# Nucleic acid sequence analysis

Presented by Guangyong Zheng
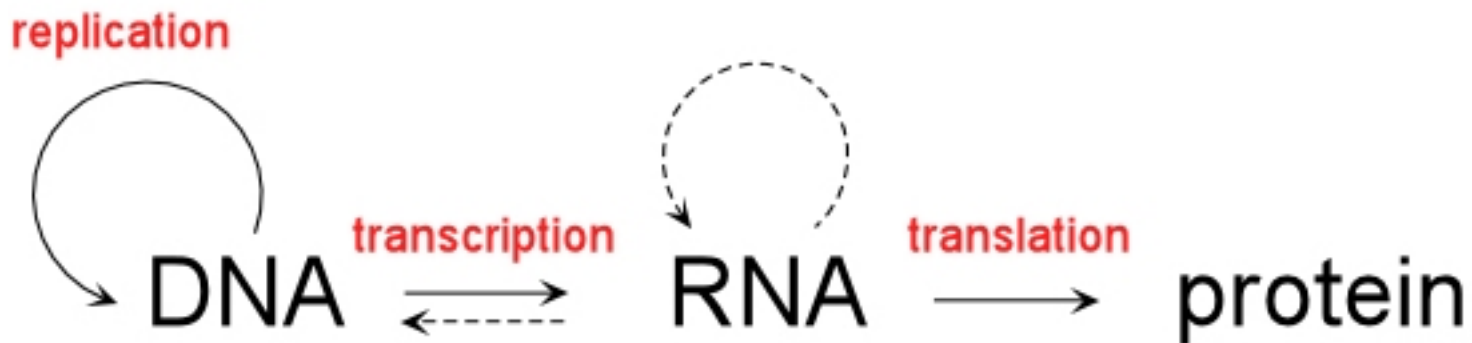
# Outline

## Nucleic acid sequence analysis

➢ Famous database of nucleic acid sequence

➢ Analysis of nucleic acid sequence

➢ BLAST software

- Basic concept of BLAST
- BLAST use - web manners
- BLAST use – local manners

# PART I  Famous database of nucleic acid sequence

# Database of nucleic acid sequence

- Nucleic acid sequence is important material of bioinformatics research.
- Nucleic acid sequences provide the fundamental starting point for describing and understanding the structure, function, and development of genetically diverse organisms.

# Database of nucleic acid sequence

http://www.insdc.org/

# Database of nucleic acid sequence

Feature table

http://www.insdc.org/documents/feature_table.html



INSDC International Nucleotide Sequence Database Collaboration

ABOUT INSDC     POLICY     ADVISORS     DOCUMENTS

DDBJ

ENA
European Nucleotide Archive

NCBI

**The DDBJ/ENA/GenBank Feature Table Definition**

```
The DDBJ/ENA/GenBank Feature Table Definition
Feature Table:
Definition

Version 10.6 November 2016



DNA Data Bank of Japan, Mishima, Japan.
EMBL-EBI, European Nucleotide Archive, Cambridge, UK.
GenBank, NCBI, Bethesda, MD, USA.
```

```
2.1 Format Design

The format design is based on a tabular approach and consists of the following
items:

Feature key - a single word or abbreviation indicating functional group
Location - instructions for finding the feature
Qualifiers - auxiliary information about a feature
```

# Database of nucleic acid sequence

http://www.ncbi.nlm.nih.gov/genbank/

# Database of nucleic acid sequence

http://www.ddbj.nig.ac.jp

# Database of nucleic acid sequence

http://www.ebi.ac.uk/ena/

# Database of nucleic acid sequence

## Submit sequence (http://www.ncbi.nlm.nih.gov/genbank/submit_types)



**NCBI** Resources ☑ How To ☑

**GenBank**    [Nucleotide ▾] [ ]

| GenBank ▾ | Submit ▾ | Genomes ▾ | WGS ▾ | Metagenomes ▾ | TPA ▾ | TSA ▾ | INSDC ▾ | Other ▾ |

### GenBank Submission Types

#### Standard

GenBank accepts mRNA or genomic sequence data directly determined by the submitter. The submission must include information about the source organism and annotation provided by the submitter. More details about adding annotation and sample files can be found in the GenBank Submissions Handbook. If you have any questions about the best method for submitting your data, please contact our user services group at: info@ncbi.nlm.nih.gov.

The following data is not accepted by GenBank:

- Noncontiguous sequences
- Primer sequences
- Protein sequences with no underlying nucleotide submission
- Sequence containing a mix of genomic and mRNA sequence
- Sequences without a physical counterpart (consensus sequences)
- Sequences with length less than 200 nucleotides

Raw sequence reads from next generation sequencing platforms should be submitted to the Sequence Read Archive (SRA).

Sequence data not directly obtained by the submitter may be acceptable for the Third Party Annotation database.

# Getting sequence from the GenBank database

**http://www.ncbi.nlm.nih.gov/nucleotide/**

# Getting sequence from the GenBank database

Data format of nucleic acid (I)

## Homo sapiens heat shock transcription factor 4 (HSF4), transcript variant 2, mRNA

NCBI Reference Sequence: NM_001040667.2

FASTA    Graphics

Go to: ☑

```
LOCUS       NM_001040667            2622 bp    mRNA    linear   PRI 03-MAY-2014
DEFINITION  Homo sapiens heat shock transcription factor 4 (HSF4), transcript
            variant 2, mRNA.
ACCESSION   NM_001040667 XM_005255925
VERSION     NM_001040667.2  GI:194440740
KEYWORDS    RefSeq.
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 2622)
  AUTHORS   Merath K, Ronchetti A and Sidjanin DJ.
  TITLE     Functional analysis of HSF4 mutations found in patients with
            autosomal recessive congenital cataracts
  JOURNAL   Invest. Ophthalmol. Vis. Sci. 54 (10), 6646-6654 (2013)
   PUBMED   24045990
  REMARK    GeneRIF: the transcriptional activation of HSF4 is mediated by
            interactions between activator and repressor domains within the
            C-terminal end.
            Publication Status: Online-Only
```

# Getting sequence from the GenBank database

## Data format of nucleic acid (II)

```
FEATURES                Location/Qualifiers
     source             1..2622
                        /organism="Homo sapiens"
                        /mol_type="mRNA"
                        /db_xref="taxon:9606"
                        /chromosome="16"
                        /map="16q21"
     gene               1..2622
                        /gene="HSF4"
                        /gene_synonym="CTM; CTRCT5"
                        /note="heat shock transcription factor 4"
                        /db_xref="GeneID:3299"
                        /db_xref="HGNC:5227"
                        /db_xref="MIM:602438"
     exon               1..468
                        /gene="HSF4"
                        /gene_synonym="CTM; CTRCT5"
                        /inference="alignment:Splign:1.39.8"
     exon               469..594
                        /gene="HSF4"
                        /gene_synonym="CTM; CTRCT5"
                        /inference="alignment:Splign:1.39.8"
     exon               595..1088
                        /gene="HSF4"
                        /gene_synonym="CTM; CTRCT5"
                        /inference="alignment:Splign:1.39.8"
     misc_feature       849..851
                        /gene="HSF4"
                        /gene_synonym="CTM; CTRCT5"
                        /note="upstream in-frame stop codon"
     STS                898..2569
                        /gene="HSF4"
                        /gene_synonym="CTM; CTRCT5"
                        /db_xref="UniSTS:494796"
     CDS                966..2444
                        /gene="HSF4"
                        /gene_synonym="CTM; CTRCT5"
                        /note="isoform b is encoded by transcript variant 2; heat
                        shock factor protein 4; HSF 4; hHSF4; HSTF 4"
```

# Getting sequence from the GenBank database

Data format of nucleic acid (III)

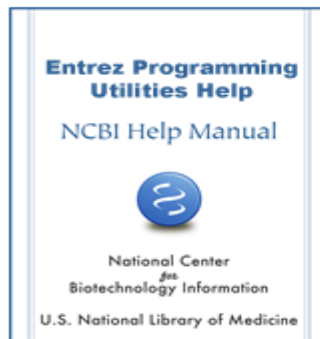## Homo sapiens heat shock transcription factor 4 (HSF4), transcript variant 2, mRNA

NCBI Reference Sequence: NM_001040667.2

GenBank    Graphics

```
>gi|194440740|ref|NM_001040667.2| Homo sapiens heat shock transcription factor 4 (HSF4),
transcript variant 2, mRNA
AGATGCACGCGCGTCCCCGCTGCCCAACGAAGCCTGGGTCGCGTTGCGCCGCCGCCACCCTGGGCTGGCA
GTGGAGCTGGAGCTGGAGCCCGCGCTGCCCGCTGAGAGCGTGACGCGCGTCCTGCAGCCAGCCGTCCCCG
TGGCTGCGCTGCGCCTCAACCTCTCAGGCGACACCGTAGGCCCAGTGCGCTTCGCAGCACACCACTACGC
CGCAACCCTGTGCGCGCTCGAGGTGCGCGCAGCCGCTTCGGCCGAGCTGAACGCCGCGCTGGAGGAGCTG
GCGGCGCGCTGCGCGGCCCTGCGCGAGGTGCATTGTTTCTGCGTGGTGAGCCACTCGGTGCTGGACGCCT
TCCGCGCGCACTGCCCGCGCCTGCGCACCTATACCCTCAAGCTCACGCGCGAGCCGCATCCCTGGAGGCC
TACGCTCGTGGCGTGATTGGGCGACTTCTCTCCCCCGTCCCCGTGGACACGCCCCACCCGCTCGGTCCTG
GACACACTGCCCCCCTCTCTTGCCTCCACCCCTCTGCGGACTCTGCAGCTCCGCGGCCCCGGCGCAGGGA
GAGGGAGGGCACGGGCGCGGGCCGGGCCTCAAGGGACTTGCCCAGCCCACACCAGGTCGCGCACCGGCGA
TTTCTCCTGTAGAACAAAGAAGGAAATAGAGGGACCGAGAGGGGTGGGACTCGAACCCAAGTCTCCCACT
CATCTCACCCCACCCCACCCCACTCCACCCCACCCCACCCCTCCACTCCACTCCACTCCACACATCCCAT
CCAGCCAGCCTTTTCTGCCTGCTGGTGCCTCGGCCGCTGTCCGAGCCCCGCCCCGCGGGCTTGCACGTGG
CCCCCGCCTGACCCGGCGCCCCGGGGCGGAGTGGGCGGAGCGGGCGGGCGGGCAAACGCAGCACTTTCCGCGGC
TTTGACGAGCCCGCAGCGGCCGGGCCCGAGCGCAGAGCCGGGCCGAGACTGCACCATGCAGGAAGCGCCA
GCTGCGCTGCCCACGGAGCCAGGCCCCAGCCCCGTGCCTGCCTTCCTCGGCAAGCTATGGGCGCTGGTGG
GGGACCCAGGCACAGACCACCTGATCCGCTGGAGCCCGAGCGGGACCAGTTTCCTCGTAAGCGACCAGAG
CCGTTTCGCCAAGGAAGTGCTGCCCCAGTATTTCAAGCATAGCAACATGGCGAGCTTCGTGCGCCAACTC
AACATGTACGGTTTTCGGAAGGTGGTGAGCATCGAGCAGGGCGGCCTGCTTAGGCCGGAGCGCGACCACG
TCGAGTTCCAGCACCCGAGCTTCGTGCGCGGCCGCGAGCAGCTACTGGAGCGCGTGCGGCGCAAGGTGCC
CGCGCTGCGCGGCGACGACGGCCGCTGGCGCCCGGAGGACCTGGGTCGACTACTGGGCGAGGTGCAGGCT
TTGCGGGGAGTGCAGGAGAGCACCGAGGCGCGGCTGCGGGAGCTCAGGCAGCAGAACGAGATCTTGTGGC
GGGAGGTGGTGACACTTCGGCAGAGCCACGGTCAGCAGCACCGGGTCATTGGCAAGCTGATCCAGTGTCT
CTTTGGGCCACTTCAGGCGGGGCCGAGCAATGCAGGAGGCAAGAGAAAGCTGTCCCTGATGCTGGATGAG
GGGAGCTCATGCCCAACACCTGCCAAGTTCAACACCTGCCCTCTACCTGGTGCCCTTCTGCAGGACCCCT
ACTTCATCCAGTCGCCTCTCCCAGAGACAAATTTGGGCCTTAGCCCTCACAGGGCCAGGGGCCCCATCAT
CTCTGACATCCCAGAAGACTCTCCATCCCCTGAGGGGACCAGGCTTTCTCCCTCCAGTGATGGCAGGAGG
GAGAAGGGCCTGGCACTGCTCAAAGAAGAGCCGGCCAGTCCAGGGGGGGATGGCGAGGCCGGGCTGGCCC
TGGCCCCAAACGAGTGTGACTTCTGCGTGACAGCCCCCCCGCCCACTGCCTGTGGCTGTGGTGCAGGCCAT
CCTGGAAGGGAAAGGGAAGCTTCAGCCCCGAGGGGCCCAGGAATGCCCAACAGCCTGAACCAGGGGATCCC
AGGGAGATACCTGACAGGGGGCCTCTGGGCCTGGAAAGCGGGGACAGGAGCCCAGAGAGTCTGCTGCCTC
CGATGCTGCTTCAGCCCCCTCAAGAAAGTGTGGAACCTGCAGGGCCTCTAGATGTGCTGGGCCCCAGTCT
CCAAGGGCGAGAATGGACCCTGATGGACTTGGACATGGAGCTGTCCTTGATGCAGCCCTTGGTTCCAGAG
CGGGGTGAGCCTGAGCTGGCGGTCAAGGGGTTAAATTCTCCAAGCCCAGGGAAGGACCCCACGCTCGGGG
CCCCACTCCTGCTGGATGTCCAGGCGGCCTTGGGAGGCCCAGCCTGGGCCTGCCTGGGGCTTTAACCAT
TTATAGCACTCCTGAGAGCCGGACTGCCTCCTACTTGGGCCCGGAAGCCAGTCCCTCCCCCTAAGACCCC
GCGCCTCTGAAGGGGCTTGGAACCAGTCCGCCGCTGCACATCCTTCTTGGCTTCCTGGCGCCCCCTATCG
GGGGTGAGCGAAGCCCCCACTACTAAATGGCCTCTCTCCACTACCCCGACTATCCCTGCACATAAACTCC
GTTTTTTTTTTTTCAAAAAAAAAAAAAAAAAAA
```

# Getting sequence from the GenBank database

**http://www.ncbi.nlm.nih.gov/books/NBK25501/**



**Entrez Programming Utilities Help**

< Prev    Next >

Bethesda (MD): National Center for Biotechnology Information (US); 2010-.

Copyright and Permissions

Search this book

**Introduction to the E-utilities**

- You Tube E-utilities Introduction

- Please see the Release Notes for details and changes.

The Entrez Programming Utilities (E-utilities) are a set of eight server-side programs that provide a stable interface into the Entrez query and database system at the National Center for Biotechnology Information (NCBI). The E-utilities use a fixed URL syntax that translates a standard set of input parameters into the values necessary for various NCBI software components to search for and retrieve the requested data. The E-utilities are therefore the structured interface to the Entrez system, which currently includes 38 databases covering a variety of biomedical data, including nucleotide and protein sequences, gene records, three-dimensional molecular structures, and the biomedical literature.

# Getting sequence from the GenBank database

## EFetch

### Base URL

http://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi

### Functions

- Returns formatted data records for a list of input UIDs
- Returns formatted data records for a set of UIDs stored on the Entrez History server

### Required Parameters

**db**

Database from which to retrieve records. The value must be a valid Entrez database name (default = pubmed). Currently EFetch does not support all Entrez databases. Please see Table 1 in Chapter 2 for a list of available databases.

### Required Parameter – Used only when input is from a UID list

**id**

UID list. Either a single UID or a comma-delimited list of UIDs may be provided. All of the UIDs must be from the database specified by **db**. There is no set maximum for the number of UIDs that can be passed to EFetch, but if more than about 200 UIDs are to be provided, the request should be made using the HTTP POST method.

```
efetch.fcgi?db=protein&id=15718680,157427902,119703751
```

# Getting sequence from the GenBank database

## Nucleotide/Nuccore

Fetch the first 100 bases of the plus strand of GI 21614549 in FASTA format:

http://eutils.ncbi.nlm.nih.gov/entrez/eutils/efetch.fcgi?db=nuccore&id=21614549&strand=1&seq_start=1&
seq_stop=100&rettype=fasta&retmode=text

```
sequence.fasta  ✕

     0        1,0        2,0        3,0        4,0        5,0        6,0        7,0
1 >gi|21614549:1-100 Homo sapiens IL2-inducible T-cell kinase (ITK), mRNA
2 TGCATTCTTTGCCCCAAAACTCTTTCCTTTGGTTGTGCTAAGAGGTGATGCCCAAGGTGCACCACCTTTC
3 AAGAACTGGATCATGAACAACTTTATCCTC
4
5 |
```

# PART II Analysis of nucleic acid sequence

# Analysis of nucleic acid sequence

(1) Getting basic information of a nucleic acid sequence

(2) Primer design

(3) Tow sequences alignment

(4) Multi sequences alignment

(5) Finding open reading frame of a nucleic acid sequence

(6) Gene prediction

(7) Sequence localization in genome

(8) Sequence assembly

# Obtain basic information of nucleic acid sequence

**Software: BioEdit (http://www.mbio.ncsu.edu/bioedit/bioedit.html)**
**menu: File -> open -> sequence -> nucleic acid -> nucleotide composition**

# Basic transition of nucleic acid sequence

- DNA -> RNA
- Sequence -> Reverse complement
- DNA -> protein

# Basic transition of nucleic acid sequence

**menu: File -> open -> sequence -> nucleic acid -> DNA-RNA**

# Basic transition of nucleic acid sequence

**menu: File -> open -> sequence -> nucleic acid -> Reverse Complement**

# Basic transition of nucleic acid sequence

**menu: File -> open -> sequence -> nucleic acid -> translate**

# Analysis of enzyme mapping

**menu: File -> open -> sequence -> nucleic acid -> Restriction Map**

# Primer design

**NCBI website (http://www.ncbi.nlm.nih.gov/tools/primer-blast/)**

# Primer design

# Primer design

## Primer3 (http://primer3.wi.mit.edu/)

Primer3web version 4.0.0 - Pick primers from a DNA sequence.

disclaimer | code
cautions

Select the Task for primer selection [generic ▾]

Paste source sequence below (5'->3', string of ACGTNacgtn -- other letters treated as N -- numbers and blanks ignored). FASTA format ok. Please N-out undesirable sequence (vector, ALUs, LINEs, etc.) or use a Mispriming Library (repeat library) [NONE ▾]

☑ Pick left primer, or use left primer below | ☐ Pick hybridization probe (internal oligo), or use oligo below | ☑ Pick right primer, or use right primer below (5' to 3' on opposite strand)

[Pick Primers] [Download Settings] [Reset Form]

## Primer3 codes (http://sourceforge.net/projects/primer3/)

Home / Browse / Science & Engineering / Bio-Informatics / Primer3 – PCR primer design tool

# Primer3 – PCR primer design tool
Brought to you by: brantfaircloth, steverozen, untergasser

Summary | Files | Reviews | Support | Wiki | Feature Requests | News | Code

★ 5.0 Stars (14)
↓ 255 Downloads (This Week)
📅 Last Update: 2013-10-28

Download
primer3-src-2.3.6.tar.gz

Browse All Files

## Description

Design PCR primers from DNA sequence. Widely used (190k Google hits for "primer3"). From mispriming libraries to sequence quality data to the generation of internal oligos, primer3 does it. C&perl. Developers/testers/documenters needed.

# Sequences alignment (two sequence)

NCBI homepage -> Blast homepage -> specialized searches -> Global Align

# Sequences alignment (two sequence)

# Sequences comparison (multi-sequences)

**http://www.ebi.ac.uk/Tools/msa/clustalo/**

# Sequences comparison (multi-sequences)

# Sequences comparison (multi-sequences)

**Software: BioEdit**
 **File -> open -> sequence -> Accessory Application -> clustalw**

# Sequences comparison (multi-sequences)

**http://www.clustal.org/clustal2/**

# Sequences comparison (multi-sequences)

Software: clustalw / clustalX        menu : load sequence

# Sequences comparison (multi-sequences)

Software: clustalw / clustalX        menu : Alignment -> do complete alignment

# Analysis of open reading frame

https://www.ncbi.nlm.nih.gov/orffinder/

# Analysis of open reading frame

# Gene prediction

[http://genes.mit.edu/GENSCAN.html](http://genes.mit.edu/GENSCAN.html)

# Gene prediction

**GENSCAN Output**

View gene model output: PS | PDF

GENSCAN 1.0     Date run: 29-Aug-110     Time: 05:08:50

Sequence /tmp/08_29_10-05:08:50.fasta : 2622 bp : 66.32% C+G : Isochore 4 (57 - 100 C+G%)

Parameter matrix: HumanIso.smat

Predicted genes/exons:

| Gn.Ex | Type | S | .Begin | ...End | .Len | Fr | Ph | I/Ac | Do/T | CodRg | P.... | Tscr.. |
|-------|------|---|--------|--------|------|----|----|------|------|-------|-------|--------|
| 1.01 | Term | + | 72 | 436 | 365 | 1 | 2 | 15 | 38 | 583 | 0.999 | 41.79 |
| 1.02 | PlyA | + | 643 | 648 | 6 | | | | | | | -3.64 |
| 2.00 | Prom | + | 905 | 944 | 40 | | | | | | | -12.52 |
| 2.01 | Sngl | + | 966 | 2444 | 1479 | 2 | 0 | 91 | 47 | 1338 | 0.993 | 124.01 |
| 2.02 | PlyA | + | 2605 | 2610 | 6 | | | | | | | -0.45 |

Suboptimal exons with probability > 1.000

| Exnum | Type | S | .Begin | ...End | .Len | Fr | Ph | B/Ac | Do/T | CodRg | P.... | Tscr.. |
|-------|------|---|--------|--------|------|----|----|------|------|-------|-------|--------|

NO EXONS FOUND AT GIVEN PROBABILITY CUTOFF

Predicted peptide sequence(s):

>/tmp/08_29_10-05:08:50.fasta|GENSCAN_predicted_peptide_1|121_aa
XELELEPALPAESVTRVLQPAVPVAALRLNLSGDTVGFVRFAAHHYAATLCALEVRAAAS
AELNAALEELAARCAALREVHCFCVVSHSVLDAFRAHCPRLRTYTLKLTREPHPWRPTLV
A

>/tmp/08_29_10-05:08:50.fasta|GENSCAN_predicted_peptide_2|492_aa
MQEAPAALPTEPGPSPVPAFLGKLWALVGDFGTDHLIRWSPSGTSFLVSDQSRFAKEVLF
QYFKHSNMASFVRQLNMYGFRKVVSIEQGGLLRPERDHVEFQHPSFVRGREQLLERVRRK
VPALRGDDGRWRPEDLGRLLGEVQALRGVQESTEARLRELRQQNEILWREVVTLRQSHGQ
QHRVIGKLIQCLFGPLQAGPSNAGGKRKLSLMLDEGSSCPTPAKFNTCPLPGALLQDPYF
IQSPLPETNLGLSPHRARGPIISDIPEDSPSPEGTRLSPSSDGRREKGLALLKEEPASPG
GDGEAGLALAPNECDFCVTAPPPLPVAVVQAILEGKGSFSPEGPRNAQQPEPGDPREIPD
RGPLGLESGDRSPESLLPPMLLQPPQESVEPAGPLDVLGPSLQGREWTLMDLDMELSLMQ
PLVPERGEPELAVKGLNSPSPGKDPTLGAPLLLDVQAALGGPALGLPGALTIYSTPESRT
ASYLGPEASPSP

Back to GENSCAN

# Sequence localization in genome

**NCBI homepage -> Blast homepage -> BLAST Genomes -> Human**

# Sequence localization in genome

# Sequence localization in genome

Blast result -> view report -> human genome view

# Sequence localization in genome

# Sequence assembly

**Software BioEdit**

**File -> open -> sequence -> Accessory Application -> contig assembly**

# Sequence assembly

**CAP3 (contig assembly program)**
**http://seq.cs.iastate.edu/cap3.html**

## CAP3 Assembly Program

- A version of CAP3 for a 32-bit Linux system with an Intel processor [download tar file](download tar file)

- A version of CAP3 for a 64-bit Linux system with an Opterron processor: [download tar file](download tar file)

- A version of CAP3 for a 64-bit Linux system with an Intel processor: [download tar file](download tar file)

- A version of CAP3 for an old version (2009) of 64-bit Linux system with an Intel processor: [download tar file](download tar file)

- A version of CAP3 for a 32-bit MacOSX system with an Intel processor: [download tar file](download tar file)

- A version of CAP3 for a 64-bit MacOSX system with an Intel processor: [download tar file](download tar file)

- A version of CAP3 for a 64-bit Solaris system with an Opterron processor: [download tar file](download tar file)

- A version of CAP3 for a 32-bit Cygwin simulator on Windows: [download tar file](download tar file)
You need to download and install: [Cygwin](Cygwin)

# PART III  insight into BLAST

# Concept of BLAST

**Concept:** The Basic Local Alignment Search Tool (BLAST) finds regions of local similarity between sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance of matches. BLAST can be used to infer functional and evolutionary relationships between sequences as well as help identify members of gene families.

**Purpose:** predict function

infer evolutional tree

construct families

# Concept of BLAST

BLAST :  word size -> **high-scoring segment (HSP)**

**query sequence** :  (nucleotide/translated nucleotide, protein)

**database**: (nucleotide/translated nucleotide, protein)

# BLAST HOME

**http://blast.ncbi.nlm.nih.gov/Blast.cgi**

# Basic program of BLAST

| | |
|---|---|
| **nucleotide blast** | Search a **nucleotide** database using a **nucleotide** query<br>*Algorithms*: blastn, megablast, discontiguous megablast |
| **protein blast** | Search **protein** database using a **protein** query<br>*Algorithms*: blastp, psi-blast, phi-blast |
| **blastx** | Search **protein** database using a **translated nucleotide** query |
| **tblastn** | Search **translated nucleotide** database using a **protein** query |
| **tblastx** | Search **translated nucleotide** database using a **translated nucleotide** query |

# PSI-BLAST

Position-Specific Iterated BLAST (PSI-BLAST, family protein)

(1) PSI-BLAST takes as an input a single protein sequence and compares it to a protein database, using the gapped BLAST program

(2) The program constructs a multiple alignment, and then a profile, from any significant local alignments found.

(3) The profile is compared to the protein database, again seeking local alignments.

(4) PSI-BLAST estimates the statistical significance of the local alignments found.

(5) Finally, PSI-BLAST iterates, by returning to step (2), an arbitrary number of times or until convergence.

# Summary of BLAST

- **Traditional BLAST (formerly blastall) nucleotide, protein, translations**
    - **blastn** nucleotide query  vs. nucleotide database
    - **blastp** protein query  vs. protein database
    - **blastx** nucleotide query  vs. protein database
    - **tblastn** protein query  vs. translated nucleotide database
    - **tblastx** translated nucleotide query  vs. translated nucleotide database

- **Position Specific BLAST Programs protein only**
    - **Position Specific Iterative BLAST (PSI-BLAST)**
      **Automatically generates a position specific score matrix (PSSM)**

# BLAST use – web means (nucleotide)

# BLAST use – web means (nucleotide)

# BLAST use – web means (taxonomy report)

# BLAST use – infer evolution tree

# BLAST use – predict function

# BLAST use – web manners (alignment report)

# BLAST use – web manners (protein)

# BLAST use – result summary

# BLAST use – multi alignment

# BLAST use – conserved domain

# BLAST use

## Batch BLAST jobs

(1)  input "batches" of sequences into one form and retrieve the results

# BLAST use

Batch BLAST jobs

(2) Utilize the standalone BLAST binaries.

You can retrieve BLAST execute files from NCBI ftp sites
ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/LATEST/

# BLAST+

**ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/LATEST/**

# BLAST use – command line means

List of the BLAST directory

(1) bin – binary files for BLAST execute

(2) doc – documents of BLAST

(3) ncbi_package_info – version information of BLAST

# BLAST+

BLAST use

(1) Make a formatted database to use

execute command : **makeblastdb**

input: fasta format sequences (database sequences)
output: formatted database , used by BLAST program

# BLAST+

**Argument of makeblastdb**

```
*** Input options
-in <File_In>
  Input file/database name
  Default = `-'
-input_type <String, `asn1_bin', `asn1_txt', `blastdb', `fasta'>
  Type of the data specified in input_file
  Default = `fasta'


*** Output options
-out <String>
  Name of BLAST database to be created
  Default = input file name provided to -in argumentRequired if multiple
  file(s)/database(s) are provided as input
-max_file_sz <String>
  Maximum file size for BLAST database files
  Default = `1GB'
```

# BLAST+

BLAST use
(2) Carry out BLAST program

 execute command : **blastn, blastp**

input: fasta sequences (query sequences), database
output : query result file

# BLAST+

## Usage of blastn

```
blastn [-h] [-help] [-import_search_strategy filename]
   [-export_search_strategy filename] [-task task_name] [-db database_name]
   [-dbsize num_letters] [-gilist filename] [-seqidlist filename]
   [-negative_gilist filename] [-entrez_query entrez_query]
   [-db_soft_mask filtering_algorithm] [-db_hard_mask filtering_algorithm]
   [-subject subject_input_file] [-subject_loc range] [-query input_file]
   [-out output_file] [-evalue evalue] [-word_size int_value]
   [-gapopen open_penalty] [-gapextend extend_penalty]
   [-perc_identity float_value] [-xdrop_ungap float_value]
   [-xdrop_gap float_value] [-xdrop_gap_final float_value]
   [-searchsp int_value] [-max_hsps int_value] [-sum_statistics]
   [-penalty penalty] [-reward reward] [-no_greedy]
   [-min_raw_gapped_score int_value] [-template_type type]
   [-template_length int_value] [-dust DUST_options]
   [-filtering_db filtering_database]
   [-window_masker_taxid window_masker_taxid]
   [-window_masker_db window_masker_db] [-soft_masking soft_masking]
   [-ungapped] [-culling_limit int_value] [-best_hit_overhang float_value]
   [-best_hit_score_edge float_value] [-window_size int_value]
   [-off_diagonal_range int_value] [-use_index boolean] [-index_name string]
   [-lcase_masking] [-query_loc range] [-strand strand] [-parse_deflines]
   [-outfmt format] [-show_gis] [-num_descriptions int_value]
   [-num_alignments int_value] [-html] [-max_target_seqs num_sequences]
   [-num_threads int_value] [-remote] [-version]
```

# BLAST+

**Usage of blastp**

```
blastp [-h] [-help] [-import_search_strategy filename]
   [-export_search_strategy filename] [-task task_name] [-db database_name]
   [-dbsize num_letters] [-gilist filename] [-seqidlist filename]
   [-negative_gilist filename] [-entrez_query entrez_query]
   [-db_soft_mask filtering_algorithm] [-db_hard_mask filtering_algorithm]
   [-subject subject_input_file] [-subject_loc range] [-query input_file]
   [-out output_file] [-evalue evalue] [-word_size int_value]
   [-gapopen open_penalty] [-gapextend extend_penalty]
   [-xdrop_ungap float_value] [-xdrop_gap float_value]
   [-xdrop_gap_final float_value] [-searchsp int_value] [-max_hsps int_value]
   [-sum_statistics] [-seg SEG_options] [-soft_masking soft_masking]
   [-matrix matrix_name] [-threshold float_value] [-culling_limit int_value]
   [-best_hit_overhang float_value] [-best_hit_score_edge float_value]
   [-window_size int_value] [-lcase_masking] [-query_loc range]
   [-parse_deflines] [-outfmt format] [-show_gis]
   [-num_descriptions int_value] [-num_alignments int_value] [-html]
   [-max_target_seqs num_sequences] [-num_threads int_value] [-ungapped]
   [-remote] [-comp_based_stats compo] [-use_sw_tback] [-version]
```

# THANK YOU