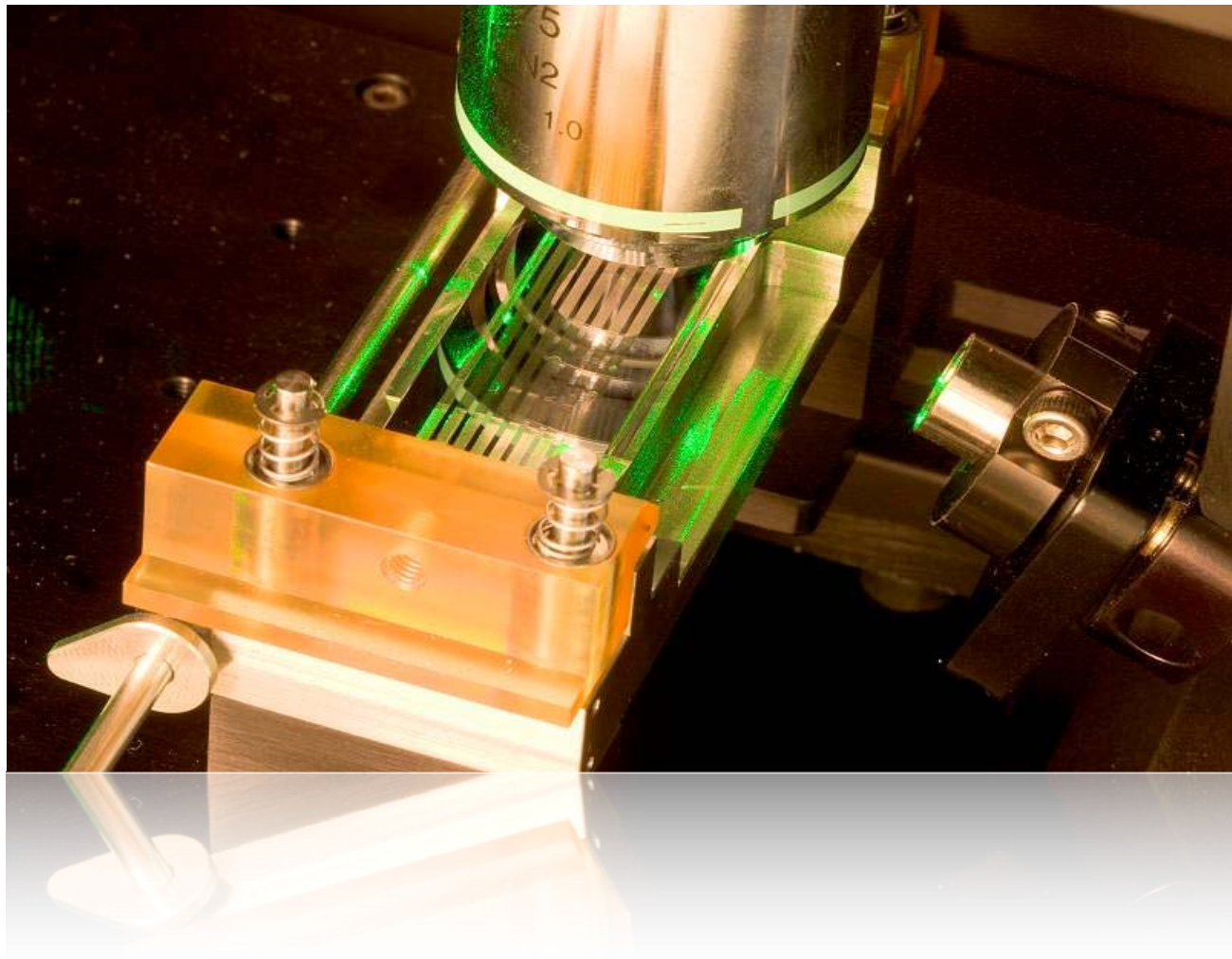
ALIGN THE NEW DATA TO A REFERENCE AND IDENTIFY SEQUENCE DIFFERENCES.

© 2007, Illumina Inc. All rights reserved.

Illumina/Solexa instrument

- Laser-based TIRF* Optics
- 4-colour Detection
- CCD camera
- 8-channel flow cell
- 1 Gb / run at launch





Data acquisition

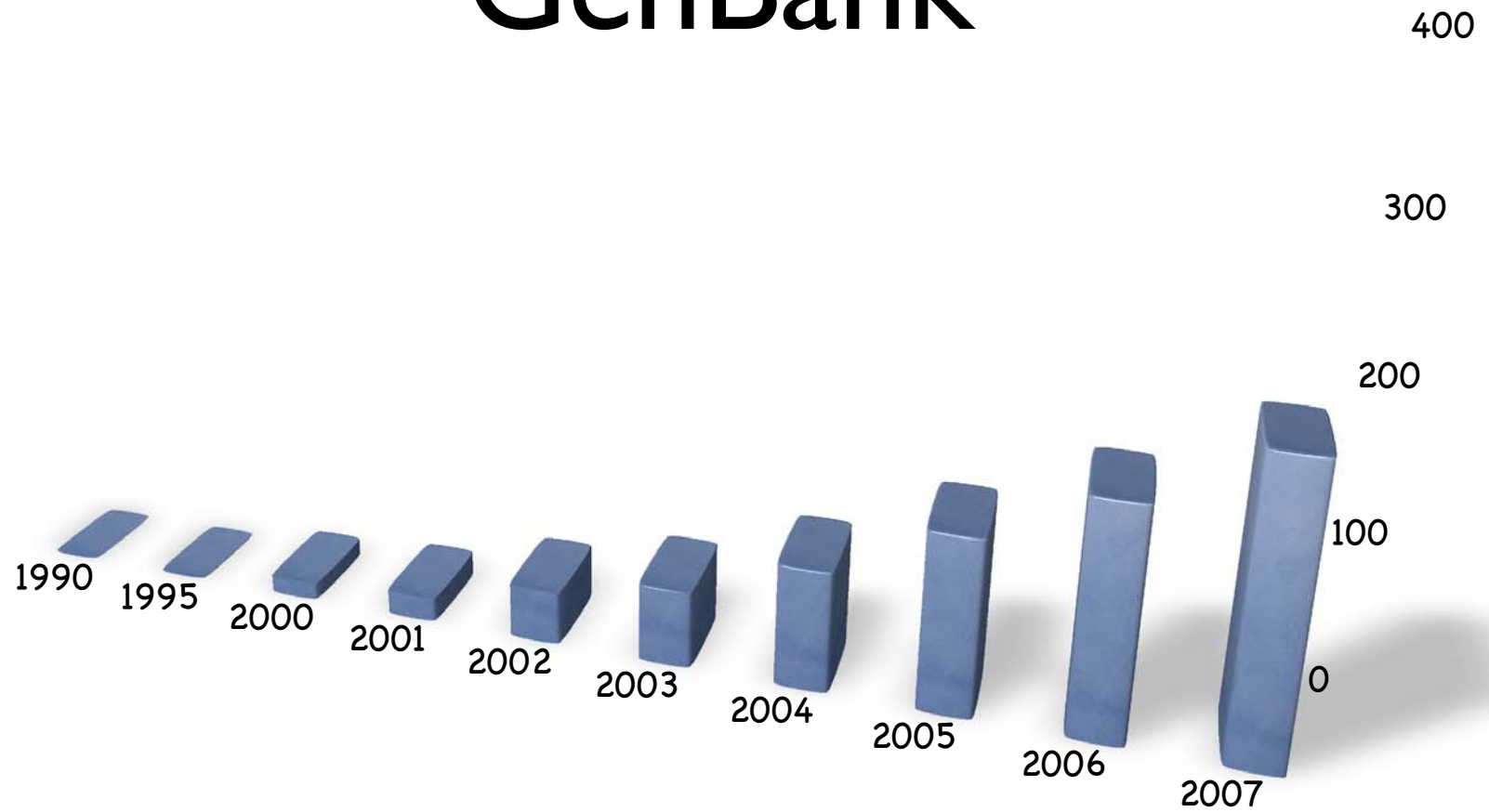
Clonal Single Molecule Array™ Technology

Up to 40 million clusters per flow cell

100 microns

20 microns

GenBank

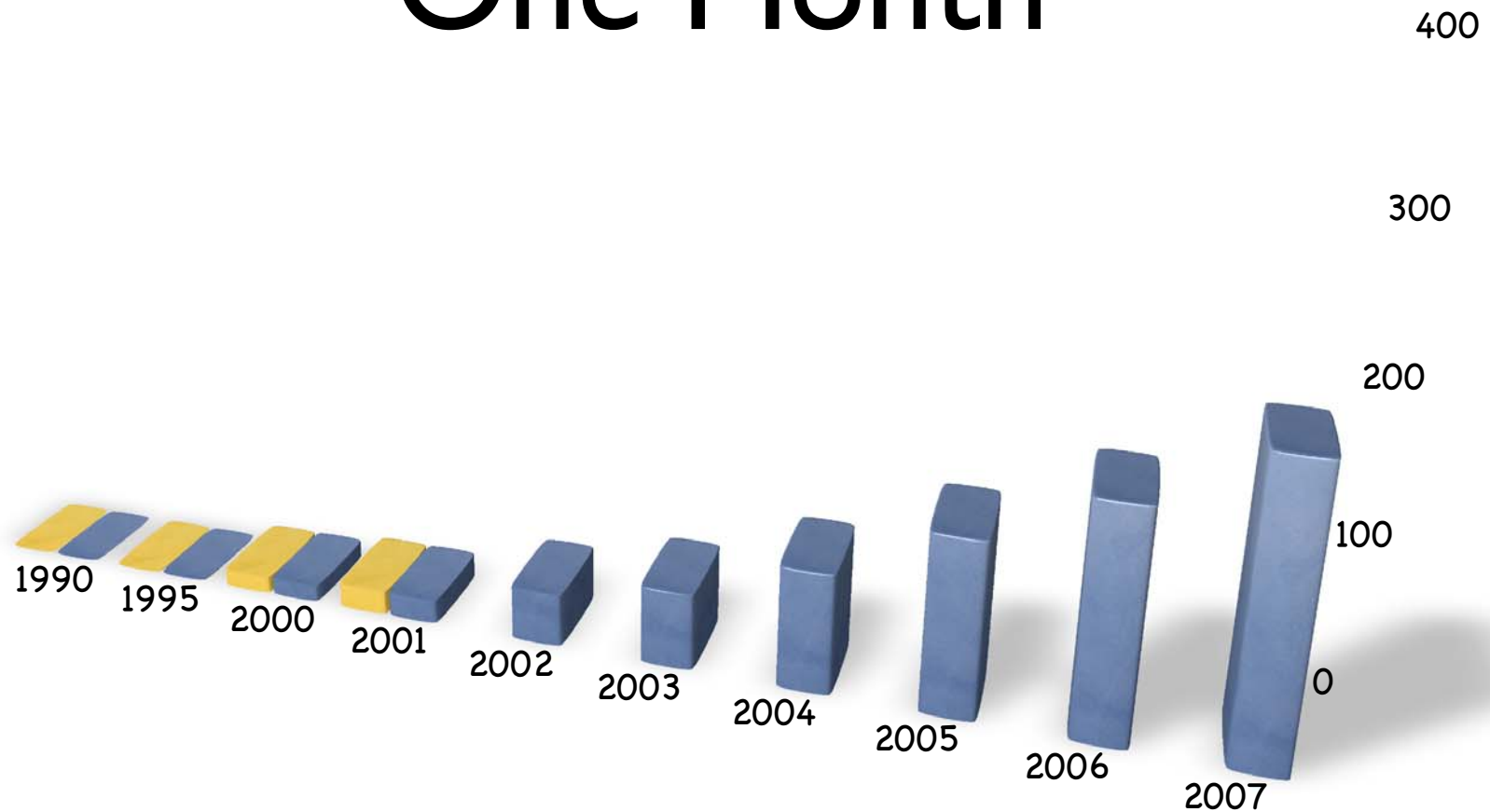


Two Days



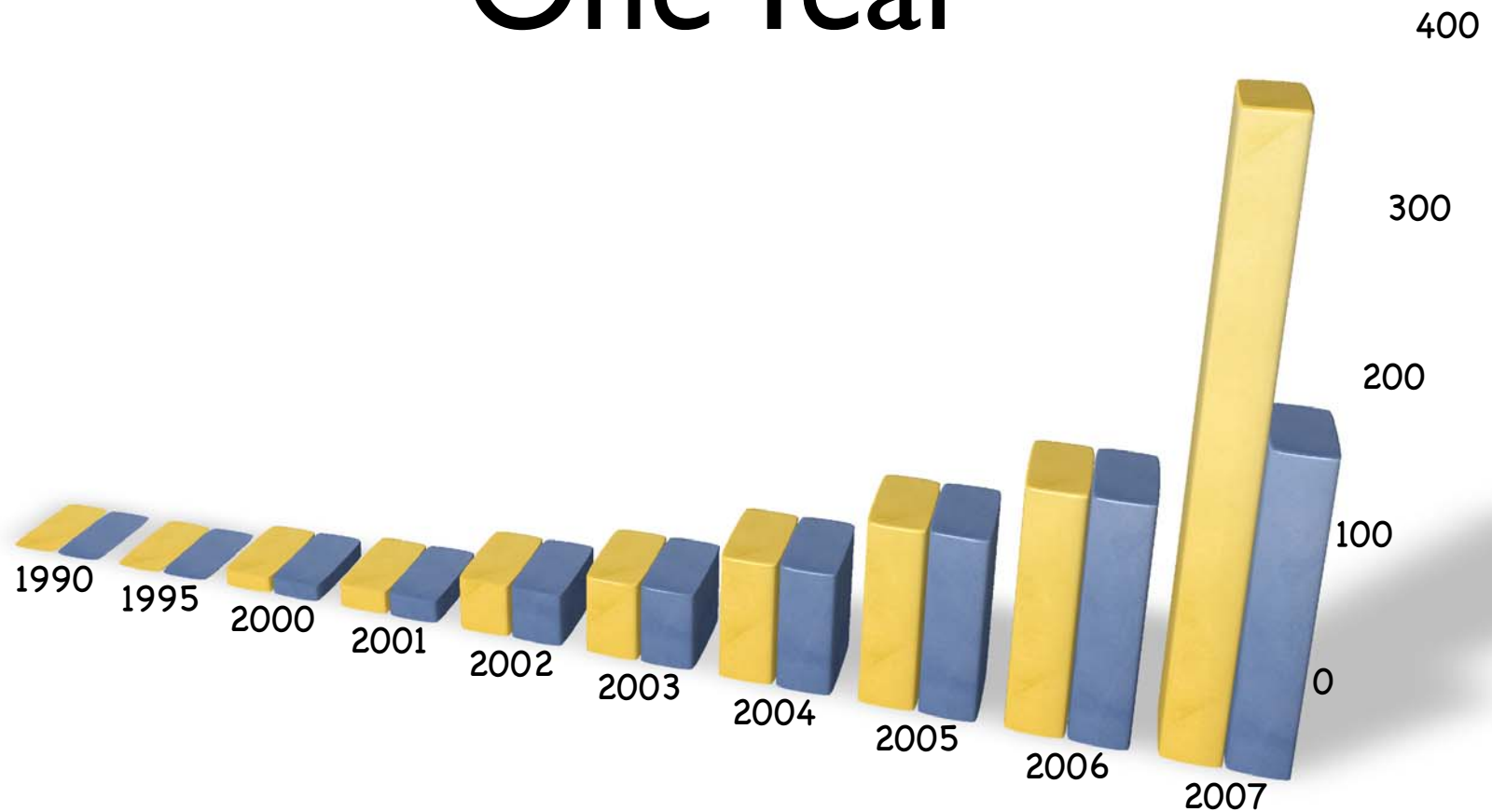
One machine, running for
two days, can generate
~1 Gb data.

One Month

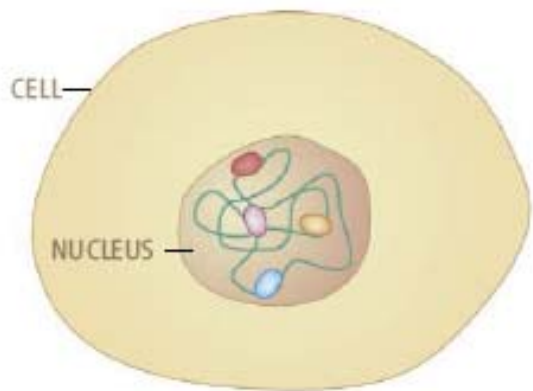


One machine, running for
one month, could generate
~15 Gb data.

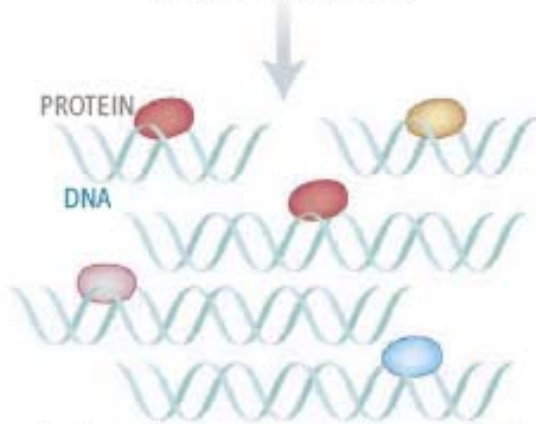
One Year



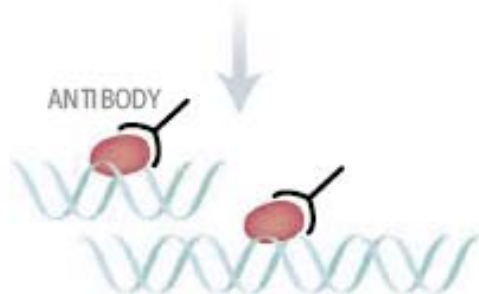
Two machines, running for
one year, could generate
~365 Gb data.



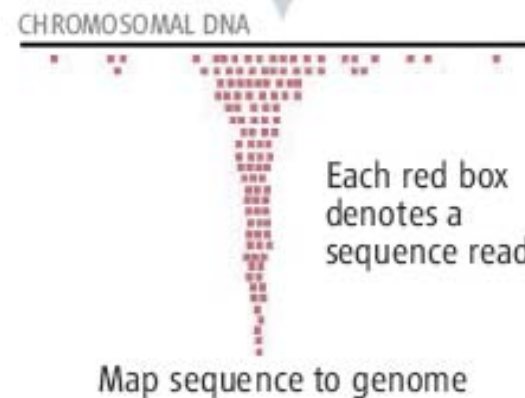
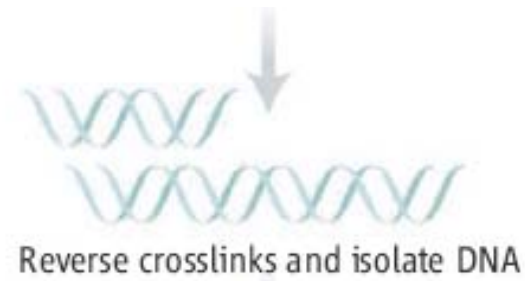
Crosslink proteins to DNA and lyse cells



Isolate chromatin and fragment it



Add a protein-specific antibody and purify protein-DNA complexes



Frontiers in Bioinformatics

- ultra high-throughput sequencing
 - DNA binding site identification
 - genome re-sequencing; SNPs, expression
 - ✓ low cost
 - ✓ whole genome
 - ✓ any genome

Credits & References

- Technology Spotlight on DNA Sequencing with Solexa Technology:

http://www.illumina.com/downloads/SS_DNAsequencing.pdf

- Dr. Steven Jones, GSC

several slides/images used with permission

- Stanley Fields, “Site-Seeing by Sequencing”, Science, 8 June 2007

Sequence Databases

Public Resources at the NCBI





The National Center for Biotechnology Information

NCBI

- **Created in 1988 as a part of the National Library of Medicine at NIH**
- Establish public databases
- Research in computational biology
- Develop software tools for sequence analysis
- Disseminate biomedical information

www.ncbi.nlm.nih.gov

NCBI
National Center for Biotechnology Information
National Library of Medicine National Institutes of Health

PubMed All Databases BLAST OMIM Books TaxBrowser Structure

Search for

SITE MAP
Alphabetical List
Resource Guide

About NCBI
An introduction to NCBI

GenBank
Sequence submission support and software

Literature databases
PubMed, OMIM, Books, and PubMed Central

Molecular databases
Sequences, structures, and taxonomy

Genomic biology
The human

What does NCBI do?
Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease. [More...](#)

Hot Spots

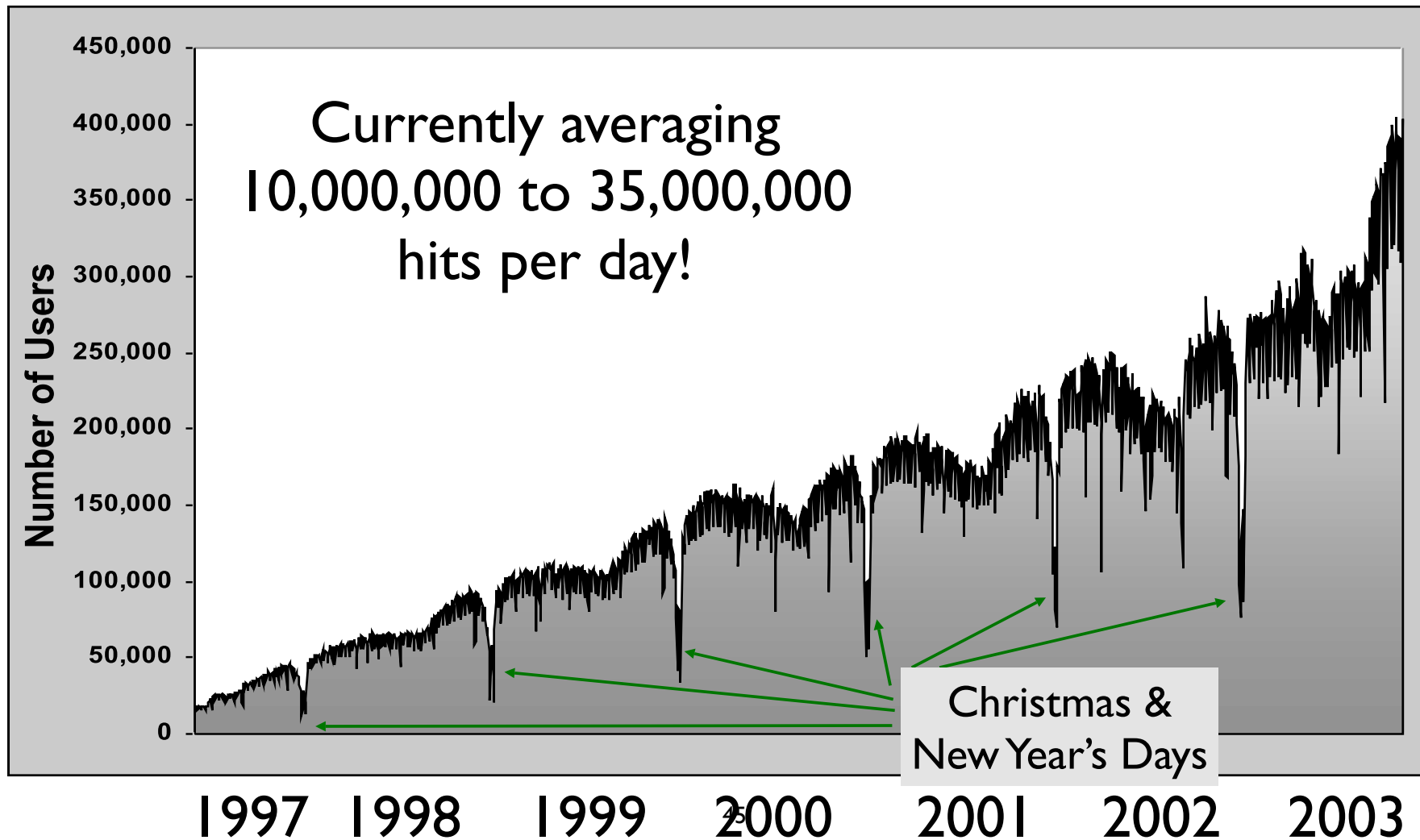
- ▶ Assembly Archive
- ▶ Clusters of orthologous groups
- ▶ Coffee Break, Genes & Disease, NCBI Handbook
- ▶ Electronic PCR
- ▶ Entrez Home
- ▶ Entrez Tools
- ▶ Gene expression omnibus (GEO)
- ▶ Human genome resources
- ▶ Influenza Virus Resource
- ▶ Map Viewer
- ▶ dbMHC

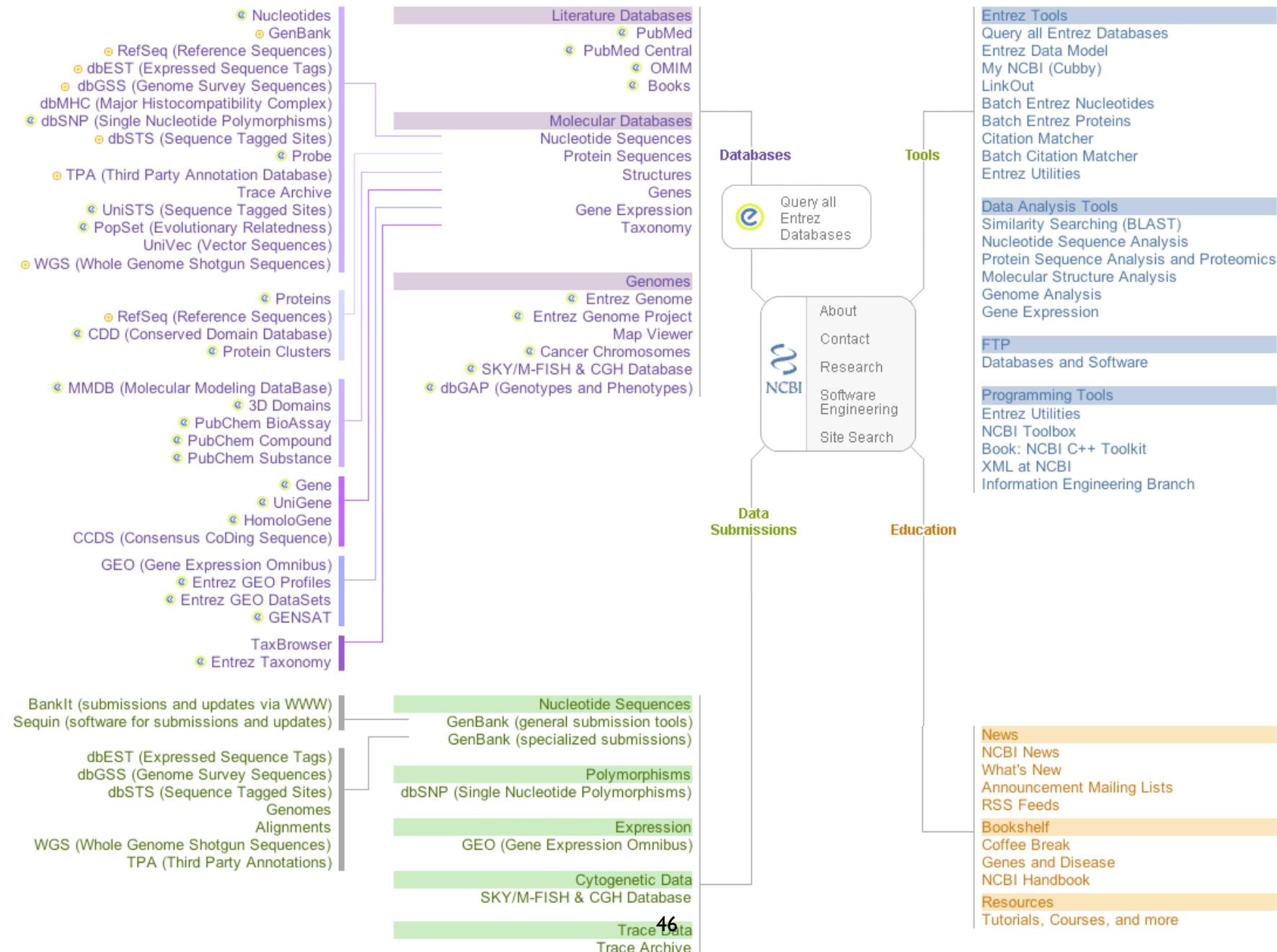
GenBank® Celebrating 25 Years
NCBI will hold a scientific meeting to celebrate the 25th anniversary of GenBank.
April 7-8, 2008
Natcher Auditorium, NIH Campus, Bethesda MD
[click here for more information](#)

New Protein Clusters
Entrez Protein Clusters database
The new Entrez Protein Clusters database is a collection of

44

Number of Users and Hits Per Day





The NCBI ftp site

NCBI **FTP site**

PubMed Entrez BLAST OMIM Books TaxBrowser Structure

Search Entrez for Go

NCBI

SITE MAP
Guide to NCBI resources

About NCBI
The science behind our resources. An introduction for researchers, educators and the public.

GenBank
sequence submission support and software

Molecular databases
sequences, structures and taxonomy

Literature databases
PubMed and OMIM

Genomic Biology

Major resources available by ftp (<ftp.ncbi.nih.gov>):

- ▶ [BLAST Basic Local Alignment Search Tool](#)
Download the BLAST database and stand-alone sequence comparison software.
- ▶ [Cn3D](#)
Download the stand-alone software for viewing 3-dimensional structures.
- ▶ [Data Repository](#)
Download collections of contributed molecular biology data.
- ▶ [GenBank](#)
Download the full release database, daily updates, or WGS files.
Note: there is a mirror site for GenBank files at Indiana University (bio-mirror.net/biomirror/genbank).
- ▶ [Gene](#)
Download gene-based information from completely sequenced genomes.
- ▶ [Genome Assembly/Annotation Projects](#)
Download complete genomes/chromosomes, contigs, mRNAs and proteins.
- ▶ [MMDB](#)

- 30,000 files per day

- 620 Gigabytes per day

NCBI Databases & Services

- GenBank **largest sequence database**
- Free public access to biomedical literature
 - PubMed **free Medline**
 - PubMed Central **full text online access**
- Entrez **integrated molecular & literature databases**
- BLAST **highest volume sequence search service**
- VAST **structure similarity searches**
- Software and Databases

Types of Databases

Primary Databases

- ✓ Original submissions by experimentalists
- ✓ Content controlled by the submitter
- ✓ Examples: GenBank, SNP, GEO

Derivative Databases

- ✓ Built from primary data
- ✓ Content controlled by third party (NCBI)
- ✓ Examples: Refseq, TPA, RefSNP, UniGene, NCBI Protein, Structure, Conserved Domain

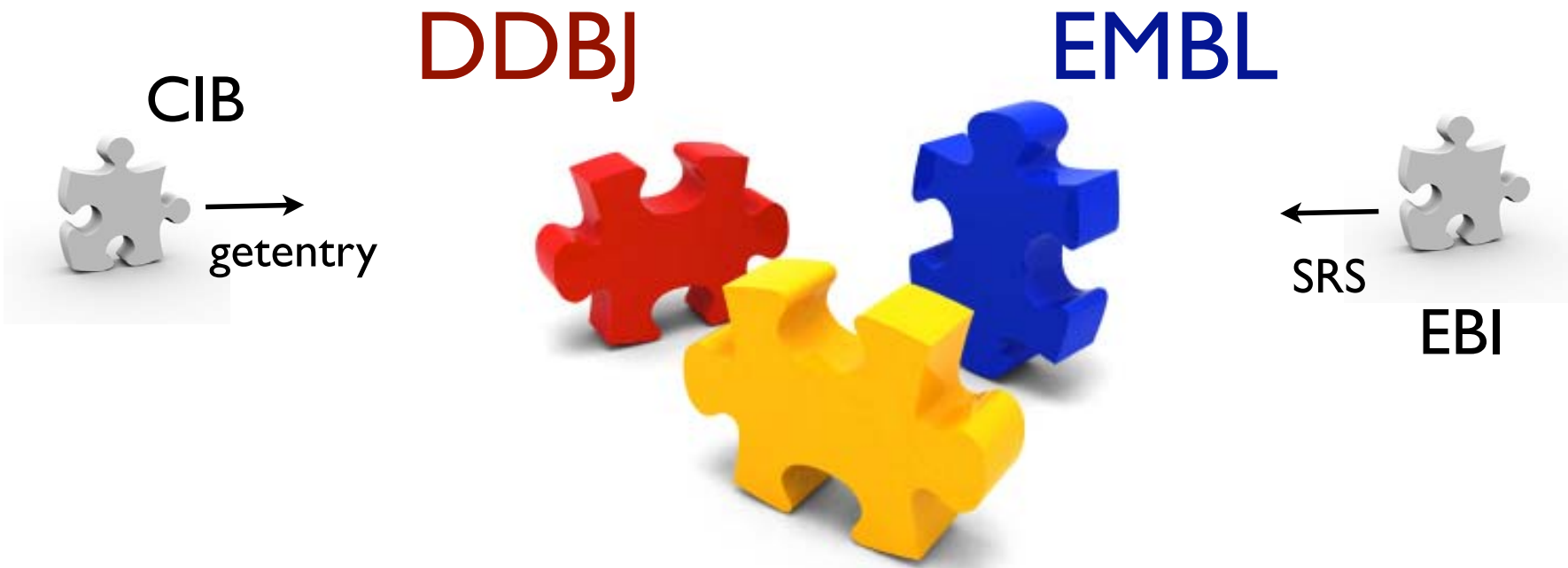
What is GenBank?

NCBI's Primary Sequence Database

- Nucleotide only sequence database
- Archival in nature
- Historical
- Reflective of submitter point of view (subjective)
- Redundant

GenBank Data

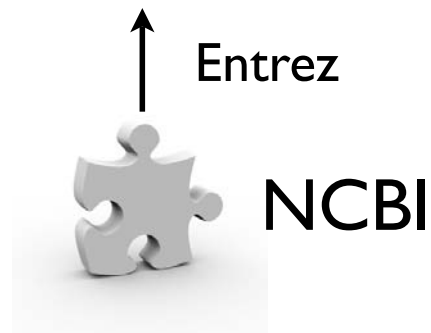
- ✓ Direct submissions (traditional records)
- ✓ Batch submissions (EST, GSS, STS)
- ✓ ftp accounts (genome data)

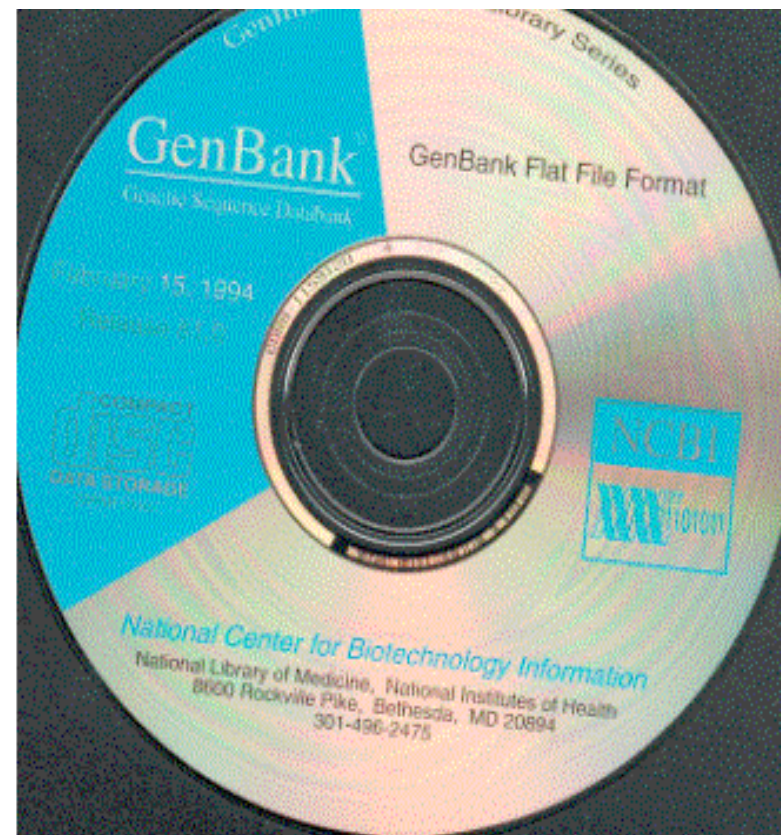
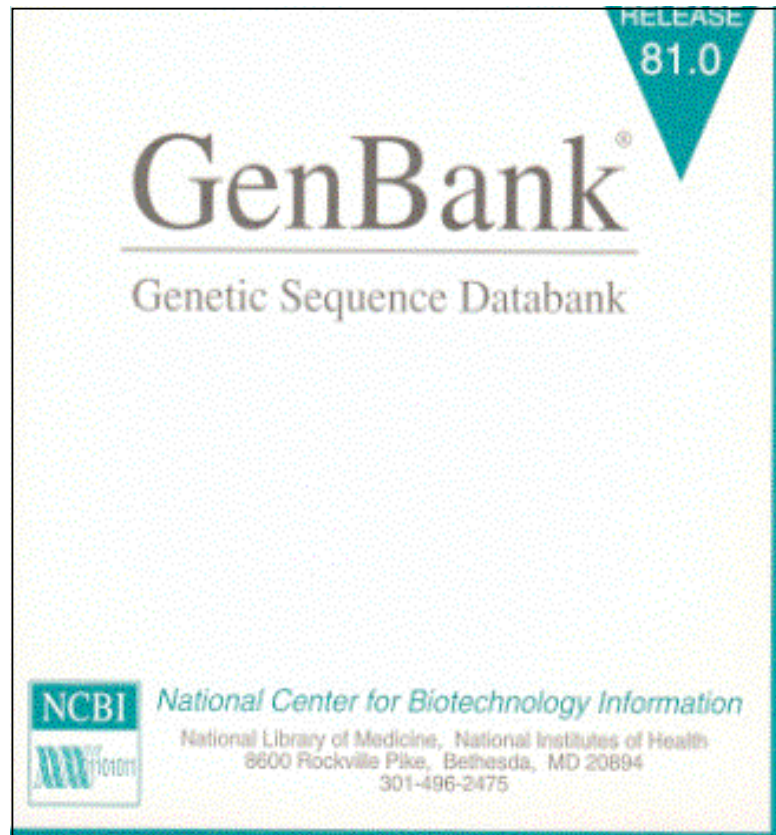


**International
Sequence
Database
Collaboration**

- submit anywhere
- daily updates

GenBank





GenBank: NCBI's Primary Sequence Database

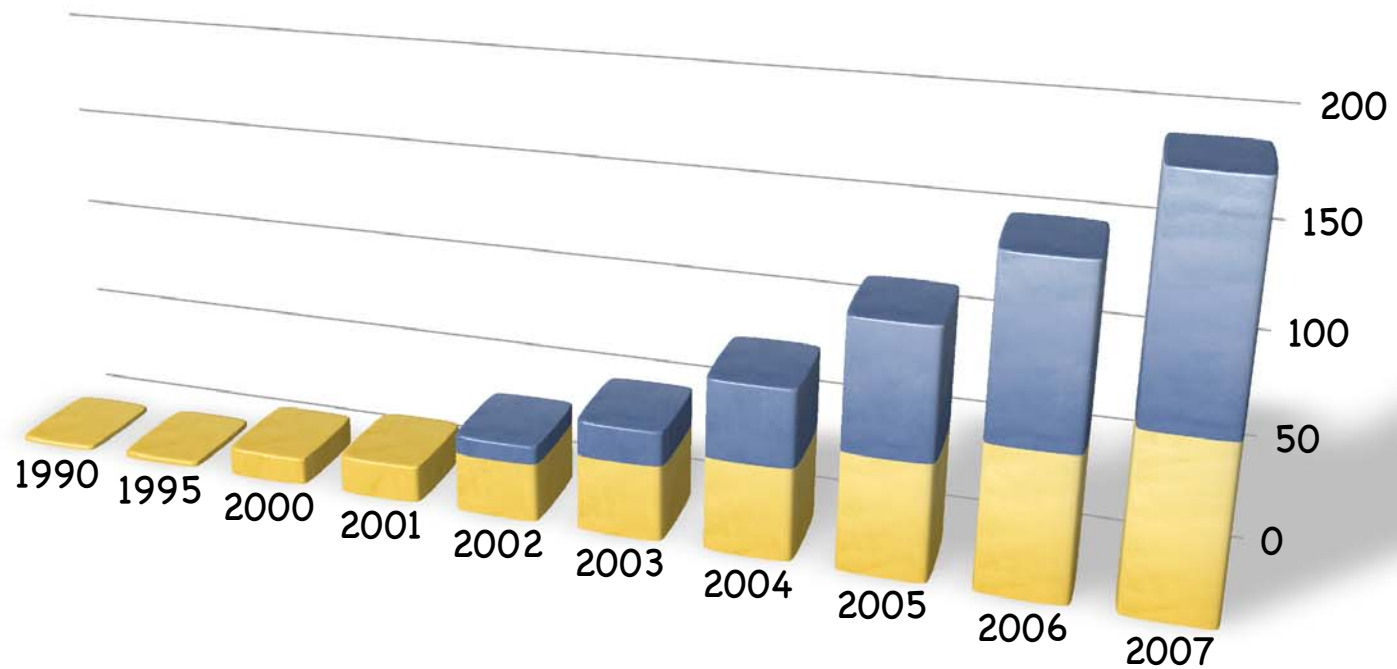
<ftp://ftp.ncbi.nih.gov/genbank/>

Release 161	August 2007
101,530,711	Records
181,489,883,388*	Total Bases

*includes WGS

- full release every two months
- incremental updates daily
- available only via ftp

Growth of GenBank



Current Release 163

Doubling time 12-14 months

GenBank WGS

Organization of GenBank

Records are divided into 18 Divisions.

12 Traditional

6 Bulk

☑ Traditional Divisions:

Direct Submissions (Sequin and BankIt)

Accurate

Well characterized

PRI Primate
PLN Plant and Fungal
BCT Bacterial and Archeal
INV Invertebrate
ROD Rodent
VRL Viral
VRT Other Vertebrate
MAM Mammalian
PHG Phage
SYN Synthetic (cloning vectors)
ENV Environmental Samples
UNA Unannotated

Organization of GenBank

Records are divided into 18 Divisions.

12 Traditional

6 Bulk

BULK Divisions:

Batch Submission (Email and FTP)

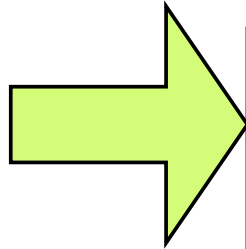
Inaccurate

Poorly characterized

EST Expressed Sequence Tag
GSS Genome Survey Sequence
HTG High Throughput Genomic
STS Sequence Tagged Site
HTC High Throughput cDNA
PAT Patent

Entrez query: `gbdiv_xxx[Properties]`

Traditional GenBank Record



```
LOCUS       HSHMLHI                2503 bp    mRNA    linear    PRI 31-MAR-1994
DEFINITION  Human DNA mismatch repair (hmlh1) mRNA, complete cds.
ACCESSION   U07418
VERSION     U07418.1  GI:466461
KEYWORDS    .
SOURCE      Homo sapiens (human)
  ORGANISM  Homo sapiens
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Mammalia;
            Mammalia; Eutheria; Primates; Catarrhini; Hominidae; Homo;
REFERENCE   1 (bases 1 to 2503)
AUTHORS     Papadopoulos,N., Nicolaides,N.C., Wei,P., Manolagas,S.C., Lippman,M.E.,
            Carter,K.C., Rosen,C.A., Haseltine,W.A., Kalish,J.M.,
            Fraser,C.M., Adams,M.D., Venter,J.C., Wilson,R.,
            Watson,P., Lynch,H.T., Peltomaki,P., O'Connell,M.P.,
            Kinzler,K.W. and Vogelstein,B.
TITLE       Mutation of a mutL homolog in hereditary non-polyposis
            colorectal cancer
JOURNAL     Science 263 (5153), 1625-1629 (1994)
MEDLINE    94174288
```

Accession

- Stable
- Reportable
- Universal

ACCESSION **U07418**

VERSION **U07418.1** **GI:466461**

Version

- Tracks changes in sequence

GI number

- NCBI internal use

```

FEATURES             Location/Qualifiers
     source           1..2503
                     /organism="Homo sapiens"
                     /db_xref="taxon:9606"
                     /chromosome="3"
                     /map="p21"
                     /tissue_type="gall bladder"
                     /dev_stage="adult"
     gene             1..2503
                     /gene="hmlh1"
     CDS              42..2312
                     /gene="hmlh1"
                     /function="DNA mismatch repair"
                     /note="human homolog of E. coli mutL gene product,
                     Swiss-Prot Accession Number P23367"
                     /codon_start=1
                     /protein_id="AAA17374.1"
                     /db_xref="GI:466462"
                     /translation="MSFVAGVIRRLDET VVNRIAAGEVIQR PANAIKEMIENCLDAKS
                     TSIQIVKBEGLKLIQIQDNGTGIRKEDLDIVCERFTTSK LQSFEDLASISTYGRGE
                     ALASISHVAHVTTITTKTADGKCA YRASYS DGLKLPKPPKPCAGNQGTQITVEDLFY NIA
                     TRRKALKNPSE EYGKILEVVGRYSVHNAGISF SVKKQGETVADVRTL PNASTVDNIRS
                     VFGNAVSRELIEIGCEDKTLAFKMNGYI SNANYSVKKCI FLLFINHRLVESTSLRKAI
                     ETVYAAYLPKNTHFFLYLSLEIS PQNV DVNVHPTKHEVHFLHEESILERVQQHIESKL
                     LGSNSRSMYFTQTLLPGLAGPSGEMVKSTTSLTSSSTSGSSDKVYAHQMVRTDSREQK
                     LDAFLQPLSKPLSQPQAI VTEDKTDISSGRARQDEEMLELPAPAEVAAKNQSL EGD
                     TTKGTSEMSEKRGPTSSNPRKRHREDS DVEMVEDDSRKEMTA ACTPRRRIINLTSVLS
                     LQEBINEQGHEVLRMLHNHSFVGC VNPQWALAQHQTKLYLLNTTKLSEELFYQIL IY
                     DFANFVGLRLSE PAFPLDLAMLALDS PEGSWTEEDGPKGLAEYIV EFLKKAEMLAD
                     YFSLEIDBEGNLIGLPLLDNYVPPLEGLPIFILRLATEVNWDEEKECFE SLSKECAM
                     FYSIRKQYISEESTLSGQSEVPGSIPNSWKWTV EHVIVYKALRSHILPPKHFTEDGNI
                     LOLANLPDLYKVFERC"

```

well annotated

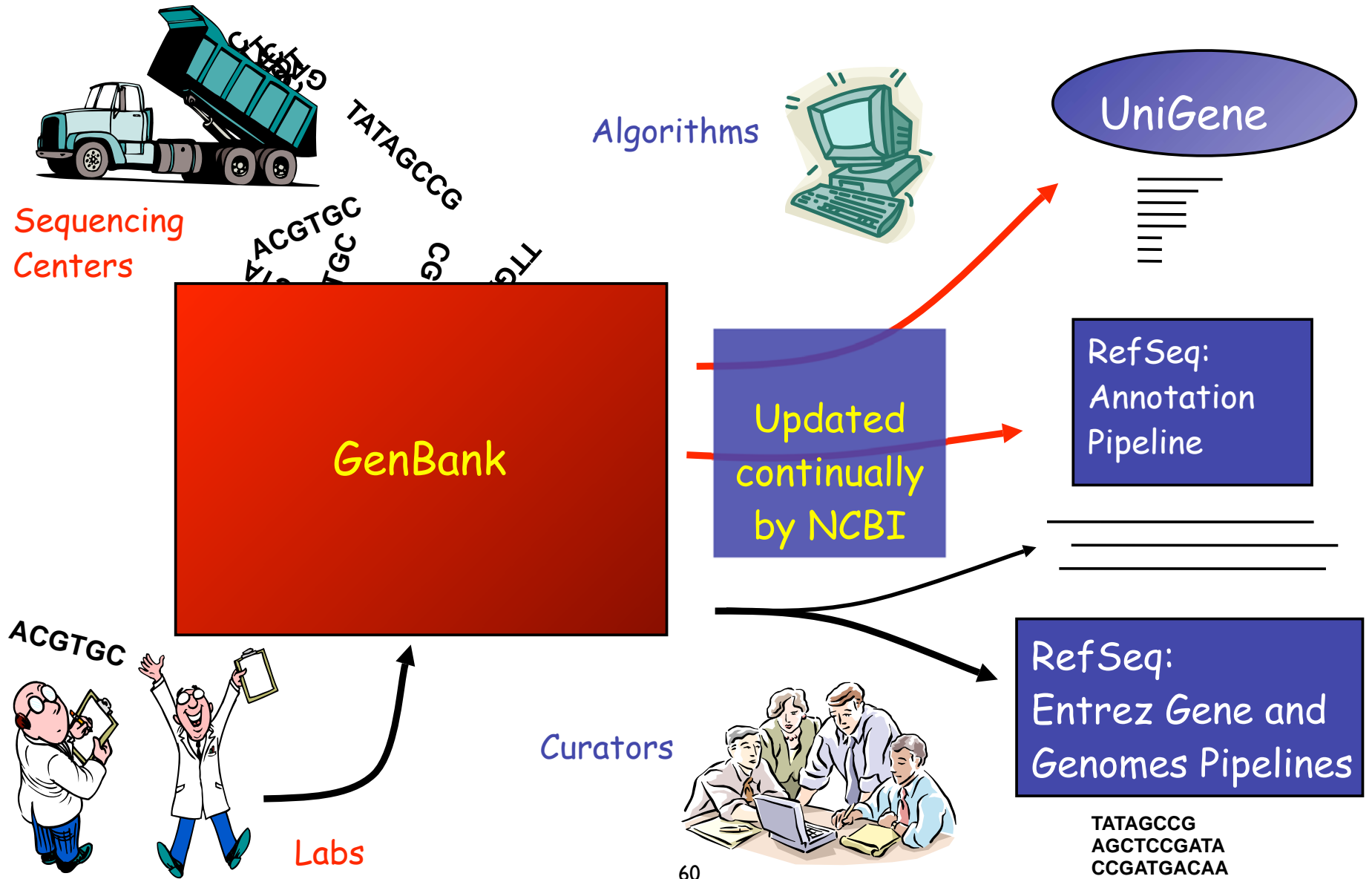
```

BASE COUNT      723 a   539 c   599 g   642 t
ORIGIN
1  gttgaacatc tagacgtttc cttggtcttt ctggcgccaa aatgctgttc gtggcagggg
61  ttattcggcg ctggacgag acagtgttga accgcatcgc ggcgggggaa gttatccagg
121  ggccagctaa tgctatcaaa gagatgattg agaactgttt agatgcaaaa tccacaagta
181  ttcaagtgat tgtaaaagag ggaggcctga agttgattca gatccaagac aatggcaccg
241  ggatcagгаа agaagatctg gatattgtat gtgaaagggt cactactagt aaactgcagt
301  cctttgagga tttagccagt atttctacct atggctttcg aggtgaggct ttggccagca
361  taagccatgt ggctcatggt actattacaa cgaaaacagc tgatgгааag tgtgcataca
421  gagcaagtta ctccagatgga aaactgaaag cccctcctaa accatgtgct ggcactcaag
481  ggaccagat cacggtggag gacctttttt acaacatagc cacgaggaga aaagctttaa
541  aaaatccaag tgaagaatat gggaaaattt tggaaagtgt tggcaggtat tcagtacaca
601  atgcaggcat tagtttctca gttaaaaaac aaggagagac agtagctgat gttaggacac
661  tacccaatgc ctcaaccgtg gacaatattc gctcctctt tggaaatgct gttagtcgag
721  aactgataga aattggatgt gaggataaaa ccctagcctt caaaatgat ggttaccatg
781  ccaatgcaaa ctactcagtg aagaagtgca tcttctact cttcatcaac catcgtctgg
841  tagaatcaac ttccttgaga aaagccatag aaacagtgta tgcagcctat ttgccccaaa
901  acacacaccc attcctgtac ctccagttag aaatcagttc ccagaatgtg gatgtaaatg
961  tgcaccccc aaagcatgaa gttcacttcc tgcacgagga gagcatcctg gacggggtgc
1021  agcagcacat cgagagcaag ctctctgggt ccaattctcc caggatgtac ttcaccagca
1081  ctttgctacc aggacttgcg ggcccctctg gggagatggt taaatccaca acaagctgaa
1141  cctcgtcttc tacttctgga agtagtgata aggtctatgc ccaccagatg gttcgtacag
1201  attcccggga acagaagctt gatgcatttc tgcagcctct gagcaaaccc ctgtccagtc
1261  agccccaggc cattgtcaca gaggataaga cagatatttc tagtggcagg gtaggcagc
1321  aagatgagga gatgctttaa ctcccagccc ctgctgaagt ggctgccaaa aatcagagct
1381  tggaggggga tacaacaagg gggacttcag aaatgtcaga gaagagagga cctactcca
1441  gcaacccag aaagagacat cgggaagatt ctgatgtgga aatggtgгаа actgattccc
1501  gaaaggaat gactgcagct tgtaccccc ggagaaggat cattaacctc actagtgttt
1561  tgagtctcca ggaagaaatt aatgagcagg gacatgaggt tctccgggag atgttgcata
1621  accactcctt cgtgggctgt gtgaatcctc agtgggcctt ggcacagcat caaaccaagt
1681  tataccttct caacaccacc aagcttagtg aagaactggt ctaccagata ctctttatg
1741  attttgccaa ttttggtggt ctccagttat cggagccagc accgctcttt gaccttgcca
1801  tgcttgctt agatagtcca gagagtggct ggacagagga agatggtccc aaagaaggad
1861  ttgctgaata cattggtgag tttctgaaga agaaggctga gatgcttcca gactatttct
1921  ctttggaaat tgatgaggaa gggaacctga ttggattacc ccttctgat gacaactatg
1981  tgccccctt ggagggactg cctatcttca ttccttgact agcccactgag gtgaattggg
2041  acgaagaaaa ggaatgtttt gaaagcctca gtaagaatg cgctatgttc tattccatcc
2101  ggaagcagta catatctgag gactcgaccc tctcaggcca gcagagtгаа gtgcctggct
2161  ccattccaaa ctccctggaag tggactgtgg aacacattgt ctataaagcc ttgcgctcad
2221  acattctgcc tcttaaacat ttcacagaag atggaaatct cctgcagctt gctaacctgc
2281  ctgatctata caaagtcttt gagaggtggt aaatatggtt atttatgcac tgtgggatgt
2341  gttctctctt ctctgtattc cgatacaaaг tgttgatca aagtgtgata tacaaagtgt
2401  accaacataa gtgttggtag cacttaagac ttatacttgc cttctgatag tattccttta
2461  tacacagtgг attgattata aataaataga tgtgtcttaa cat

```

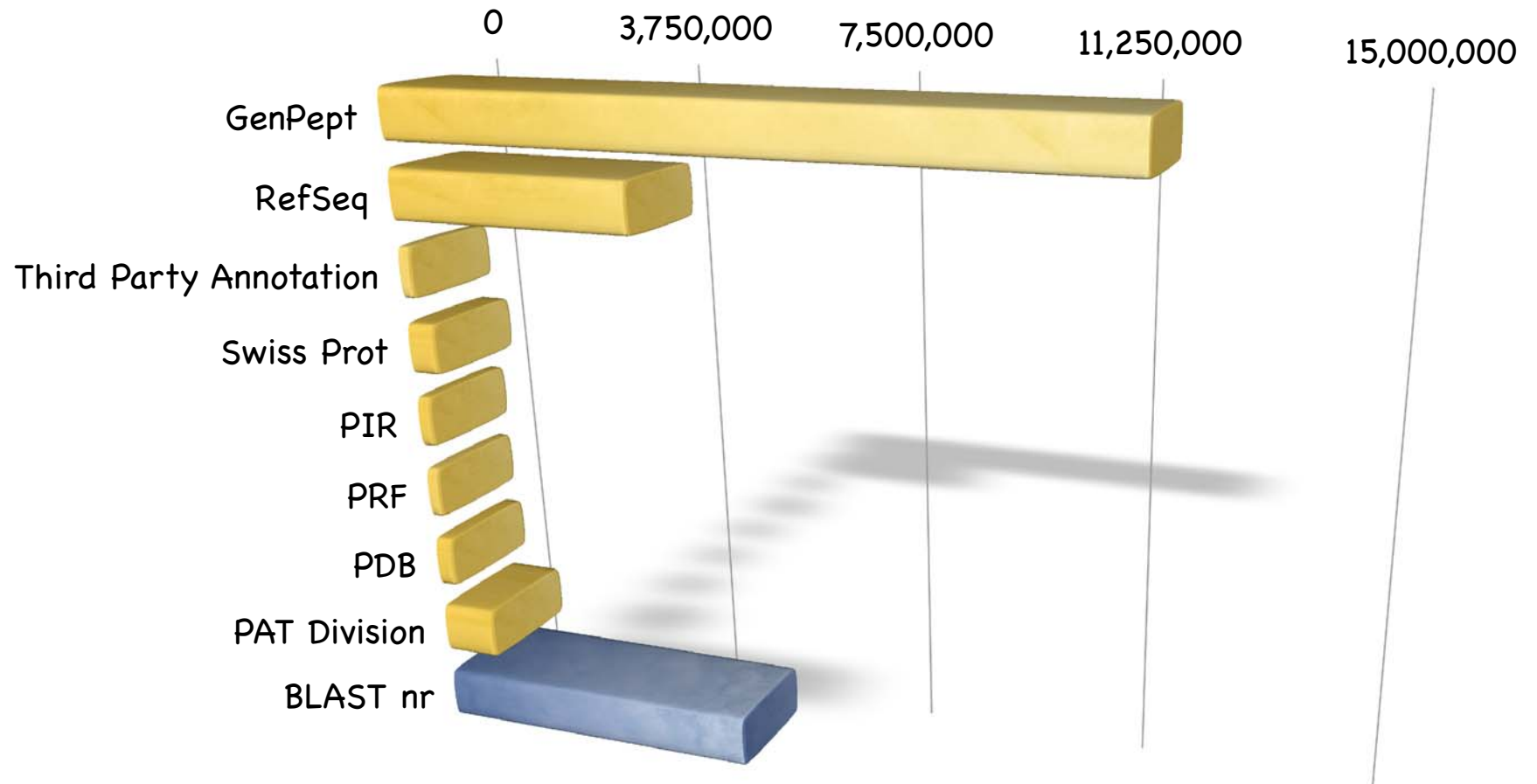
the sequence is the data

Primary vs. Derivative Databases



Derivative Databases

Entrez Protein



GenPept

- GenBank CDS translations

```
FEATURES             Location/Qualifiers
     source            1..2484
                        /organism="Homo sapiens"
                        /mol_type="mRNA"
                        /db_xref="taxon:9606"
                        /chromosome="3"
                        /map="3p22-p23"
     gene              1..2484
                        /gene="MLH1"
     CDS                22..2292
                        /gene="MLH1"
                        /note="homologous to DNA mismatch repair protein homolog (GenBank Accession
                        Number P14242), S. cerevisiae MLH1 (GenBank Accession
                        Number U07187), E. coli MUTL (Swiss-Prot Accession Number
                        P23367), Salmonella typhimurium MUTL (Swiss-Prot Accession
                        Number P14161) and Streptococcus pneumoniae (Swiss-Prot
                        Accession Number P14161)"
                        /codon_start=1
                        /product="DNA mismatch repair protein homolog"
                        /protein_id="AAC50285.1"
                        /db_xref="GI:463989"
                        /translation="MSFVAGVIRRLDETQVNRVIAAGEVIQRPANAIKEMIEENCLDAKSTSIQVIV...
                        TSIQVIVKEGGLKLIQIQDNG...RKEDLDIVCERFTTSKLQSFEDLASISTYGFRGE
                        ALASISHVAHVTTITTKTADGK...RASYSYDGLKAPPKPCAGNQGTTITVEDLFYNIA
                        TRRKALKNPSEEYGKILEVVGGRYSVHNAGISFSVKKQGETVADVRTLPNASTVDNIRS"
```

>gi|463989|gb|AAC50285.1| DNA mismatch repair prote...
MSFVAGVIRRLDETQVNRVIAAGEVIQRPANAIKEMIEENCLDAKSTSIQVIV...
EDLDIVCERFTTSKLQSFEDLASISTYGFRGEALASISHVAHVTTITTKTAD...



RefSeq

- The goal is to provide a comprehensive, standard dataset that represents sequence information for a species.
 - transcript, protein, assembled genomic (contigs), and chromosome records
 - known and predicted
 - reviewed
 - human, mouse, rat, fruit fly, zebrafish, arabidopsis, microbial genome (proteins), organelles, and more

RefSeq Accession Numbers

- mRNAs and Proteins

NM_123456	Curated mRNA
NP_123456	Curated Protein
NR_123456	Curated nc RNA
XM_123456	Predicted mRNA
XP_123456	Predicted Protein
XR_123456	Predicted nc RNA

- Genomic Records

NG_123456	Reference Genomic Sequence
-----------	----------------------------

- Chromosome

NC_123455	Microbial replicons, organelle, genomes, human chromosomes
-----------	--

- Assemblies

NT_123456	Contig
NW_123456	WGS Supercontig

Other NCBI Databases

Structure:	imported structures (PDB)	Cn3D viewer, NCBI curation
CDD:	conserved domain database	Protein families (COGs and KOGs); Single domains (PFAM, SMART, CD)
dbSNP:	nucleotide polymorphism	
Gene:	gene records	Unifies LocusLink and Microbial Genomes
HomoloGene:	homologs	neighboring function for Gene



Sequence Databases

GUIDED TOUR: Retrieving Data



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

The image shows a screenshot of the NCBI website. At the top, the URL <http://www.ncbi.nlm.nih.gov/> is displayed. Below it is the NCBI logo and the text "National Center for Biotechnology Information", "National Library of Medicine", and "National Institutes of Health". A navigation bar contains links for "PubMed", "All Databases", "BLAST", "OMIM", "Books", "TaxBrowser", and "Structure". A search bar is present with a dropdown menu set to "All Databases" and a "Go" button. On the left, a "SITE MAP" sidebar lists categories like "Alphabetical Resource Guide", "About NCE", "GenBank", "Literature databases", "Molecular databases", "Genomic biology", and "Tools". The main content area features a "What do you want to do?" section with a list of services: Assembly Archive, Clusters of orthologous groups, Coffee Break, Genes & Disease, NCBI Handbook, Electronic PCR, Entrez Home, Entrez Tools, Gene expression omnibus (GEO), Human genome resources, Influenza Virus Resource, Map Viewer, dbMHC, Mouse genome resources, and M69NCBI. A "Hot Spots" section is also visible. A blue callout box titled "Celebrating 25 Years" mentions a scientific meeting. A yellow callout box titled "New Protein Clusters" describes the Entrez Protein Clusters database. On the right side of the page, three horizontal double-headed arrows point to the text "WWW", "Entrez", and "BLAST" respectively.



WWW

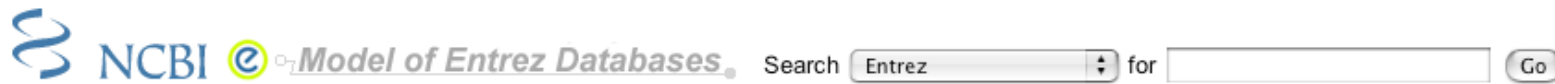


Entrez



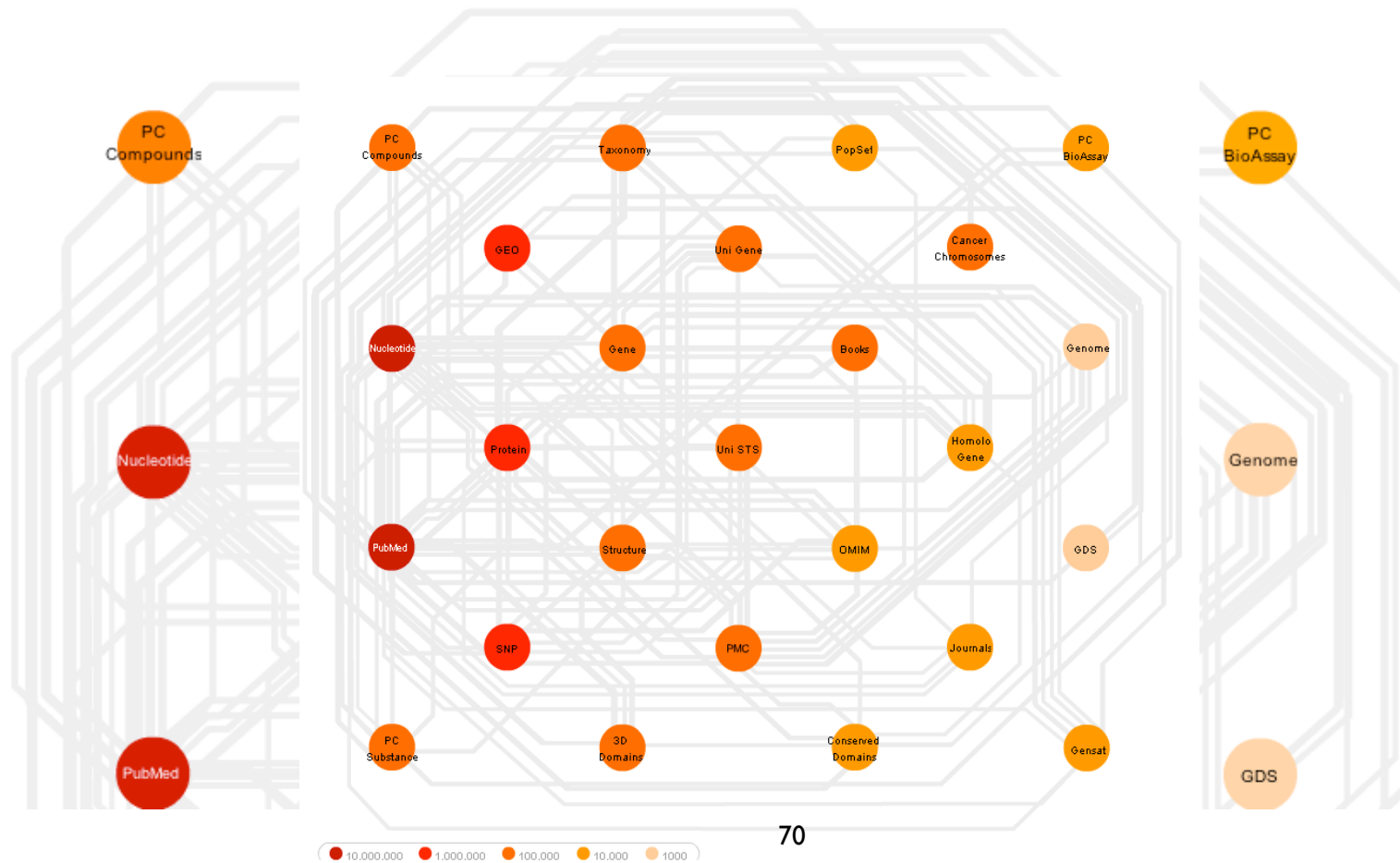
BLAST

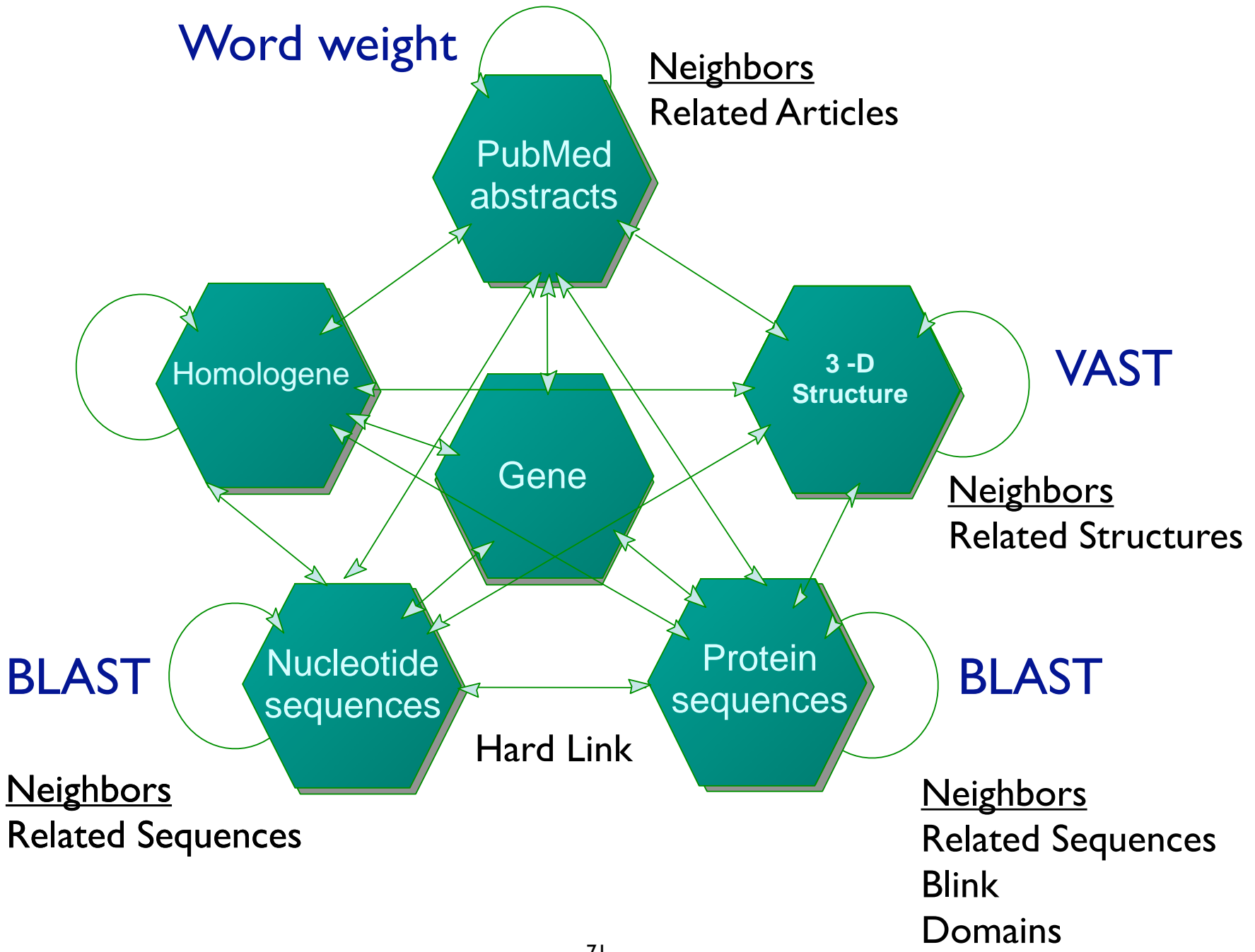
<http://www.ncbi.nih.gov/Database/datamodel>



The diagram shows the Entrez databases and the connections between them. Each database is represented by a colored circle, where the color indicates the approximate number of records in the database. Mouse over a circle to see which databases are linked to the one selected, and how many links exist between those databases.

This diagram requires [Flash](#) for viewing.





Neighbors in Entrez

1: [rs709932](#) [*Homo sapiens*]CGAP-GAI, ILLUMINA, ILLUMINA, ILLUMINA, ILLUMINA, LEE, TSC-CSHI Links SNP

1: [GDS596 record](#) | [GPL96 211298_s_at](#) [*Homo sapiens*] 158 samples Profile Neighbors, Sequence Neighbors, Links

Annotation: [ALB](#): albumin DKFZp779N1935, PRO0883, PRO0903, PRO1341 GEO

Reporter: [AF116645](#)

Exper 1: [MLH1](#) Order cDNA clone, Links

Official Symbol: MLH1 **and Name:** mutL homolog 1, colon cancer, nonpolyposis type 2 (*E. coli*) [*Homo sapiens*]

Other Aliases: COCA2, FCC2, HNPCC, HNPCC2, MGC5172, hMLH1

Other Designations: DNA mismatch repair protein Mlh1; MutL protein homolog 1

Location: 3p21.3 Gene

1: [Plotz G, Welsch C, Giron-Monzon L, Friedhoff P, Albrecht M, Piiper A, E S, Raedle J.](#) PubMed Related Articles, Links

1: [NP_000240](#). Reports MutL protein homo...[gi:4557757] BLink, Conserved Domains, Links

[Comment](#) [Features](#) [Sequence](#)

LOCUS	NP_000240	756 aa	linear	PRI 22-APR-2007
DEFINITION	MutL protein homolog 1 [<i>Homo sapiens</i>].			
ACCESSION	NP_000240			
VERSION	NP_000240.1	GI:4557757		
DBSOURCE	REFSEQ: accession NM_000249.2			

Protein

Blink & Domains

Neighbors: BLAST Link
pre-computed BLAST

1: [NP_000240](#). Reports MutL protein homo...[gi:4557757]

BLink, Conserved
Domains, Links

[Comment](#) [Features](#) [Sequence](#)

LOCUS NP_000240
DEFINITION MutL protein homolog 1 [
ACCESSION NP_000240
VERSION NP_000240.1 GI:4557757
DBSOURCE REFSEQ: accession [NM_000249.2](#)

Neighbors:
pre-computed CDD search

APR-2007

Links

1: [NP_000240](#). Reports MutL protein homo...[gi:4557757]

[Comment](#) [Features](#) [Sequence](#)

LOCUS	NP_000240	56 aa
DEFINITION	MutL prot	[no sapiens].
ACCESSION	NP_000240	
VERSION	NP_000240.1	GI:4557757
DBSOURCE	REFSEQ: accession	NM_000249.2

Neighbors

Links

- ▶ Gene
- ▶ Genome Project
- ▶ HomoloGene
- ▶ PubMed (RefSeq)
- ▶ Gene Genotype
- ▶ GeneView in dbSNP
- ▶ Related Structure
- ▶ UniGene
- ▶ Related Sequences
- ▶ Domain Relatives
- ▶ Genome
- ▶ Map Viewer
- ▶ Nucleotide
- ▶ OMIM
- ▶ PubMed
- ▶ SNP
- ▶ Taxonomy
- ▶ LinkOut

Hard Links

Database searching with Entrez

- **Scenario:** Let's find out more about the MutL gene
- Using limits and field restriction to find human MutL homolog
- Linking and neighboring with MutL



Start with a search for “colon cancer”

The screenshot shows the NCBI homepage with a search bar containing 'colon cancer' and a 'Go' button. The search results are displayed in a list format. The first result is 'What does NCBI do?' with a description of the center's mission. The second result is 'Hot Spots' with a list of featured resources. A banner for 'GenBank Celebrating 25 Years' is also visible.

NCBI
National Center for Biotechnology Information
National Library of Medicine National Institutes of Health

PubMed All Databases BLAST OMIM Books TaxBrowser Structure

Search All Databases for colon cancer Go

SITE MAP
Alphabetical List
Resource Guide

About NCBI
An introduction to
NCBI

GenBank
Sequence
submission support
and software

**Literature
databases**
PubMed, OMIM,
Books, and PubMed

What does NCBI do?

Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease. [More...](#)

Hot Spots

- ▶ Assembly Archive
- ▶ Clusters of orthologous groups
- ▶ Coffee Break, Genes & Disease, NCBI Handbook
- ▶ Electronic PCR
- ▶ Entrez Home
- ▶ Entrez Tools

GenBank® Celebrating 25 Years

NCBI will hold a scientific meeting to celebrate the 25th anniversary of GenBank.

Search across databases

[Help](#)

- Result counts displayed in gray indicate one or more terms not found

58219		PubMed: biomedical literature citations and abstracts	
7197		PubMed Central: free, full text journal articles	
7		Site Search: NCBI web and FTP sites	
894		Books: online books	
374		OMIM: online Mendelian Inheritance in Man	
none		OMIA: online Mendelian Inheritance in Animals	
19529		CoreNucleotide: Core subset of nucleotide sequence records	
1156		EST: Expressed Sequence Tag records	
none		GSS: Genome Survey Sequence records	
940		Protein: sequence database	
6		Genome: whole genome sequences	
2		Structure: three-dimensional macromolecular structures	
none		Taxonomy: organisms in GenBank	
none		SNP: single nucleotide polymorphism	
493		Gene: gene-centered information	
20		HomoloGene: eukaryotic homology groups	
2		dbGaP: genotype and phenotype	
160		UniGene: gene-oriented clusters of transcript sequences	
6		CDD: conserved protein domain database	
19		3D Domains: domains from Entrez Structure	
34		UniSTS: markers and mapping data	
2		PopSet: population study data sets	
109008		GEO Profiles: expression and molecular abundance profiles	
83		GEO DataSets: experimental sets of GEO data	
123		Cancer Chromosomes: cytogenetic databases	
4		PubChem BioAssay: bioactivity screens of chemical substances	

Human Disease Genes

The screenshot shows the OMIM website interface. At the top, there is a navigation bar with tabs for 'All Databases', 'PubMed', 'Nucleotide', 'Protein', 'Genome', 'Structure', 'PMC', and 'OMIM'. The search bar contains 'OMIM' and 'for' followed by a search input field. Below the search bar are buttons for 'Limits', 'Preview/Index', 'History', 'Clipboard', and 'Details'. The display settings show 'Detailed' view, 'Show 20' items, and a 'Send to' dropdown menu. The main content area displays the entry for *120436, MutL, E. COLI, HOMOLOG OF, 1; MLH1. The entry includes a 'GeneTests, Links' link, a 'Gene map locus 3p21.3' link, and sections for 'TEXT', 'DESCRIPTION', and 'CLONING'. The 'DESCRIPTION' section states that MLH is homologous to the E. coli MutL gene and is involved in DNA mismatch repair. The 'CLONING' section describes the discovery of human MMR genes, including MLH1, and mentions the work of Genuardi et al. (1998) on alternative splicing of the MLH1 gene.

NCBI

OMIM
Online Mendelian Inheritance in Man

Johns Hopkins University

My NCBI
[Sign In] [Reg]

All Databases PubMed Nucleotide Protein Genome Structure PMC OMIM

Search OMIM for [] Go Clear

Limits Preview/Index History Clipboard Details

Display Detailed Show 20 Send to

[*120436](#) GeneTests, Links

MutL, E. COLI, HOMOLOG OF, 1; MLH1

Gene map locus [3p21.3](#)

TEXT

DESCRIPTION

MLH is homologous to the E. coli MutL gene and is involved in DNA mismatch repair. Heterozygous mutations in the MLH1 gene result in hereditary nonpolyposis colorectal cancer-2 (HNPCC2; [609310](#)) ([Papadopoulos et al., 1994](#)).

CLONING

After human homologs of the mutS gene of bacteria and yeast were found to have mutations responsible for hereditary nonpolyposis colorectal cancer (HNPCC1; [120435](#)), [Papadopoulos et al. \(1994\)](#) searched for other human mismatch repair (MMR) genes. A survey of EST databases derived from random cDNA clones revealed 3 additional human MMR genes, all related to the bacterial mutL gene. One of these genes was MLH1. The other 2 genes had a slightly greater similarity to the yeast mutL homolog PMS1 and were therefore denoted PMS1 ([600258](#)) and PMS2 ([600259](#)), respectively. 💡

[Genuardi et al. \(1998\)](#) characterized the normal alternative splicing of the MLH1 gene and reported a number of splice variants that exist in various tissue types. They observed splice variants lacking exons 6/9, 9, 9/10, 9/10/11, 10/11, 12, 16, and 17. The level of

78

Search CoreNucleotide

The screenshot shows the NCBI Nucleotide search interface. The search query is 'CoreNucleotide' for 'colon cancer'. The results show 19994 nucleotide sequences, with 18929 CoreNucleotide and 1065 EST sequences. The results are displayed in a table with columns for accession number, title, and reports. An arrow points from the 'CoreNucleotide [18929]' link to a text box explaining the database structure.

NCBI Nucleotide

Search CoreNucleotide for colon cancer

Found 19994 nucleotide sequences. CoreNucleotide [18929] EST [1065]

Display Summary Show 20 Sort by Send to

All: 18929 Bacteria: 5 RefSeq: 387 mRNA: 642

Items 1 - 20 of 18929 Page 1 of 947 Next

- 1: [NC_009045](#) Reports Links
Pichia stipitis CBS 6054 chromosome 5, complete sequence
gi126212632|ref|NC_009045.1|[126212632]
- 2: [XM_001385073](#) Reports Links
Pichia stipitis CBS 6054 highly conserved
gi126137172|ref|XM_001385073.1|[126137172]
- 3: [NM_008361](#) Reports Links
Mus musculus interleukin 1 beta (Il1b), mRNA
gi118130747|ref|NM_008361.3|[118130747]
- 4: [NM_018828](#) Reports Links
Mus musculus formin binding protein 4 (Fnm1)
gi118130721|ref|NM_018828.2|[118130721]
- 5: [NM_008628](#) Reports Links
Mus musculus mutS homolog 2 (E. coli) (Msh2), mRNA
gi118130707|ref|NM_008628.2|[118130707]

Nucleotide database now three parts:
EST expressed sequence tags
GSS genome survey sequences
CoreNucleotide everything else

79

Advanced Search Options

The screenshot shows a search results page with several annotations. A yellow box labeled "Tabs" points to the "Summary" dropdown menu. A yellow arrow points to the "bacteria: 4" tab. A yellow highlight covers the filter bar showing "CoreNucleotide (18225), EST (1061), GSS (0)".

Limits Preview/Index History Clipboard Details **Tabs**

Display: Summary Show: 20

19286 bacteria: 4 mRNA: 1658 RefSeq: 342

Show only records from: [CoreNucleotide](#) (18225), [EST](#) (1061), [GSS](#) (0). [\[What's this?\]](#)

Items 1 - 20 of 19286 Page 1 of 965 Next

- 1: [AM270351](#) Reports Links
Aspergillus niger contig An15c0240, complete genome
gi|134082757|emb|AM270351.1|[134082757]
- 2: [AM270300](#) Reports Links
Aspergillus niger contig An13c0060, complete genome
gi|134081008|emb|AM270300.1|[134081008]
- 3: [AM270178](#) Reports Links
Aspergillus niger contig An08c0230, complete genome
gi|134077487|emb|AM270178.1|[134077487]
- 4: [NM_007831](#) Reports Links
Mus musculus deleted in colorectal carcinoma (Dcc), mRNA
gi|133778956|ref|NM_007831.3|[133778956]
- 5: [NM_014059](#) Reports Links
Homo sapiens response gene to complement 32 (RGC32), mRNA
gi|132626810|ref|NM_014059.2|[132626810]

80

All Databases PubMed Nucleotide Protein Genome Structure PMC

Search CoreNucleotide for colon cancer AND nonpolyposis Go Clear

About Entrez

Entrez Nucleotide Help | FAQ

Entrez Tools

Check sequence revision history

LinkOut

My NCBI (Cubby)

Related resources BLAST

Reference sequence project

Search for Genes

Submit to GenBank

Search for full length cDNAs

Limits Preview/Index History Clipboard Details

Field: Title

- Use All Fields pull-down menu to specify a field.
- If search fields tags are used enclose in square brackets, e.g., rubella [ti].
- More help on using limits is available [here](#).

Limited to:

Fields

Title

EC/RN Number

Feature key

Filter

Gene Name

Genome Project

Issue

Journal

Keyword

Modification Date

Organism

Page Number

Primary Accession

Properties

Protein Name

Publication Date

SeqID String

Sequence Length

Substance Name

Text Word

Title

draft TPA patents

Gene Location: Any

Only from: Any

Write to the Help Desk
NCBI | NLM | NIH

colon cancer[Title] AND nonpolyposis[Title]

All Databases PubMed Nucleotide Protein Genome Structure PMC Taxonomy Books

Search CoreNucleotide for colon cancer AND nonpolyposis Go Clear

Limits Preview/Index History Clipboard Details

Field: Title

- Use All Fields pull-down menu to specify a field.
- If search fields tags are used enclose in square brackets, e.g., rubella [ti].
- More help on using limits is available [here](#).

Limited to:

Fields: [Title]

Exclude: STSs working draft TPA patents

Molecule: [mRNA] Gene Location: [Any]

Segmented Sequences: [Any] Only from: [RefSeq]

Published in the last: [Any Date]

Modified in the last: [Any Date]

About Entrez

Entrez Nucleotide Help | FAQ

Entrez Tools

Check sequence revision history

LinkOut

My NCBI (Cubby)

Related resources BLAST

Reference sequence project

Search for Genes

Submit to GenBank

Search for full length cDNAs

colon cancer[Title] AND nonpolyposis[Title] AND biomol_mrna [Properties] AND srcdb_refseq[Properties]

Advanced Search Options

The screenshot shows a search results page with a navigation bar at the top containing tabs: Limits, Preview/Index, History, Clipboard, Details, and a highlighted 'Tabs' tab. Below the navigation bar, there is a 'Display' dropdown set to 'Summary' and a 'Show' dropdown set to '20'. A yellow arrow points to the 'Summary' dropdown. Below this, there are filters for 'All: 19286', 'bacteria: 4', 'mRNA: 1658', and 'RefSeq: 342'. A yellow highlight covers the text 'Show only results from: CoreNucleotide (18225), EST (1061), GSS (0). [What's this?]'. Below the highlight, it says 'Items 1 - 20 of 19286' and 'Page 1 of 965 Next'. The main content area lists five search results, each with a checkbox, a link to the report, a title, and a 'Links' link.

Display: Summary Show: 20

All: 19286 bacteria: 4 mRNA: 1658 RefSeq: 342

Show only results from: CoreNucleotide (18225), EST (1061), GSS (0). [What's this?]

Items 1 - 20 of 19286 Page 1 of 965 Next

- 1: [AM270351](#) Reports Links
Aspergillus niger contig An15c0240, complete genome
gi|134082757|emb|AM270351.1|[134082757]
- 2: [AM270300](#) Reports Links
Aspergillus niger contig An13c0060, complete genome
gi|134081008|emb|AM270300.1|[134081008]
- 3: [AM270178](#) Reports Links
Aspergillus niger contig An08c0230, complete genome
gi|134077487|emb|AM270178.1|[134077487]
- 4: [NM_007831](#) Reports Links
Mus musculus deleted in colorectal carcinoma (Dcc), mRNA
gi|133778956|ref|NM_007831.3|[133778956]
- 5: [NM_014059](#) Reports Links
Homo sapiens response gene to complement 32 (RGC32), mRNA
gi|132626810|ref|NM_014059.2|[132626810]

All Databases PubMed Nucleotide Protein Genome Structure PMC Taxonomy Books

Search CoreNucleotide for colon cancer AND nonpolyposis AND human[Organism] Preview Go Clear

Limits Preview/Index History Clipboard Details

Field: Title Limits: mRNA, RefSeq, RefSeq

- Enter terms and click Preview to see only the number of search results.
- To save search indefinitely, click query # and select Save in My NCBI.
- To combine searches use #search, e.g., #2 AND #3 or click query # for more options.

Search	Most Recent Queries	Time	Result
#7	Search colon cancer AND nonpolyposis Field: Title Limits: mRNA, RefSeq, RefSeq	13:52:17	8
#4	Search colon cancer AND nonpolyposis Field: Title	13:39:18	21
#1	Search colon cancer	13:38:59	19006

Add Term(s) to Query or View Index:

- Enter a term in the text box; use the pull-down menu to specify a search field.
- Click Preview to add terms to the query box and see the number of search results, or click Index to view terms within a field.

Organism | human | Preview Index

Click **AND** OR NOT to add a term to the query box

- Organism
- Accession
- All Fields
- Author
- EC/RN Number
- Feature key
- Filter
- Gene Name
- Genome Project
- Issue
- Journal
- Keyword
- Modification Date
- Organization
- Page Number
- Primary Accession
- Properties
- Protein Name
- Publication Date
- SeqID String
- Sequence Length

Refining your Search

The screenshot shows a search results interface with the following elements:

- Navigation buttons: Limits (checked), Preview/Index, History, Clipboard, Details.
- Field: Title, Limits: mRNA, RefSeq.
- Display: Summary, Show: 20, Sort by, Send to.
- Summary: All: 2, Bacteria: 0, mRNA: 2, RefSeq: 2.
- Items 1 - 2 of 2, One page.
- Result 1: 1: [NM_000249](#) Reports [Links](#)
Homo sapiens mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli) (MLH1), mRNA
gi|28559089|ref|NM_000249.2||[28559089]
- Result 2: 2: [NM_000251](#) Reports [Links](#)
Homo sapiens mutS homolog 2, colon cancer, nonpolyposis type 1 (E. coli) (MSH2), mRNA
gi|4557760|ref|NM_000251.1||[4557760]

colon cancer[Title] AND nonpolyposis[Title]
AND human[Organism] AND biomol_mrna
[Properties] AND srcdb_refseq[Properties]

Useful Field Restrictions

- **[Title]:** Definition line in GenBank / GenPept format shown in Summary format
glyceraldehyde 3 phosphate dehydrogenase[Title]
- **[Organism]:** NCBI's taxonomy. Organizing system for molecular databases
mouse[organism]; green plants[organism]; Streptomyces coelicolor
[organism]
- **[Properties]:** molecule type, location, database source
biomol_mrna[properties]; biomol_genomic[properties];
gene_in_mitochondrion[properties]; srcdb_pdb[properties]
- **[Filter]:** subsets of data, Entrez links
all[filter]; nucleotide mapview[filter]; nucleotide_omim[filter]

Search CoreNucleotide for colon cancer AND nonpolyposis AND human[Organism] Go Clear [Save Search](#)

- About Entrez
- Entrez Nucleotide Help | FAQ
- Entrez Tools
- Check sequence revision history
- LinkOut
- My NCBI (Cubby)
- Related resources BLAST
- Reference sequence project
- Search for Genes
- Submit to GenBank
- Search for full length cDNAs

Limits Preview/Index History Clipboard Details

Field: **Title** Limits: **mRNA, RefSeq, RefSeq**

Display Summary Show 20 Sort by Send to

All: 2 Bacteria: 0 RefSeq: 2 mRNA: 2

Items 1 - 2 of 2

One page.

- 1: [NM_000249](#) Reports Order cDNA clone, Links
Homo sapiens mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli) (MLH1), mRNA
gil28559089|refNM_000249.2|[28559089]
- 2: [NM_000251](#) Reports Order cDNA clone, Links
Homo sapiens mutS homolog 2, colon cancer, nonpolyposis type 1 (E. coli) (MSH2), mRNA
gil4557760|refNM_000251.1|[4557760]

PubMed Nucleotide Protein Genome Structure PMC Taxonomy OMIM Books

Search for

Limits Preview/Index History Clipboard Details

Display Show Send to Hide: sequence all but gene, CDS and mRNA features

Range: from to Reverse complemented strand Features: SNP STS Exon

1: [NM_000249](#). Reports Homo sapiens mutL...[gi:28559089]

[Comment](#)
[Features](#)
[Sequence](#)

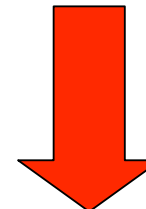
LOCUS NM_000249 2524 bp mRNA linear PRI 20-AUG-2007
DEFINITION Homo sapiens mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli) (MLH1), mRNA.
ACCESSION NM_000249
VERSION NM_000249.2 GI:28559089
KEYWORDS .
SOURCE Homo sapiens (human)
ORGANISM [Homo sapiens](#)
 Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominidae; Homo.

REFERENCE 1 (bases 1 to 2524)
AUTHORS Perri, F., Cotugno, R., Piepoli, A., Merla, A., Quitadamo, M., Gentile, A., Pilotto, A., Annese, V. and Andriulli, A.
TITLE Aberrant DNA methylation in non-neoplastic gastric mucosa of H. Pylori infected patients and effect of eradication
JOURNAL Am. J. Gastroenterol. 102 (7), 1361-1371 (2007)
PUBMED [17509026](#)
REMARK GeneRIF: While CDH1 methylation seems to be an early event in Hp gastritis, MLH1 methylation occurs late along with IM.

REFERENCE 2 (bases 1 to 2524)
AUTHORS Bettstetter, M., Dechant, S., Ruummele, P., Grabowski, M., Keller, G., Holinski-Peder, E., Hartmann, A., Hofstaedter, F. and Dietmaier, W.
TITLE Distinction of hereditary nonpolyposis colorectal cancer and sporadic microsatellite-unstable colorectal cancer through quantification of MLH1 methylation by real-time PCR
JOURNAL Clin. Cancer Res. 13 (11), 3221-3228 (2007)
PUBMED [17545526](#)
REMARK GeneRIF: quantitative MLH1 methylation analysis in MSI-H CRC is a valuable molecular tool to distinguish between HNPCC and sporadic MSI-H CRC

REFERENCE 3 (bases 1 to 2524)
AUTHORS Takahashi, M., Shimodaira, H., Andreutti-Zaugg, C., Iggo, R., Kolodner, R.D. and Ishioka, C.
TITLE Functional analysis of human MLH1 variants using yeast and in vitro mismatch repair assays
JOURNAL Cancer Res. 67 (10), 4595-4604 (2007)
PUBMED [17510385](#)
REMARK GeneRIF: The 101 MLH1 variants were examined for the dominant

- Order DNA clone
- Gene
 - HomoloGene
 - Genome
 - Genome Project
 - Master
 - Full text in PMC
 - Probe
 - Protein
 - PubMed
 - PubMed (RefSeq)
 - Gene Genotype
 - GeneView in dbSNP
 - Taxonomy
 - Related Sequences
 - Map Viewer
 - OMIM
 - GEO Profiles
 - SNP
 - UniGene
 - UniSTS
 - LinkOut



Search **CoreNucleotide** for **colon cancer AND nonpolyposis AND human[Organism]** [Go](#) [Clear](#) [Save Search](#)

- About Entrez
- Entrez Nucleotide Help | FAQ
- Entrez Tools
 - Check sequence revision history
 - LinkOut
 - My NCBI (Cubby)
 - Related resources
 - BLAST
 - Reference sequence project
 - Search for Genes
 - Submit to GenBank
 - Search for full length cDNAs

Limits Preview/Index History Clipboard Details

Field: **Title** Limits: **mRNA, RefSeq, RefSeq**

Found 2 nucleotide sequences. CoreNucleotide [2]

Display Summary Show 20 Sort by Send to

All: 2 Bacteria: 0 RefSeq: 2 mRNA: 2

Items 1 - 2 of 2

One page.

1: [NM_000249](#) Reports
Homo sapiens mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli) (MLH1), mRNA
gil28559089|reflNM_000249.2|[28559089]

2: [NM_000251](#) Reports
Homo sapiens mutS homolog 2, colon cancer, nonpolyposis type 1 (E. coli) (MSH2), mRNA
gil4557760|reflNM_000251.1|[4557760]

- Full text in PMC
- Gene
- Gene Genotype
- GeneView in dbSNP
- Genome
- Genome Project
- HomoloGene
- Master
- Probe
- Protein
- PubMed
- PubMed (RefSeq)**
- Taxonomy**
- Related Sequences
- Map Viewer
- OMIM
- GEO Profiles
- SNP
- UniGene
- UniSTS
- LinkOut

[Write to the Help Desk](#)
[NCBI](#) | [NLM](#) | [NIH](#)
[Department of Health & Human Services](#)
[Privacy Statement](#) | [Freedom of Information Act](#) | [Disclaimer](#)

Taxonomy

NCBI Entrez PubMed Nucleotide

Search for

Display 3 levels using filter:

Nucleotide Protein Structure
 3D Domains Domains GEO Datasets
 Gene HomoloGene MapView

Lineage (full): [root](#); [cellular organisms](#); [Chordata](#); [Craniata](#); [Vertebrata](#); [Gnathostomata](#); [Eutheria](#); [Euarchontoglires](#); [Primates](#);

◊ [Homo sapiens](#) (human) 11,643,469
 Click on organism name to get more information

- [Homo sapiens neanderthalensis](#)

Homo sapiens

Taxonomy ID: 9606

Genbank common name: **human**

Rank: species

Genetic code: [Translation table 1 \(Standard\)](#)

Mitochondrial genetic code: [Translation table 2 \(Vertebrate Mitochondrial\)](#)

Other names:

Other names:

common name: **man**

Lineage (full)

[cellular organisms](#); [Eukaryota](#); [Fungi/Metazoa group](#); [Metazoa](#); [Eumetazoa](#); [Bilateria](#); [Coelomata](#); [Deuterostomia](#); [Chordata](#); [Craniata](#); [Vertebrata](#); [Gnathostomata](#); [Teleostomi](#); [Euteleostomi](#); [Sarcopterygii](#); [Tetrapoda](#); [Amniota](#); [Mammalia](#); [Theria](#); [Eutheria](#); [Euarchontoglires](#); [Primates](#); [Haplorhini](#); [Simiiformes](#); [Catarrhini](#); [Hominoidea](#); [Hominidae](#); [Homo/Pan/Gorilla group](#); [Homo](#)

Entrez records		
Database name	Subtree links	Direct links
Nucleotide	11,643,469	11,642,134
Protein	392,990	392,989
Structure	9,472	9,472
Genome Sequences	51	51
Genome Projects	1	1
Popset	20,878	20,878
SNP	11,870,024	11,870,024
3D Domains	35,848	35,848
Domains	19	19
GEO Datasets	3,525	3,525
GEO Expressions	10,649,715	10,649,715
UniGene	124,179	124,179
UniSTS	322,789	322,789
PubMed Central	3,586	3,586
Gene	38,624	38,624
HomoloGene	20,167	20,167
Taxonomy	2	1

All molecular databases

Genome Information

[See the NCBI Genome homepage](#)

[Go to NCBI genomic BLAST page for Homo sapiens](#)

Genome view: 24 chromosomes																								
Names	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	X	Y

Goal: Find MLH1 homologs

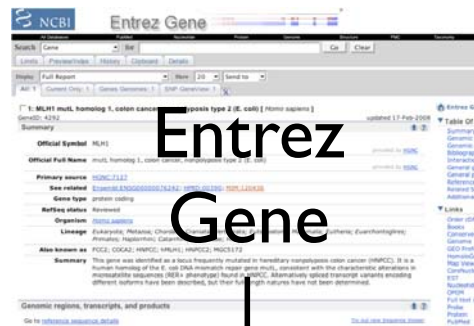
- **Tip:** Use Entrez Gene as your hub to connect to everything else!



Protein

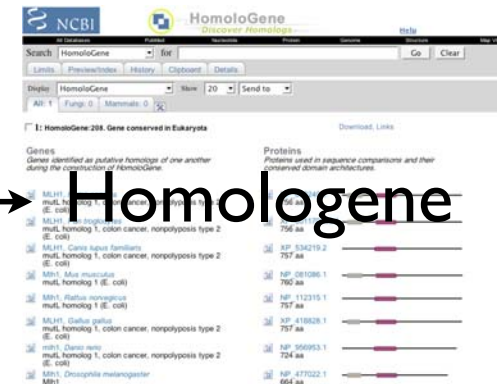


BLink



Entrez
Gene

Other Entrez
Databases



Homologene

Gene neighbors



Search CoreNucleotide for colon cancer AND nonpolyposis AND human[Organism] Go Clear Save Search

- About Entrez
- Entrez Nucleotide Help | FAQ
- Entrez Tools
 - Check sequence revision history
 - LinkOut
 - My NCBI (Cubby)
 - Related resources
 - BLAST
 - Reference sequence project
 - Search for Genes
 - Submit to GenBank
 - Search for full length cDNAs

Limits Preview/Index History Clipboard Details

Field: Title Limits: mRNA, RefSeq, RefSeq
Found 2 nucleotide sequences. CoreNucleotide [2]

Display Summary Show 20 Sort by Send to

All: 2 Bacteria: 0 RefSeq: 2 mRNA: 2

Items 1 - 2 of 2

One page.

- 1: [NM_000249](#) Reports
Homo sapiens mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli) (MLH1), mRNA
gil28559089|reflNM_000249.2|[28559089]
- 2: [NM_000251](#) Reports
Homo sapiens mutS homolog 2, colon cancer, nonpolyposis type 1 (E. coli) (MSH2), mRNA
gil4557760|reflNM_000251.1|[4557760]

- Links
- Full text in PMC
 - Gene
 - Gene Genotype
 - GeneView in dbSNP
 - Genome
 - Genome Project
 - HomoloGene
 - Master
 - Probe
 - Protein
 - PubMed
 - PubMed (RefSeq)
 - Taxonomy
 - Related Sequences
 - Map Viewer
 - OMIM
 - GEO Profiles
 - SNP
 - UniGene
 - UniSTS
 - LinkOut

Write to the Help Desk
NCBI | NLM | NIH
Department of Health & Human Services
Privacy Statement | Freedom of Information Act | Disclaimer

MLH1 Gene Record

1: MLH1 mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli) [*Homo sapiens*]

GeneID: 4292

updated 10-Apr-2007

Summary

Official Symbol MLH1

provided by [HGNC](#)

Official Full Name mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli)

provided by [HGNC](#)

Primary source [HGNC:7127](#)

See related [HPRD:0039](#)

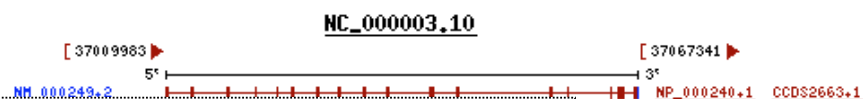
Gene type protein coding

RefSeq status Reviewed

Organism [Homo sapiens](#)

Genomic regions, transcripts, and products

Go to [reference sequence details](#)

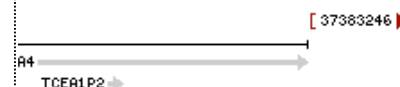


GeneRIFs: Gene References Into Function

[What's a GeneRIF?](#)

1. Results confirmed complete exon skipping for the mutations of MLH1 in hereditary nonpolyposis colorectal cancer patients.
2. hMLH1 may have a role in development of secondary carcinoma in the gastrointestinal tract in patients (stomach and colorectal carcinoma)
3. Inactivation of MLH1 gene is associated with head and neck squamous cell carcinoma tumors and leukoplakia
4. In three adenocarcinomas, microsatellite instability and lack of the MLH1 protein expression were detected.
5. MLH1 is associated with longevity.
6. The identification of residues whose mutation disrupts MutL-MutS interaction and affects mismatch repair activity, suggesting a mechanism by which hereditary mutations in this region can produce a cancer predisposition.
7. These results indicate that an age-related increase of medullary-type tumors in poorly differentiated adenocarcinoma may play an important

[See MLH1 in MapViewer](#)



Interactions + GO

Interactions

Description					
Product	Interactant	Other Gene	Complex	Source	P
E2F1 interacts with the MLH1 promoter.					
NC_000003.9	NP_005216.1	E2F1		BIND	
E2F4 interacts with the MLH1 promoter region.					
NC_000003.9	NP_001941.2	E2F4		BIND	
NP_000240.1	NP_000048.1	BLM		HPRD	
MLH1 interacts with BLM.					
NP_000240.1	NP_000048.1	BLM		BIND	
NP_000240.1	NP_009225.1	BRCA1		HPRD	
The exonuclease HEX1 interacts with the mismatch repair protein hMLH1.					
NP_000240.1	NP_003677.3	EXO1		BIND	
The exonuclease hEXO1b interacts with the mismatch repair protein hMLH1.					
NP_000240.1	NP_006018.3	EXO1		BIND	
NP_000240.1	NP_569082.1	EXO1		HPRD	
NP_000240.1	NP_003916.1	MBD4		HPRD	
MLH1 and interacts with MED1.					
NP_000240.1	NP_003916.1	MBD4		BIND	
NP_000240.1	BAA92353.1	MLH3		HPRD	

GeneOntology

Provided by [GOA](#)

Function	Evidence
ATP binding	IEA
contributes_to MutSalpha complex binding	IDA Pubmed
guanine/thymine mispair binding	IMP Pubmed
guanine/thymine mispair binding	IEA
mismatched DNA binding	IEA
protein binding	IPI Pubmed
contributes_to single-stranded DNA binding	IDA Pubmed
Process	Evidence
DNA damage response, signal transduction resulting in induction of apoptosis	IEA
cell cycle	IEA
male meiosis chromosome segregation	IEA
meiotic recombination	IEA
mismatch repair	IEA
mismatch repair	TAS Pubmed
negative regulation of mitotic recombination	IEA
negative regulation of progression through cell cycle	IEA
Component	Evidence
MutLalpha complex	IEA
condensed chromosome	IEA
nucleus	IC Pubmed
nucleus	IEA
synaptonemal complex	IEA

MLH1: Sequence Links

Genomic regions, transcripts, and products ↑ ?

Go to [reference sequence details](#)

NC_000003.10

5' 3'

[NM_000249.2](#) [NP_000240.1](#) [CCDS2663.1](#)

■ - coding region ■ - untranslated region

Links

mRNA LINKS

- ▶ FASTA
- ▶ GENBANK

Links

PROTEIN LINKS

- ▶ FASTA
- ▶ GENPEPT
- ▶ Blink
- ▶ Conserved Domains

chromosome: 3; Location: 3p21.3

[36992791 ▶ [37383246 ▶

LOC645571 ▶ LRRFIP2 ← GOLGA4 TCEA1P2 →

EPM2AIP1 ← MLH1 →

▼ **Links** [Explain](#)

- [Order cDNA clone](#)
- [Books](#)
- [Conserved Domains](#)
- [Genome](#)
- [GEO Profiles](#)
- [HomoloGene](#)
- [Map Viewer](#)
- [Nucleotide](#)
- [OMIM](#)
- [Full text in PMC](#)
- [Probe](#)
- [Protein](#)
- [PubMed](#)
- [PubMed \(GeneRIF\)](#)
- [SNP](#)
- [SNP: Genotype](#)
- [SNP: GeneView](#)
- [Taxonomy](#)
- [UniSTS](#)
- [AceView](#)
- [CCDS](#)
- [Colon.html](#)
- [Evidence Viewer](#)
- [GDB](#)
- [GeneTests for MIM: 120436](#)
- [HGMD](#)
- [HGNC](#)
- [HPRD](#)
- [KEGG](#)
- [MGC](#)
- [ModelMaker](#)
- [PharmGKB](#)
- [UniGene](#)
- [LinkOut](#)



Search Gene for

Display Full Report Show 20 Send to

All: 1 Current Only: 1 Genes Genomes: 1 SNP GeneView: 1

1: **MLH1 mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli) [*Homo sapiens*]**

GeneID: 4292 updated 16-Sep-2007

Summary

Official Symbol	MLH1	<small>provided by HGNC</small>
Official Full Name	mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli)	<small>provided by HGNC</small>
Primary source	HGNC:7127	
See related	Ensembl:ENSG00000076242 ; HPRD:00390 ; MIM:120436	
Gene type	protein coding	
RefSeq status	Reviewed	
Organism	Homo sapiens	
Lineage	<i>Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominidae; Homo</i>	
Also known as	FCC2; COCA2; HNPCC; hMLH1; HNPCC2; MGC5172	
Summary	This gene was identified as a locus frequently mutated in hereditary nonpolyposis colon cancer (HNPCC). It is a human homolog of the E. coli DNA mismatch repair gene mutL, consistent with the characteristic alterations in microsatellite sequences (RER+ phenotype) found in HNPCC. Alternatively spliced transcript variants encoding different isoforms have been described, but their full-length natures have not been determined.	

[Entrez Gene Home](#)

Table Of Contents

- [Summary](#)
- [Genomic regions, transcripts...](#)
- [Genomic context](#)
- [Bibliography](#)
- [Interactions](#)
- [General gene information](#)
- [General protein information](#)
- [Reference Sequences](#)
- [Related Sequences](#)
- [Additional Links](#)

Links [Explain](#)

- [Order cDNA clone](#)
- [Books](#)
- [Conserved Domains](#)
- [Genome](#)
- [GEO Profiles](#)
- [HomoloGene](#)
- [Map Viewer](#)
- [CoreNucleotide](#)
- [EST](#)
- [Nucleotide](#)
- [OMIM](#)
- [Full text in PMC](#)
- [Probe](#)
- [Protein](#)
- [PubMed](#)
- [PubMed \(GeneRIF\)](#)
- [SNP](#)
- [SNP: Genotype](#)

Genomic regions, transcripts, and products

Go to [reference sequence details](#)

NC_000003.10

Finding Homologs:

All: 1 Fungi: 0 Mammals: 0 ✕

1: HomoloGene:208. Gene conserved in Eukaryota

H.sapiens	MLH1	mutL homolog 1, colon
P.troglodytes	MLH1	MutL protein homolog 1
C.familiaris	LOC477019	similar to MutL protein
M.musculus	Mlh1	mutL homolog 1 (E. col
R.norvegicus	Mlh1	mutL homolog 1 (E. col
G.gallus	MLH1	mutL homolog 1, colon
D.melanogaster	Mlh1	Mlh1
A.gambiae	AgaP_ENSA...	ENSANGP0000001401
A.gambiae	ENSANGG00...	ENSANGP0000001348
S.pombe	SPBC1703.04	hypothetical protein
S.cerevisiae	MLH1	Mlh1p
K.lactis	KLLA0D099...	mRNA gene KLLA0D09
E.gossypii	GeneID:27...	Eremothecium gossypii
N.crassa	NCU08309.1	hypothetical protein
A.thaliana	ATMLH1	ATMLH1
O.sativa	Os01g0958...	mRNA gene Os01g0958

HomoloGene Downloader

[Homologene:208](#). Gene conserved in Eukaryota

Download **Protein** sequences (in FASTA format)

Include bp upstream of gene

Include bp downstream of gene

Select which sequences should be included

Species	Gene			
<input checked="" type="checkbox"/>	H.sapiens	MLH1	NM_00...	
<input checked="" type="checkbox"/>	P.troglodytes	MLH1	XM_00...	
<input checked="" type="checkbox"/>	C.familiaris	LOC477019	XM_534...	
<input checked="" type="checkbox"/>	M.musculus	Mlh1	NM_02...	
<input checked="" type="checkbox"/>	R.norvegicus	Mlh1	NM_031053.1	NP_112315.1
<input checked="" type="checkbox"/>	G.gallus	MLH1	XM_418828.1	XP_418828.1
<input checked="" type="checkbox"/>	D.melanogaster	Mlh1	NM_057674.2	NP_477022.1
<input checked="" type="checkbox"/>	A.gambiae	AgaP_ENSANGG00000011527	XM_320342.2	XP_320342.2
<input checked="" type="checkbox"/>	A.gambiae	ENSANGG00000010995	XM_307435.2	XP_307435.2
<input checked="" type="checkbox"/>	S.pombe	SPBC1703.04	NM_001022118.1	NP_596199.1
<input checked="" type="checkbox"/>	S.cerevisiae	MLH1	MLH1_6323819	NP_013890.1
<input checked="" type="checkbox"/>	K.lactis	KLLA0D09955g	XM_453504.1	XP_453504.1
<input checked="" type="checkbox"/>	E.gossypii	GeneID:2757243	NM_210705.1	NP_985351.1
<input checked="" type="checkbox"/>	N.crassa	NCU08309.1	XM_329014.1	XP_329015.1
<input checked="" type="checkbox"/>	A.thaliana	ATMLH1	NM_116983.2	NP_567345.2
<input checked="" type="checkbox"/>	O.sativa	Os01g0958900	NM_001051992.1	NP_001045457.1

Protein
mRNA
Genomic

HomoloGene Cluster

1: HomoloGene:208. Gene conserved in Eukaryota

[Download](#), [Links](#)

Genes

Genes identified as putative homologs of one another during the construction of HomoloGene.

- [H.sapiens MLH1](#)
mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli)
- [P.trogodytes MLH1](#)

Proteins

Proteins used in sequence comparisons and their conserved domain architectures.

- [NP_000240.1](#)
756 aa
- [XP_001170433.1](#)

[M.musculus Mlh1](#)

Links

- Conserved Domains
- Genome
- GEO Profiles
- Nucleotide
- Order cDNA clone
- OMIM
- Full text in PMC
- Probe
- Protein
- PubMed
- PubMed (GeneRIF)
- SNP
- Gene Genotype
- GeneView in dbSNP
- Taxonomy
- UniGene
- UniSTS
- MapViewer

1 (E. coli)

- [mutL homolog 1 \(E. coli\)](#)
- [R.norvegicus Mlh1](#)
mutL homolog 1 (E. coli)
- [G.gallus MLH1](#)
mutL homolog 1, colon cancer, nonpolyposis type 2 (E. coli)
- [D.melanogaster Mlh1](#)
Mlh1
- [A.gambiae AgaP_ENSANGG00000014016](#)
ENSANGP00000014016
- [A.gambiae ENSANGG00000010995](#)
ENSANGP00000013484
- [S.pombe SPBC1703.04](#)
hypothetical protein
- [S.cerevisiae MLH1](#)
Mlh1p
- [K.lactis KLLA0D09955g](#)
mRNA gene KLLA0D09955g
- [E.gossypii GeneID:2757243](#)
Eremothecium gossypii AFL199C gene
- [N.crassa NCU08309.1](#)
hypothetical protein
- [A.thaliana ATMLH1](#)
ATMLH1
- [O.sativa Os01g0958900](#)
mRNA gene Os01g0958900

Links

- Conserved Domains
- Gene
- Genome Project
- Nucleotide
- Genome
- OMIM
- Full text in PMC
- Related Sequences
- Domain Relatives
- PubMed
- PubMed (RefSeq)
- SNP
- Gene Genotype
- GeneView in dbSNP
- Related Structure
- Taxonomy
- UniGene
- BLink
- Domains

[NP_081086.1](#)

760 aa

- [760 aa](#)
- [NP_112315.1](#)
757 aa
- [XP_418828.1](#)
757 aa
- [NP_477022.1](#)
664 aa
- [XP_320342.2](#)
671 aa
- [XP_307435.2](#)
395 aa
- [NP_596199.1](#)
684 aa
- [NP_013890.1](#)
769 aa
- [XP_453504.1](#)
724 aa
- [NP_985351.1](#)
771 aa
- [NP_001045457.1](#)
724 aa

Gene Links

Protein Links

Finding Homologs 2: BLink

Genomic regions, transcripts, and products ↑ ?

Go to [reference sequence details](#)

NC_000003.10

[37009983] 5' ————— [37067341] 3'

[NM_000249.2](#) ■ - coding region ■ - untranslated region [NP_000240.1](#)

Links

PROTEIN LINKS

- ▶ [FASTA](#)
- ▶ [GENPEPT](#)
- ▶ [Blink](#)
- ▶ [Conserved Domains](#)

1: [NP_000240](#). Reports MutL protein homo...[gi:4557757] ▶ BLink, Conserved Domains, Links

[Comment](#) [Features](#) [Sequence](#)

LOCUS NP_000240 756 aa linear PRI 08-APR-2007

DEFINITION MutL protein homolog 1 [Homo sapiens].

ACCESSION NP_000240

VERSION NP_000240.1 GI:4557757

DBSOURCE REFSEQ: accession [NM_000249.2](#)

KEYWORDS .

SOURCE Homo sapiens (human)

ORGANISM [Homo sapiens](#)
 Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
 Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
 Catarrhini; Hominidae; Homo.

REFERENCE 1 (residues 1 to 756)

AUTHORS Marmo,R., Rotondano,G., Riccio,G., D'Angella,R., Rescinito,M., Rescinito,A., Bianco,M.A. and Cipolletta,L.

TITLE Small-bowel adenocarcinoma diagnosed via capsule endoscopy in a patient found to have hereditary nonpolyposis colorectal cancer

JOURNAL Gastrointest. Endosc. 65 (3), 524-525 (2007)

PUBMED [17208239](#)

BLink: BLAST Link (Best Hits)

Query: gi|4557757 MutL protein homolog 1 [Homo sapiens]

Matching gi: [33738032](#), [119584889](#), [27805155](#), [53932122](#), [13905126](#), [14107168](#), [14120083](#), [463989](#), [720920](#), [741682](#), [1070207](#), [01122084](#), [75916970](#), [40200246](#), [31688772](#)

COG0323 assigned by Cognitor (5 best hits)

Show identical All hits Common Tree Taxonomy Report 3D structures CD

Redundant Proteins
First 200 only

198 BLAST hits to 7 selected species [Sort by taxonomy proximity](#)

3 Archaea 72 Bacteria 51 Metazoa 39 Fungi 9 Plants 0 Viruses 22 Other Eukaryotae

Keep only Cut-Off 100 Select Reset

New search by GI: Go

756 aa

SCORE	P	ACCESSION	GI	N	ORGANISM
Conserved Domain Database hits					
3869	1	AAQ02400	33303773	1	synthetic
3868	31	AAA17374	466462	8	Homo sapiens
3860	29	XP_001...	114585960	8	Proglodytes
3615	21	XP_534219	73989704	1	Canis familiaris
3442	22	BAE40671	74223060	6	Mus musculus
3380	22	P97679	13872071	1	Rattus norvegicus
3111	18	XP_418828	5732924	1	Gallus gallus
3009	20	XP_001...	126336756	2	Monodelphis domestica
2915	17	AAI24967	117167959	1	Xenopus laevis
2893	26	XP_001...	109042257	2	Macaca mulatta
2633	15	NP_956953	41054934	1	Danio rerio
2393	9	XP_001...	115932300	1	Strongylocentrotus purpuratus
2331	15	CAG04734	47216556	1	Tetraodon nigroviridis
1892	8	XP_001...	110763401	1	Apis mellifera
1690	8	EAT42626	108878401	1	Aedes aegypti
1679	8	AAA00135	116116919	1	Anopheles gambiae str. PEST
1663	8	XP_001...	125807138	1	Drosophila pseudoobscura
1662	8	AAF59117	7304079	3	Drosophila melanogaster
1558	3	XP_637285	66807125	1	Dictyostelium discoideum AX4
1503	4	XP_962522	85108177	1	Neurospora crassa OR74A

BLAST

Opossum homolog



Sequence Databases

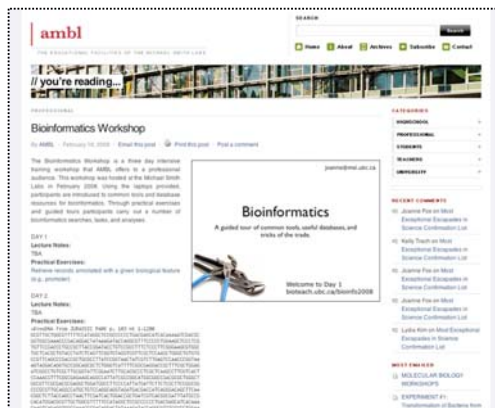
PRACTICAL EXERCISES: Navigating Links, Retrieving Data
with Entrez, and Searching PubMed



I am studying the regulation of cancer genes and would like to retrieve all human sequence records associated with cancer that contain a promoter region.

navigate to:
bioteach.ubc.ca/bioinfo2008

Let's compare
our results



Follow link to practical exercise
page at the NCBI where you'll find
step-by-step instructions



Use the preview tab and feature keys

Strategy #1:
search nt

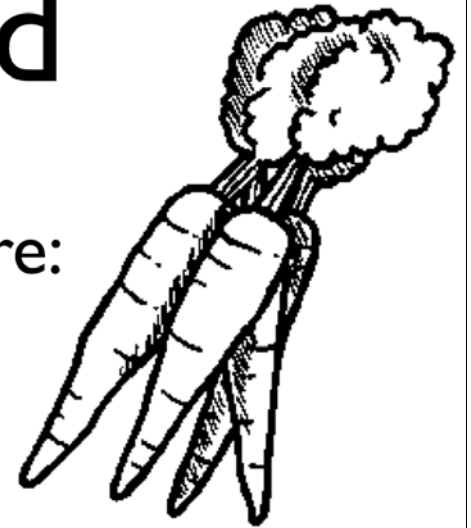
Strategy #2: search
entrez gene

Check your History

Search	Most Recent Queries	Result
#5	Search #3 NOT #1 (unique hits from Approach B: Entrez Gene to CoreNucleotide)	286
#4	Search #1 NOT #3 (unique hits from Approach A: straight to Entrez CoreNucleotide search)	183
#3	Search #2 AND promoter[Feature key] (limit Approach B search to records with promoter annotated)	329
#2	CoreNucleotide Links for Gene (Search human [Organism] AND cancer[Text Word] AND gene_nucleotide[Filter]) (Approach B: Entrez gene follow link to CoreNucleotide)	53844
#1	Search human[Organism] AND cancer[Text Word] AND promoter[Feature key] (Approach A: Entrez CoreNucleotide search)	226

Searching PubMed

- How many papers in PubMed are there:
 - about cancer?
 - about carrots?
- Using Entrez PubMed, can you see if there is any scientific links between carrots and cancer?
- How many papers are there about “carrots AND cancer”?
- What is the active chemical substance in carrots that may play a role in cancers?



About Entrez
Text Version
Entrez PubMed
Overview
Help | FAQ
Tutorials
New/Noteworthy
E-Utilities
PubMed Services
Journals Database
MeSH Database
Single Citation
Matcher
Batch Citation
Matcher
Clinical Queries
Special Queries
LinkOut
My NCBI
Related Resources
Order Documents
NLM Mobile
NLM Catalog
NLM Gateway
TOXNET
Consumer Health
Clinical Alerts
ClinicalTrials.gov
PubMed Central

- 1:** [Collins AR, Gaivao I.](#) Related Articles, Links
 DNA base excision repair as a biomarker in molecular epidemiology studies.
 Mol Aspects Med. 2007 Jun 2; [Epub ahead of print]
 PMID: 17659329 [PubMed - as supplied by publisher]
- 2:** [Young JF, Duthie SJ, Milne L, Christensen LP, Duthie GG, Bestwick CS.](#) Related Articles, Links
 Biphasic effect of falcarinol on caco-2 cell proliferation, DNA damage, and apoptosis.
 J Agric Food Chem. 2007 Feb 7;55(3):618-23.
 PMID: 17263451 [PubMed - indexed for MEDLINE]
- 3:** [Galeone C, Negri E, Pelucchi C, La Vecchia C, Bosetti C, Hu J.](#) Related Articles, Links
 Dietary intake of fruit and vegetable and lung cancer risk: a case-control study in Harbin, northeast China.
 Ann Oncol. 2007 Feb;18(2):388-92. Epub 2006 Oct 23.
 PMID: 17060488 [PubMed - indexed for MEDLINE]
- 4:** [Roumanas ED, Garrett N, Blackwell KE, Freymiller E, Abemayor E, Wong WK, Beumer J 3rd, Fueki K, Fueki W, Kapur KK.](#) Related Articles, Links
 Masticatory and swallowing threshold performances with conventional and implant-supported prostheses after mandibular fibula free-flap reconstruction.
 J Prosthet Dent. 2006 Oct;96(4):289-97.
 PMID: 17052474 [PubMed - indexed for MEDLINE]
- 5:** [Simon HB.](#) Related Articles, Links
 On call. My 77-year-old father is healthy, but his older brother has just been diagnosed with prostate cancer. Dad says he read that carrot juice will prevent prostate cancer, and he's now drinking it every day. Is he just kidding himself?
 Harv Mens Health Watch. 2006 Jun;10(11):8. No abstract available.
 PMID: 16775868 [PubMed - indexed for MEDLINE]

- Search History will be lost after eight hours of inactivity.
- Search numbers may not be continuous; all searches are represented.
- To save search indefinitely, click query # and select Save in My NCBI.
- To combine searches use #search, e.g., #2 AND #3 or click query # for more options.

Search	Most Recent Queries	Time	Result
#22	Search cancer AND carrots	17:18:07	115
#21	Search carrots	17:17:56	1419
#20	Search cancer	17:17:48	1957409

Clear History

Special Queries
LinkOut
My NCBI

Related Resources
Order Documents
NLM Mobile
NLM Catalog
NLM Gateway
TOXNET
Consumer Health
Clinical Alerts
ClinicalTrials.gov
PubMed Central

Humans or Animals CLEAR

Humans Animals

Gender CLEAR

Male Female

Languages CLEAR

English
 French
 German
 Italian
 Japanese
 Russian
 Spanish

More Languages

Afrikaans
 Albanian

Type of Article

Clinical Trial
 Editorial
 Letter
 Meta-Analysis
 Practice Guideline
 Randomized
 Review

More Publication

Addresses
 Bibliography

Grant Number
Issue
Journal
Language
Last Author
MeSH Date
MeSH Major Topic
MeSH Subheading
MeSH Terms
Pagination
Pharmacological Action
Publication Date
Publication Type
Secondary Source ID
Substance Name
Text Word

Tag Terms

Title
Title/Abstract
Transliterated Title
Volume

Default Tag:

All Fields

Subsets CLEAR

Journal Groups

Core clinical journals
 Dental journals
 Nursing journals

Topics

AIDS
 Bioethics
 Cancer
 Complementary Medicine
 History of Medicine

Ages CLEAR

All Infant: birth-23 months
 All Child: 0-18 years
 All Adult: 19+ years
 Newborn: birth-1 month
 Infant: 1-23 months
 Preschool Child: 2-5 years
 Child: 6-12 years
 Adolescent: 13-18 years
 Adult: 19-44 years
 Middle Aged: 45-64 years

GO

Clear All Limits

[Write to the Help Desk](#)
[NCBI](#) | [NLM](#) | [NIH](#)
Department of Health & Human Services
[Privacy Statement](#) | [Freedom of Information Act](#) | [Disclaimer](#)

Search PubMed for cancer AND carrot Go Clear [Save Search](#) Limits Preview/Index History Clipboard Details

Field: Title

Display Summary Show 20 Sort By Send to

All: 8 Review: 0

Items 1 - 8 of 8

One page.

- 1: [Konety BR.](#) [Related Articles, Links](#)
Bladder cancer prevention--could a carrot be the stick?
J Urol. 2006 Sep;176(3):864-5. No abstract available.
PMID: 16890640 [PubMed - indexed for MEDLINE]
- 2: [Simon HB.](#) [Related Articles, Links](#)
On call. My 77-year-old father is healthy, but his older brother has just been diagnosed with prostate cancer. Dad says he read that carrot juice will prevent prostate cancer, and he's now drinking it every day. Is he just kidding himself?
Harv Mens Health Watch. 2006 Jun;10(11):8. No abstract available.
PMID: 16775868 [PubMed - indexed for MEDLINE]
- 3: [Ambrosini GL.](#) [Related Articles, Links](#)
Does drinking carrot juice affect cancer of the prostate?
Med J Aust. 2001 Jul 2;175(1):53; author reply 53-4. No abstract available.
PMID: 11476210 [PubMed - indexed for MEDLINE]
- 4: [Vitetta L, Sali A, Reavley NJ.](#) [Related Articles, Links](#)
Does drinking carrot juice affect cancer of the prostate?
Med J Aust. 2001 Jul 2;175(1):52-3; author reply 53-4. No abstract available.
PMID: 11476209 [PubMed - indexed for MEDLINE]
- 5: [Campbell GR.](#) [Related Articles, Links](#)
Does drinking carrot juice affect cancer of the prostate?
Med J Aust. 2001 Jul 2;175(1):51; author reply 53-4. No abstract available.
PMID: 11476208 [PubMed - indexed for MEDLINE]

[Links](#)

- ▶ Compound via MeSH
- ▶ Substance via MeSH
- ▶ LinkOut

Search PubMed for cancer and carrots Go Clear [Save Search](#)

Limits Preview/Index History Clipboard Details

Field: Title


Display Summary Show 20 Sort By Send to


All: 3 Review: 0


Items 1 - 3 of 3

One page.

- 1: [Longnecker MP, Newcomb PA, Mittendorf R, Greenberg ER, Willett WC.](#) Related Articles, Links

 Intake of carrots, spinach, and supplements containing vitamin A in relation to risk of breast cancer. *Cancer Epidemiol Biomarkers Prev.* 1997 Nov;6(11):887-92. PMID: 9367061 [PubMed - indexed for MEDLINE]
- 2: [Jacobsen BK.](#) Related Articles, Links

 [Vegetables and prevention of cancer. Carrots are still good for you] *Tidsskr Nor Laegeforen.* 1988 Oct 30;108(30):2744-6. Norwegian. No abstract available. PMID: 3206486 [PubMed - indexed for MEDLINE]
- 3: [Pisani P, Berrino F, Macaluso M, Pastorino U, Crosignani P, Baldasseroni A.](#) Related Articles, Links

 Carrots, green vegetables and lung cancer: a case-control study. *Int J Epidemiol.* 1986 Dec;15(4):463-8. PMID: 3818153 [PubMed - indexed for MEDLINE]

- About Entrez
- Text Version
- Entrez PubMed
 - Overview
 - Help | FAQ
 - Tutorials
 - New/Noteworthy 
 - E-Utilities
- PubMed Services
 - Journals Database
 - MeSH Database
 - Single Citation Matcher
 - Batch Citation Matcher
 - Clinical Queries
 - Special Queries
 - LinkOut
 - My NCBI
- Related Resources
 - Order Documents

Search PubMed for

Limits

Field: Title

Display AbstractPlus Show 20 Sort By Send to

All: 1 Review: 0

1: [Cancer Epidemiol Biomarkers Prev.](#) 1997 Nov;6(11):887-92.

Full Text **FREE** [Links](#)
Cancer Epid Biomark

Intake of carrots, spinach, and supplements containing vitamin A in relation to risk of breast cancer.

[Longnecker MP](#), [Newcomb PA](#), [Mittendorf R](#), [Greenberg ER](#), [Willett WC](#).

Epidemiology Branch, National Institute of Environmental Health Sciences, Research Triangle Park, North Carolina 27709, USA.

Intake of fruits, vegetables, vitamin A, and related compounds are associated with a decreased risk of breast cancer in some studies, but additional data are needed. To estimate intake of beta-carotene and vitamin A, the authors included nine questions on food and supplement use in a population-based case-control study of breast cancer risk conducted in Maine, Massachusetts, New Hampshire, and Wisconsin in 1988-1991. Multivariate-adjusted models were fit to data for 3543 cases and 9406 controls. Eating carrots or spinach more than twice weekly, compared with no intake, was associated with an odds ratio of 0.56 (95% confidence interval 0.34-0.91). Estimated intake of preformed vitamin A from all evaluated foods and supplements showed no trend or monotonic decrease in risk across categories of intake. These data do not allow us to distinguish among several potential explanations for the protective association observed between intake of carrots and spinach and risk of breast cancer. The findings are, however, consistent with a diet rich in these foods having a modest protective effect.

PMID: 9367061 [PubMed - indexed for MEDLINE]

Display AbstractPlus Show 20 Sort By Send to

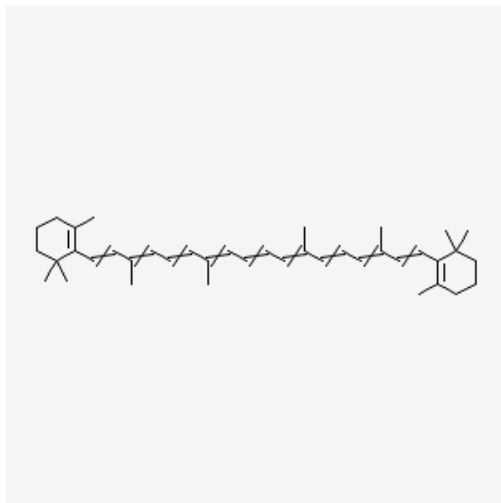
Related Links

- ▶ [Vitamins C and E, retinol, beta-carotene and dietary fibre in relation to breast cancer risk \[Br J Cancer. 1997\]](#)
- ▶ [Population attributable risk for breast cancer: diet, nutrition, and physical exercise \[J Natl Cancer Inst. 1998\]](#)
- ▶ [Fruits, vegetables, and micronutrients in relation to breast cancer m- \[Cancer Epidemiol Biomarkers Prev. 2004\]](#)
- ▶ [Intake of dietary fat and vitamin in relation to breast cancer risk in Korean women: : \[J Korean Med Sci. 2003\]](#)
- ▶ [Dietary carotenoids and vitamins A, C, and E and risk of breast cancer. \[J Natl Cancer Inst. 1999\]](#)

[See all Related Articles...](#)

Substance Summary:

Compound Displayed



 **SID:** [4266322](#) 
 **CID:** [573](#) 


 **Related Substances:** 
 Same: [11 Links](#)
 Same, Connectivity: [12 Links](#)

 **Structure Search** 

Source: [LipidMAPS \(LMPR01070001\)](#) 

[MeSH](#) | [Synonyms](#) | [Properties](#) | [Descriptors](#) | [Exports](#)

 **Medical Subject Annotations:** (Total:1) 

 **beta Carotene**
 A carotenoid that is a precursor of VITAMIN A. It is administered to reduce the severity of photosensitivity reactions in patients with erythropoietic protoporphyria (PORPHYRIA, ERYTHROPOIETIC). (From Reynolds JEF(Ed): Martindale: The Extra Pharmacopoeia (electronic version). Micromedex, Inc, Engewood, CO, 1995.)

[Show MeSH Tree Structure](#)

Pharmacological Action:
[Antioxidants](#)

Search

All Databases



for

Go

SITE MAPAlphabetical List
Resource Guide**About NCBI**An introduction to
NCBI**GenBank**Sequence
submission support
and software**Literature
databases**PubMed, OMIM,
Books, and PubMed
Central**Molecular
databases**Sequences,
structures, and
taxonomy**Genomic
biology**The human
genome, whole
genomes, and
related resources**Tools**

Data mining

▶ What does NCBI do?

Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease. [More...](#)

Hot Spots

- ▶ Assembly Archive
- ▶ Clusters of orthologous groups
- ▶ Coffee Break, Genes & Disease, NCBI Handbook
- ▶ Electronic PCR
- ▶ Entrez Home
- ▶ Entrez Tools
- ▶ Gene expression omnibus (GEO)
- ▶ Human genome resources
- ▶ Influenza Virus Resource
- ▶ Map Viewer
- ▶ dbMHC
- ▶ Mouse genome resources
- ▶ My NCBI

New**Protein Clusters**

Entrez Protein Clusters database

The new Entrez Protein Clusters database is a collection of Reference Sequence (RefSeq) proteins, from the complete genomes of prokaryotes, plasmids, and organelles, that have been grouped and annotated based on sequence similarity and protein function. Click here to find out more about the [Protein Clusters](#) database.

**1 Billion Live Traces**

The Trace Archive of sequencing traces has reached 1 billion live traces from over 480 organisms. For more information about the Trace Archive database [click here](#).

PubMed Central

Search Books for Go Clear [Save Search](#)

Limits Preview/Index History Clipboard Details

Display Books Show 20 Send to

All: 196 **Figures: 2**

- About Entrez
- Books**
- Overview
- Using the books
- Information for authors and publishers
- Contact us
- Mailing list
- Project background
- FAQ
- My NCBI
- Privacy Policy

- 

171 items in **Health Services/Technology Assessment Text (HSTAT)**
Bethesda (MD): [National Library of Medicine](#) (US), 2003 Oct.
- 

12 items in **Cancer Medicine**. 6th ed.
Kufe, Donald W.; Pollock, Raphael E.; Weichselbaum, Ralph R.; Bast, Robert C., Jr.; Gansler, Ted S.; Holland, James F.; Frei III, Emil, editors.
Hamilton (Canada): [BC Decker Inc](#); c2003.
- 

2 items in **Molecular Cell Biology**. 4th ed.
Lodish, Harvey; Berk, Arnold; Zipursky, S. Lawrence; Matsudaira, Paul; Baltimore, David; Darnell, James E.
New York: [W. H. Freeman & Co.](#); c2000.
- 

2 items in **Disease Control Priorities in Developing Countries** 2nd ed.
Dean T. Jamison, Joel G. Breman, Anthony R. Measham, George Alleyne, Mariam Claeson, David B. Evans, Prabhat Jha, Anne Mills, Philip Musgrove, editors
Washington (DC): [IBRD/The World Bank and Oxford University Press](#); 2006
- 

2 items in **Alternative Medicine and Rehabilitation: A Guide for Practitioners**
Wainapel, Stanley F.; Fast, Avital, editors
New York: [Demos Medical Publishing](#); c2003
- 

1 item in **The Organization of the Retina and Visual System**
Kolb, Helga; Fernandex, Eduardo; Nelson, Raph, editors
Salt Lake City (UT): [Moran Eye Center, University of Utah](#); 2000
- 

1 item in **Surgical Treatment**
Holzheimer, Rene G.; Mannick, John A., editors.
Munich: [Zuckschwerdt Publishers](#); c2001.
- 

1 item in **Parkinson's Disease: Diagnosis and Clinical Management**
Factor, Stewart A.; Weiner, William J.
New York: [Demos Medical Publishing, Inc.](#); c2002

Navigation

About this book

Section 4: Cancer Epidemiology, Prevention, and Screening

29. Nutrition in the Etiology and Prevention of Cancer

Methodologic Issues in Diet, Nutrition, and Cancer Studies

Public Health Guidelines for Cancer Prevention

Summary of Research Efforts Focusing on Specific Cancers

→ **Current Research**

Survivorship: Diet and Nutritional Guidance During and Following Cancer Treatment

References

Cancer Medicine → **Section 4: Cancer Epidemiology, Prevention, and Screening** → 29. Nutrition in the Etiology and Prevention of Cancer

Current Research

Specific Foods, Nutrients, and Dietary Components Frequently Associated with Cancer Prevention

Many people at risk of cancer focus their attention upon specific foods or nutrients in part because of the extensive marketing of products and publicity generated by the popular press. This tendency is facilitated by the news media when science reporters publicize results of single studies or preliminary findings, often confusing readers with contradictory and conflicting results. The following section briefly summarizes data regarding selected food components or nutrients and may assist the medical practitioner in responding to specific inquiries from individuals.

[↑ TOP](#)

Vitamins

Vitamin A

Vitamin A is essential for the normal growth and development of epithelial tissues. Vitamin A deficiency is common in many parts of the developing world, but is extremely rare in Americans. Vitamin A is provided in the diet as retinol and its esters, primarily from milk and organ meats, and as β -carotene and a few other provitamin A carotenoids in yellow and leafy green vegetables. Interest in vitamin A and related compounds in the etiology, prevention, and treatment of cancer is rapidly expanding. A protective effect of consuming foods rich in vitamin A has been hypothesized for several types of cancer^{1,2,11,16,18,188}; at this time, however, there is no clear evidence that vitamin A supplementation will decrease the risk of cancer in populations or individuals consuming a healthy diet. Although many studies in laboratory models indicate that vitamin A deficiency increases the susceptibility of

β Carotene

Foods rich in β -carotene, such as many fruits and vegetables, are associated with a lower risk of cancer. However, recent intervention trials with β -carotene clearly question the validity of that hypothesis that the benefits of a plant-based diet can be produced through β -carotene supplements. Two large intervention studies of β -carotene demonstrated a higher risk of lung cancer in smokers.^{124,203} Although β -carotene is a potential antioxidant and source of vitamin A, supplements should be discouraged for cancer prevention and dietary sources should be encouraged. [↑ TOP](#)

risk will be observed
A excess has not
s pharmacologic
s an important area of
s is lacking, the
a large body of

MOLECULAR CELL BIOLOGY

Lodish Berk Zipursky Matsudaira Baltimore Darnell

W H FREEMAN AND COMPANY

[Short Contents](#) | [Full Contents](#)

[Other books @ NCBI](#)

[Molecular Cell Biology](#) → **16. Cellular Energetics: Glycolysis, Aerobic Oxidation, and Photosynthesis** → 16.3. Photosynthetic Stages and Light-Absorbing Pigments

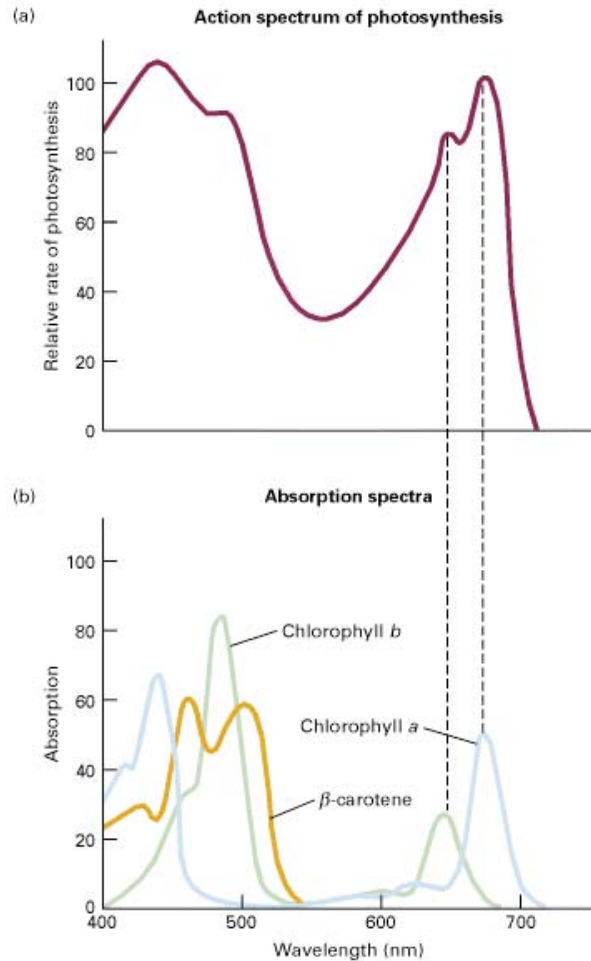


Figure 16-37. Photosynthesis at different wavelengths. (a) The action spectrum of photosynthesis in plants; that is, the ability of light of different wavelengths to support photosynthesis. (b) The absorption spectra for three photosynthetic pigments: chlorophyll *a*, chlorophyll *b*, and β -carotene. Each spectrum shows how well light of different wavelengths is absorbed by one of the pigments. A comparison of the action spectrum with the individual absorption spectra suggests that photosynthesis at 680 nm is primarily due to light absorbed in the

Navigation

[About this book](#)

16. Cellular Energetics: Glycolysis, Aerobic Oxidation, and Photosynthesis

16.1. Oxidation of Glucose and Fatty Acids to CO₂

16.2. Electron Transport and Oxidative Phosphorylation

➔ 16.3. Photosynthetic Stages and Light-Absorbing Pigments

16.4. Molecular Analysis of Photosystems

16.5. CO₂ Metabolism during Photosynthesis

PERSPECTIVES for the Future

PERSPECTIVES in the Literature

Testing Yourself on the Concepts

MCAT/GRE-Style Questions

[References](#)

Search

Go

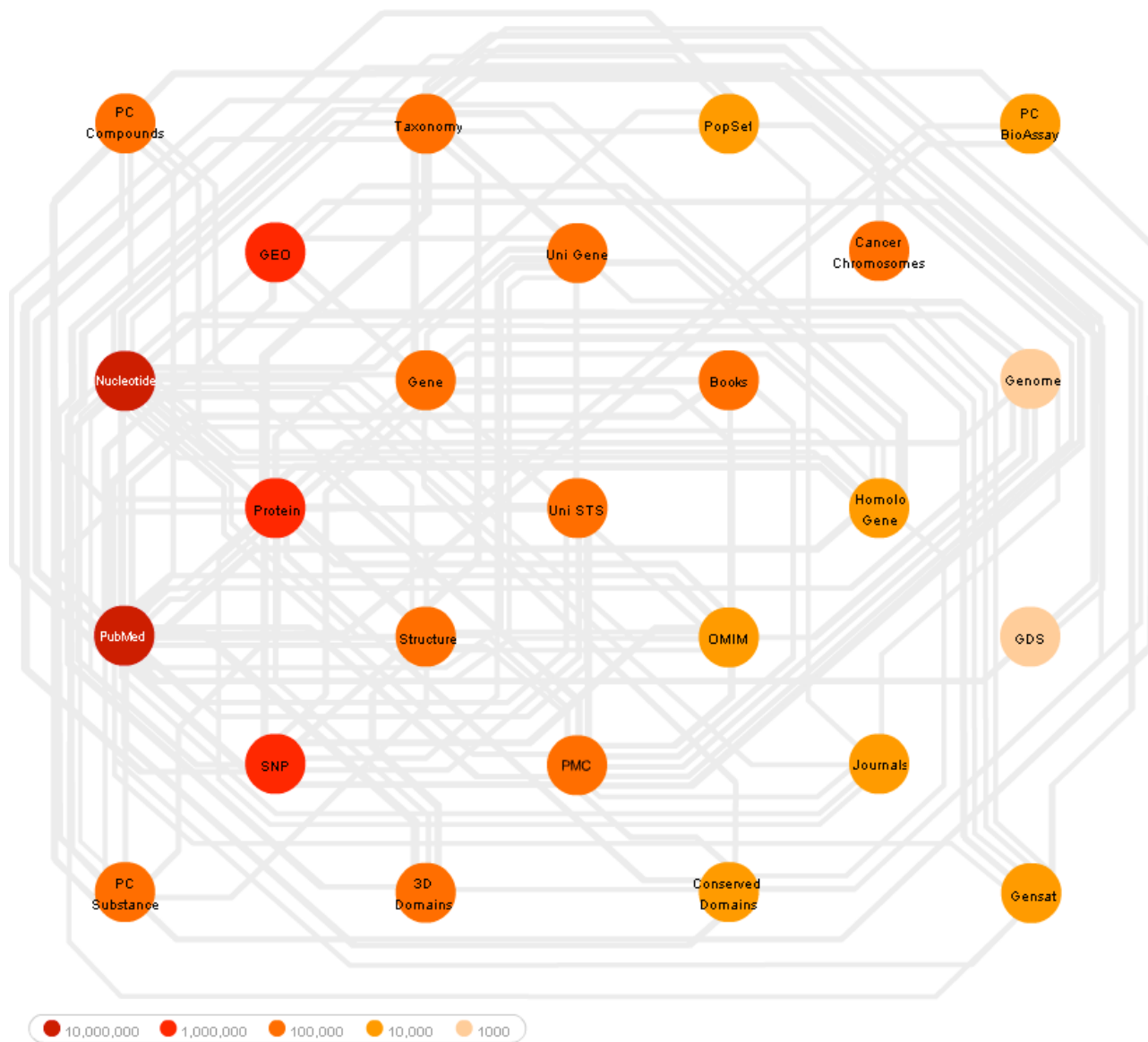
This book

All books

PubMed

You can make up your
own examples, to search
Pubmed...
the Bookshelf...





Web References

- The “About Entrez” page at the NCBI

<http://www.ncbi.nlm.nih.gov/Database/index.html>

- Model of Entrez Databases from NCBI

<http://www.ncbi.nih.gov/Database/datamodel/index.html>

- PubMed Tutorial from NLM

http://www.nlm.nih.gov/bsd/pubmed_tutorial/m1001.html

Credits

- Materials for this presentation have been adapted from the following sources:

NCBI HelpDesk - Field Guide Course Materials

Bioinformatics: A practical guide to the analysis of genes and proteins

- Questions? Please contact:

Dr. Joanne Fox
Michael Smith Laboratories
joanne@mssl.ubc.ca

Let's start at 9:30am

BLAST background, guided tour & practical exercises

