# Sequence Based Function Annotation

**Qi Sun**

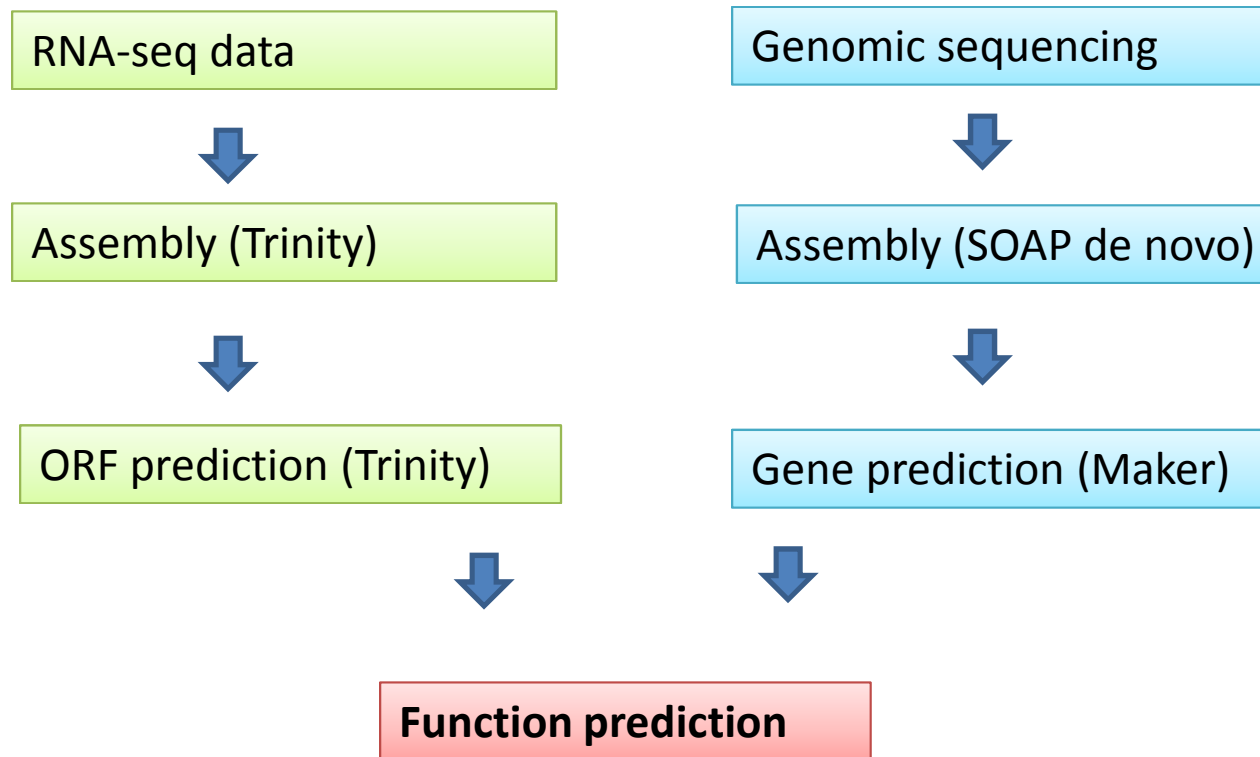**Bioinformatics Facility**

**Biotechnology Resource Center**

**Cornell University**

# Sequence Based Function Annotation

1. Given a sequence, how to predict its biological function?

2. How to describe the function of a gene?

3. How to works with 50,000 genes?

**1. Given a protein sequence, how to predict its function?**

```
>unknow_protein_1
MVHLTDAEKAAVSCLWGKVNSDEVGGEALGRLLVVYPWTQR
YFDSFGDLSSASAIMGNAKVKAHGKKVITAFNDGLNHLDSL
KGTFASLSELHCDKLHVDPENFRLLGNMIVIVLGHHLGKDF
TPAAQAAFQKVVAGVATALAHKYH
```

## Common approaches

- **Identify the homologous gene in a different species with good function annotation; (BLAST, et al.)**

- **Identify conserved motif; (PFAM, InterProScan, et al.)**

## Alternative approaches

- **Protein 3D structure prediction (threading methods)**

- **Co-expression network modules; (Genevestigator)**

- **Linkage or association mapping;**
- **…**

# NCBI BLAST

- **How does BLAST work?**

- **BLAST and Psi-BLAST: Position independent and position specific scoring matrix.**



```
Score =  176 bits (447),  Expect = 4e-50, Method: Compositional matrix adjust.
Identities = 98/232 (42%), Positives = 139/232 (60%), Gaps = 14/232 (6%)

Query  30   MAKVLTLELYKKLRDKETPSGFTVDDVIQTGV--DNPGHPFIMTVGCVAGDEESYEVFKE   87
            + K LT +L+++ +D+      GF+      I +G   N G       VG  AG  +SY  F
Sbjct  26   LQKCLTKDLWEQCKDRRDKYGFSFKQAIFSGSKWTNSG------VGVYAGSHDSYYAFAP   79

Query  88   LFDPIISDRHGGYKPTDKHKTDLNHENLKGG---DDLDPNYVLSSRVRTGRSIKGYTLPP  144
            D II   HG +KP+DKH + ++++ L      D D   + S+R+R  R++      L
Sbjct  80   FMDKIIEAYHG-HKPSDKHISSMDYKQLNCPPFPADED-KMINSTRIRVARNLAADPLGT  137

Query  145  HCSRGERRAVEKLSVEALNSLTGEFKGKYYPLKSMTEKEQQQLIDDHFLFDKPVSPLLLA  204
            +R ER+ +E L   AL   TGE KGKYY L++M++ E++QLI DHFLF K      L +
Sbjct  138  AVTRKERKEIEHLVTSALGEFTGELKGKYYSLETMSDAEKKQLIADHFLF-KGGDKYLQS  196

Query  205  SGMARDWPDARGIWHNDNKSFLVWVNEEDHLRVISMEKGGNMKEVFRRFCVG  256
            +G+ RDWP+ARGI+HND K+FLVWVNEED LR+ISM+ G N+ EVF+R  V
Sbjct  197  AGLERDWPEARGIFHNDAKTFLVWVNEEDQLRIISMQAGSNILEVFKRLSVA  248
```

# How does BLAST work

## Step 1. Create alignments between HSPs (High-scoring Segment Pair)



The BLAST Search Algorithm

query word (W = 3)

Step1  Query:  TGSQSLAALLNKCKTPQGQRLVNQWIKQPLMDKNRIEERLNLVEAFV

Step2  neighborhood words

| PQG | 18 |
|-----|----|
| PEG | 15 |
| PRC | 14 |
| PKG | 14 |
| PNG | 13 |
| PDG | 13 |
| PHG | 13 |
| PMG | 13 |
| PSG | 13 |

neighborhood score threshold (T = 13)

| PQA | 12 |
|-----|----|
| PQN | 12 |

etc...

Step3

Query:   325 SLAALLNKCKTPQGQRLVNQWIKQPLMDKNRIEERLNLVEA 365
         +LA++L+    TP  G  R++  +W+    P+  D    + ER    +  A
Subject: 290 TLASVLDCTVTPMGSRMLKRWLHMPVRDTRVLLERQQTIGA 330

High-scoring Segment Pair (HSP)

**How does BLAST work**

**Step 2. Score each alignment, and report the top alignments**

Number of Chance Alignments = $2 \times 10^{-73}$

```
Score =  288 bits (318),  Expect = 2e-73
Identities = 262/325 (81%), Gaps = 8/325 (2%)
Strand=Plus/Plus

Query  1923   TCAGCCTACCATGAGAATAAGAGAAAGA-AAATGAAGATCAAAAGCTTATTCATCTGTTT  1981
              ||||  |||| |||||||||||||||||| ||||||||| |  |  | |||| |||| |||
Sbjct  33774  TCAGACTACCCTGAGAATAAGAGAAAGAGAGAAATGAAGACCTAGA-CTTATCCATCTCTTT  33832

Query  1982   TTCTTTTTCGTTGGTGTAAAGCCAACACCCTGTCTAAAAAACATAAATTTCTTTAATCAT  2041
              |||| 
```

Match=+2

```
Sbjct         TG
```

Mismatch=-3

```
                                                           ACAAATTTCTTTAAATAT  33892

Query  2042   TTTGCCTCTTTTC
              ||||||||||||||
                                                           AGAATCTAATAGAGTGGT  2100
Sbjct  33893  TTTGCCTCTTTTCTCTGTGCTACAATTAATAAAAAAATGAAAAGAATCTAATTTAATTGT  33952

Query  2101   ACAGCACTGTTA-T
              |  |||||| |
```

Gap

$-(5 + 4(2)) = -13$

```
Sbjct  33953  CTATGACTGTTATT                                   GGTTCTATGA  34012

Query  2160   AAGTTCCAGTGTTC                                   TTGTGGGCTA  2219
              || |||||| | |||
Sbjct  34013  AAATTCCACTATTCTCTCTTTCGCTATTTCAATGGAGGACTTCTAGTTCCTTCTGGATTA  34072

Query  2220   AT----TAAATAAATCATTAATACT   2240
              ||     |||| || ||||||||||
Sbjct  34073  ATTGCATAAAAGAAACATTAATACT   34097
```

# BLOSUM62, a position independent matrix

| | Ala | Arg | Asn | Asp | Cys | Gln | Glu | Gly | His | Ile | Leu | Lys | Met | Phe | Pro | Ser | Thr | Trp | Tyr | Val |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Ala | 4 | | | | | | | | | | | | | | | | | | | |
| Arg | -1 | 5 | | | | | | | | | | | | | | | | | | |
| Asn | -2 | 0 | 6 | | | | | | | | | | | | | | | | | |
| Asp | -2 | -2 | 1 | 6 | | | | | | | | | | | | | | | | |
| Cys | 0 | -3 | -3 | -3 | 9 | | | | | | | | | | | | | | | |
| Gln | -1 | 1 | 0 | 0 | -3 | 5 | | | | | | | | | | | | | | |
| Glu | -1 | 0 | 0 | 2 | -4 | 2 | 5 | | | | | | | | | | | | | |
| Gly | 0 | -2 | 0 | -1 | -3 | -2 | -2 | 6 | | | | | | | | | | | | |
| His | -2 | 0 | 1 | -1 | -3 | 0 | 0 | -2 | 8 | | | | | | | | | | | |
| Ile | -1 | -3 | -3 | -3 | -1 | -3 | -3 | -4 | -3 | 4 | | | | | | | | | | |
| Leu | -1 | -2 | -3 | -4 | -1 | -2 | -3 | -4 | -3 | 2 | 4 | | | | | | | | | |
| Lys | -1 | 2 | 0 | -1 | -3 | 1 | 1 | -2 | -1 | -3 | -2 | 5 | | | | | | | | |
| Met | -1 | -1 | -2 | -3 | -1 | 0 | -2 | -3 | -2 | 1 | 2 | -1 | 5 | | | | | | | |
| Phe | -2 | -3 | -3 | -3 | -2 | -3 | -3 | -3 | -1 | 0 | 0 | -3 | 0 | 6 | | | | | | |
| Pro | -1 | -2 | -2 | -1 | -3 | -1 | -1 | -2 | -2 | -3 | -3 | -1 | -2 | -4 | 7 | | | | | |
| Ser | 1 | -1 | 1 | 0 | -1 | 0 | 0 | 0 | -1 | -2 | -2 | 0 | -1 | -2 | -1 | 4 | | | | |
| Thr | 0 | -1 | 0 | -1 | -1 | -1 | -1 | -2 | -2 | -1 | -1 | -1 | -1 | -2 | -1 | 1 | 5 | | | |
| Trp | -3 | -3 | -4 | -4 | -2 | -2 | -3 | -2 | -2 | -3 | -2 | -3 | -1 | 1 | -4 | -3 | -2 | 11 | | |
| Tyr | -2 | -2 | -2 | -3 | -2 | -1 | -2 | -3 | 2 | -1 | -1 | -2 | -1 | 3 | -3 | -2 | -2 | 2 | 7 | |
| Val | 0 | -3 | -3 | -3 | -1 | -2 | -2 | -3 | -3 | 3 | 1 | -2 | 1 | -1 | -2 | -2 | 0 | -3 | -1 | 4 |

**How does BLAST work**

**Step 2. Score each alignment – protein alignment**

Number of Chance Alignments = $4 \times 10^{-50}$

```
 Score =  176 bits (447),  Expect = 4e-50, Method: Compositional matrix adjust.
 Identities = 98/232 (42%), Positives = 139/232 (60%), Gaps = 14/232 (6%)

Query  30    MAKVLTLELYKKLRDKETPSGFTVDDVIQTGV--DNPGHPFIMTVGCVAGDEESYEVFKE   87
             + K LT +L+++ +D+    GF+    I +G    N G       VG  AG  +SY  F
Sbjct  26    LQKCLTKDLWEQCKDRRDKYGFSFKQAIFSGSKWTNSG------VGVYAGSHDSYYAFAP   79

Query  8     RH                LI          L
Sbjct  8     YH                MI          E

Query  1     AV                GI          E
             +R ER+ +E L    AL    TGE  KGKYY  L++M++
Sbjct  138   AVTRKERKEIEHLVTSALGEFTGELKGKYYSLETMSD

Query  205   SGMARDWPDARGIWHNDNKSFLVWVNEEDHLRVISMEKGGNMKEVFRRFCVG   256
             +G+  RDWP+ARGI+HND K+FLVWVNEED LR+ISM+ G N+ EVF+R  V
Sbjct  197   AGLERDWPEARGIFHNDAKTFLVWVNEEDQLRIISMQAGSNILEVFKRLSVA   248
```

| K    | K    | Q    | Gap |
|------|------|------|-----|
| K +5 | E +1 | F −3 | −(11 + 6(1)) = −18 |

Scores from BLOSUM62, a position independent matrix

- NCBI Discovery Workshops

# BLOSUM62 substitution score is position independent



```
Score =  176 bits (447),  Expect = 4e-50, Method: Compositional matrix adjust.
Identities = 98/232 (42%), Positives = 139/232 (60%), Gaps = 14/232 (6%)

Query   30    MAKVLTLELYKKLRDKETPSGFTVDDVIQTGV--DNPGHPFIMTVGCVAGDEESYEVFKE    87
              + K LT +L+++ +D+     GF+     I +G    N G        VG   AG  +SY   F
Sbjct   26    LQKCLTKDLWEQCKDRRDKYGFSFKQAIFSGSKWTNSG------VGVYAGSHDSYYAFAP    79

Query   88    LFDPIISDI HGYKPTDI HITDLNHENLKGG--DILDPNYVLSSRVRTGRSIKGYTLPP   144
              D  II    H  +KP+DI H + ++++  L       D D   + S+R+R   R++       L
Sbjct   80    FMDKIIEAYHG-HKPSDIHJSSMDYKQLNCPPFIADED-KMINSTRIRVARNLAADPLGT   137

Query   145   HCSRGERRAVEKLSVEALNSLTGEFKGKYYPLKSMTEKEQQQLIDIHFLFDKPVSPLLLA   204
              +R ER+ +E L    AL    TGE KGKYY L++M++ E++QLI IHFLF K       L +
Sbjct   138   AVTRKERKEIEHLVTSALGEFTGELKGKYYSLETMSDAEKKQLIAIHFLF-KGGDKYLQS   196

Query   205   SGMARDWPDARGIVHIDNKSFLVWVNEEDHLRVISMEKGGNMKEVFRRFCVG          256
              +G+ RDWP+ARGI HI D K+FLVWVNEED LR+ISM+ G N+ EVF+R  V
Sbjct   197   AGLERDWPEARGIVHIDAKTFLVWVNEEDQLRIISMQAGSNILEVFKRLSVA          248
```

Scores from BLOSUM62, a position independent matrix

- NCBI Discovery Workshops

# PSSM Alignment: Globins



**cd01040: globin, with user query added** ?

Globins are heme proteins, which bind and transport oxygen. This family summarizes a diverse set of homologous protein domains, including: (1) tetrameric vertebrate hemoglobins, which are the major protein component of erythrocytes and transport oxygen in the bloodstream, (2) microorganismal flavohemoglobins, which are linked to C-terminal FAD-dependend reductase domains, (3) homodimeric bacterial hemoglobins, such as from Vitreoscilla, (4) plant leghemoglobins (symbiotic hemoglobins, involved in nitrogen metabolism in plant rhizomes), (5) plant non-symbiotic hexacoordinate globins and hexacoordinate globins from bacteria and animals, such as neuroglobin, (6) invertebrate hemoglobins, which may occur in tandem-repeat arrangements, and (7) monomeric myoglobins found in animal muscle tissue.

```
Feature 1                                           ##              ##        #  #                       ##
1ASH           1 ANKTRELCMKSL.[12].QDGI ELLDRHAR RKYF.[16].FAKQGQKILLACHVLCA.[13].ELLDRHAR  99
query          2 TPAQIALVQQSF.[ 8].QAAS RLAKLHVS QPLF.[ 4].IRDQGKKLMGTLAVVVG.[13].RLAKLHVS  84
gi 13810249   18 NILQRLKVKNQW.[11].SXGT QLAHLHAQ DKFF.[12].FQAHIQRVFGGFDMCIS.[10].QLAHLHAQ 109
gi 20513982    3 SSHERSLIRKTW.[ 7].DVAF ALGGAHQA QKMF.[16].FLAQAYTILAGLNVVIQ.[13].ALGGAHQA  96
gi 22001638   14 GEEQEALVLKSW.[ 8].NLGI RLGATHLR EQMF.[15].LKTHAMSVFVMTCEAAA.[16].RLGATHLR 110
gi 22960923    8 SPADIHRVRTSF.[ 8].EMAD KLAVDHVR RTLF.[ 3].MTRMKDKFIQTLAVLVG.[13].KLAVDHVR  89
gi 25495425   21 NEIKRLKVKLQW.[11].DFED FLKAQHAP EKFF.[12].FRAFGMRVASGLDMVLS.[13].FLKAQHAP 115
gi 32417616    4 TYQQSKLVRDTI.[ 8].RITS RMCNKHCS NNYF.[ 6].NGRQPRALTAVILGFAS.[13].RMCNKHCS  88
gi 33300043   12 TQEEKNDLEHSW.[ 8].HIAC NLGRRHGK RRLF.[19].QAMRFMQVIEGAVKALD.[10].NLGRRHGK 106
gi 34447132    7 SIEDIRDIQHDW.[13].VFGQ HLSQQHKE KGVH.[ 8].FKNHVLRVLNGLDNLIN.[13].HLSQQHKE 102
```

Conserved Histidine

**- NCBI Discovery Workshops**

# PSSM Viewer



Histidine scored differently at two positions

- NCBI Discovery Workshops

# Build PSSM with PSI-BLAST

- PSI-BLAST

  1. Iteration 1: Regular BLASTP (BLOSSOM62) to identify a list of closely related proteins. Build PSSM from these proteins.

  2. Iteration 2: Use the PSSM built from Iteration 1 to score alignment in this Iteration.

  3. Repeat multiple iterations.

**- NCBI Discovery Workshops**

# Build PSSM with DELTA-BLAST

**DELTA-BLAST employs a subset of NCBI's Conserved Domain Database (CDD) to construct PSSM**

# Heme Binding Site

Conserved Histidine

**blastp**

```
TFATLSELHCDKLHVD----PENFRLLG
     S L    KLHV       P ++   +G
ILPAASRLA--KLHVSYGVQPTHYAPVG
```

**DELTA-BLAST**

```
TF---ATLSELHCDKLHVDPENFRLLG
   + L++LH       V P ++   +G
ILPAASRLAKLHVS-YGVQPTHYAPVG
```

# Heme Binding Site

Conserved Histidine

bla...                                                    ...LGI
                                                          ...G
                                                          ...GA

**BLAST is not reliable for alignment of homologous genes between distantly related species.**

DELTA-BLAST

ILPAASRLAKLHVS-YGVQPTHYAPVG.

# BLAST does Local Alignment
# (Basic Local Alignment Search Tool)

## Local Alignment    vs    Global Alignment



HSP-1

HSP-2

(Bowtie, BWA, ClustalW, et al)



Distribution of 41 Blast Hits on the Query Sequence

Mouse-over to show defline and scores.  Click to show alignments

Color Key for Alignment Scores

| <40 | 40-50 | 50-80 | 80-200 | >=200 |

# BLAST and BLAST–like programs

- **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 query  vs. translated database
- **Megablast  nucleotide only**
  - **Contiguous megablast**
    - **Nearly identical sequences**
  - **Discontiguous megablast**
    - **Cross-species comparison**
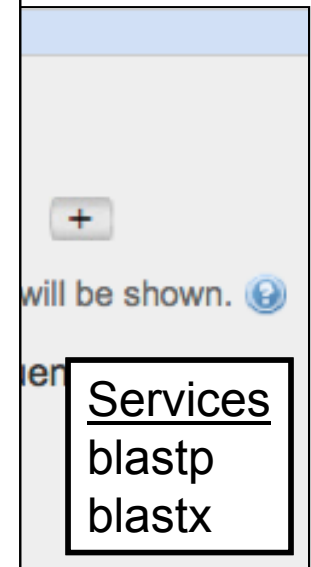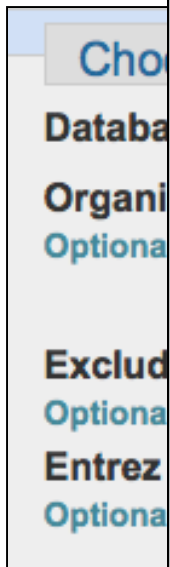
# Nucleotide Databases: List

# Non−redundant protein

**nr** (non-redundant protein sequences)

– GenBank CDS translations

– NP_, XP_ **refseq_protein**

– Outside Protein

  • PIR, **Swiss-Prot**, PRF

  • **PDB** (sequences from structures)

**pat** protein patents

**env_nr** metagenomes
  (environmental samples)

Services
blastp
blastx

# Reference Sequence Databases

**Archive Databases**

Genbank & Genpept
(NCBI nt and nr)

**Genome Based Reference**

NCBI Refseq

UCSC Genomes (Animals)

Ensembl (Animals)

Phytozome (Plants)

Ensembl Plants (Plants)
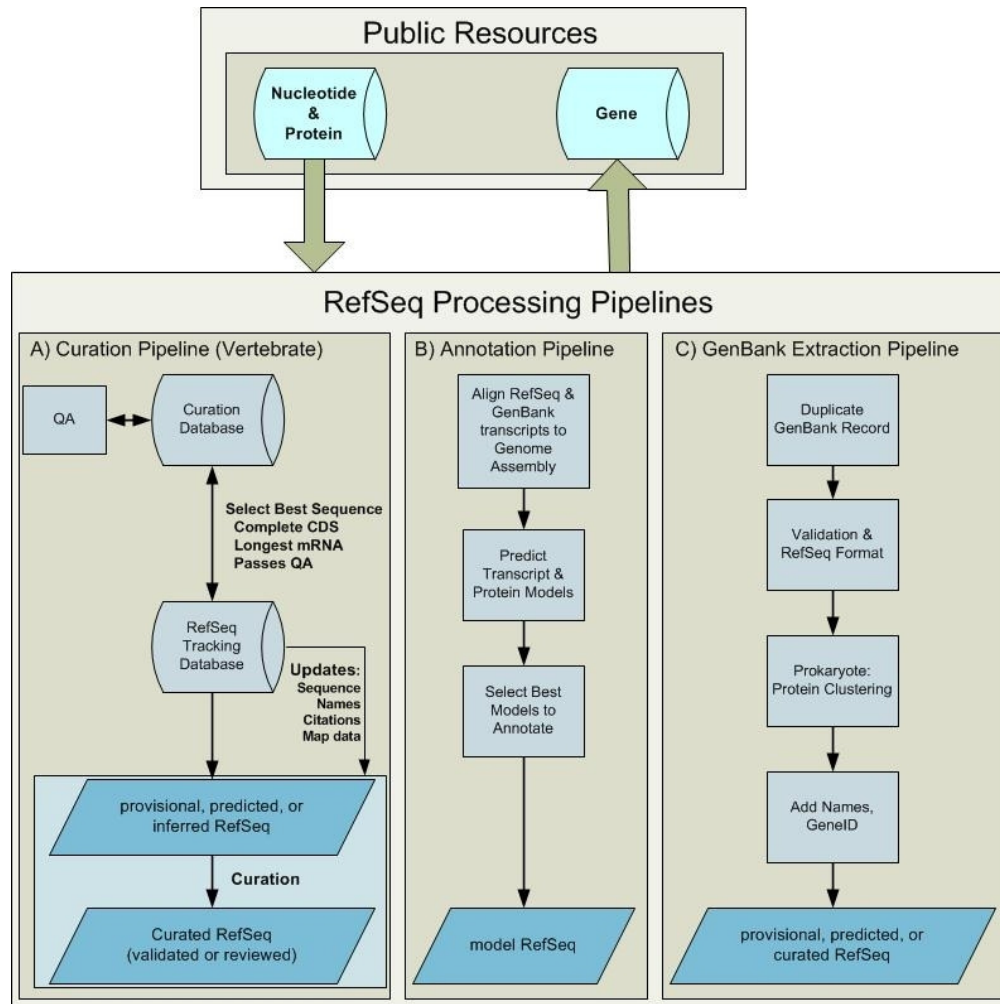
**Species Specific Databases**

Flybase
Wormbase
TAIR
et al.

# RefSeq Curation Pipeline



**Accessions:**

NM_000000     mRNA

NP_000000     Protein

NC_000000     genome assembly
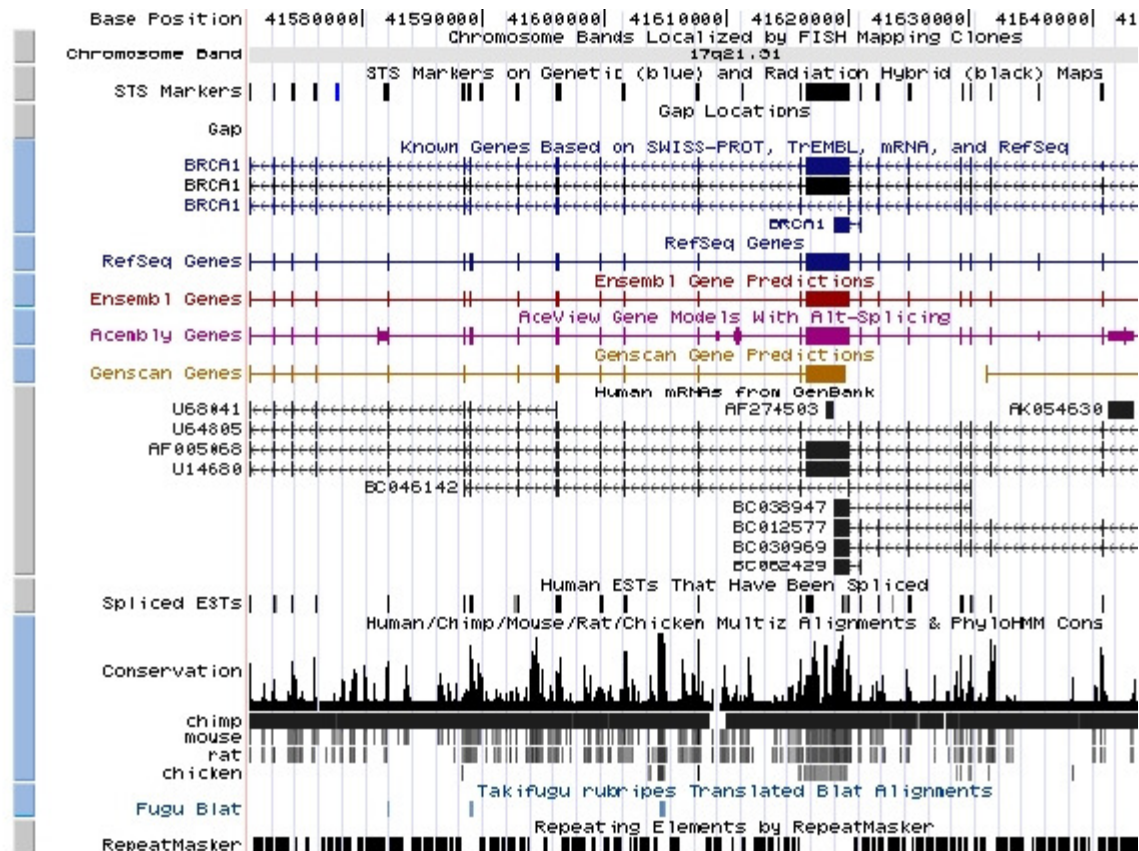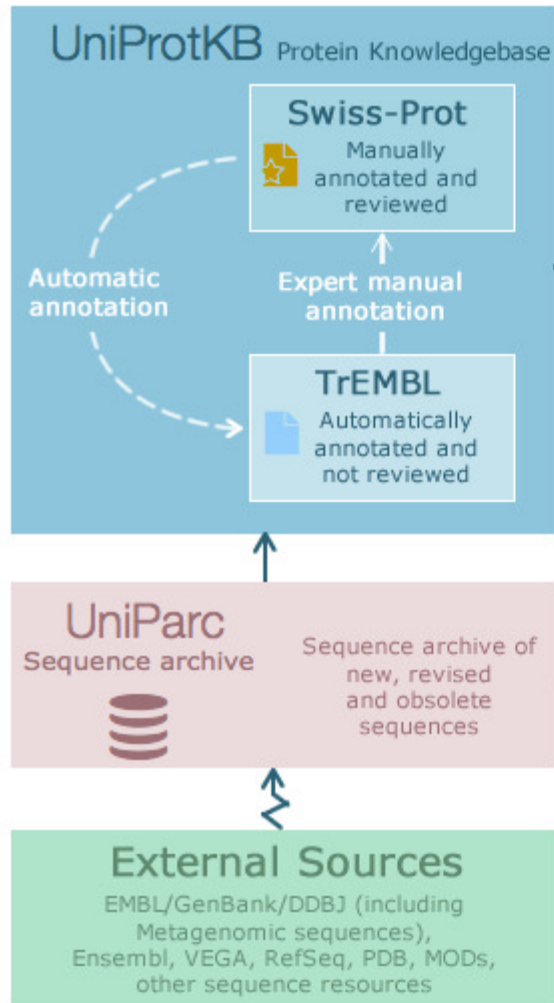
XM_000000     predicted mRNA
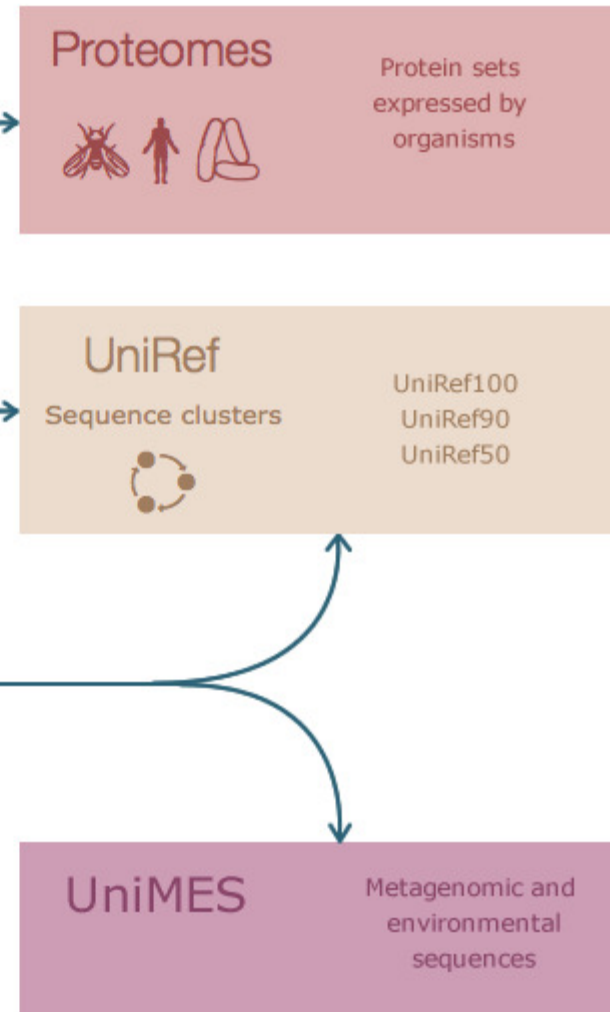
XP_000000     predicted protein

# UCSC Genome Database



RefSeq records are mapped to the genome by UCSC.

# UniProtKB/Swiss-Prot

# UniRef

**Versions of NCBI, UCSC and Ensembl**


**Versions of genome assembly:**
**NCBI:** GRCh38 (December 2013)
**UCSC:** HG38
**Ensembl:** GRCh38


**Versions of gene annotation:**
**NCBI:** release 69 (January, 2015)
**UCSC:** Daily incremental updates as NCBI
**Ensembl:** release 79 (March 2015)

# Database usage examples

1. **Identify species for a list of sequences.**
   Tool on BioHPC: fastq_species_detector  (using NCBI nt database)
   https://cbsu.tc.cornell.edu/lab/userguide.aspx?a=software&i=149#c

2. **Function annotation**
- BLAST to a closely related species. (Ensembl, Flybase, et al)
- BLAST to swiss-prot  (used by BLAST2GO) or Uni-ref (used by Trinotate).