

Finishing up  
databases

Loose ends

RefSeq accession  
numbers

NCBI tutorials

NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST

Orthologs and  
Paralogs

Gene duplication  
and mutation

Next time

Reading for  
next class

## Finish databases

## Start sequence alignment

## Lec'05'slides

## Finishing up databases

- Loose ends
- RefSeq accession  
numbers
- NCBI tutorials
- NCBI Demo

## Introduction to Pairwise Sequence Alignment

- BLAST
- Orthologs and  
Paralogs
- Gene duplication  
and mutation

## Next time

## Reading for next class

### ① Finishing up databases

- Loose ends

- RefSeq accession numbers

- NCBI tutorials

- NCBI Demo

### ② Introduction to Pairwise Sequence Alignment

- BLAST

- Orthologs and Paralogs

- Gene duplication and mutation

### ③ Next time

### ④ Reading for next class

Finishing up  
databases

Loose ends

RefSeq accession  
numbers

NCBI tutorials

NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST

Orthologs and  
Paralogs

Gene duplication  
and mutation

Next time

Reading for  
next class

## 1 Finishing up databases

Loose ends

RefSeq accession numbers

NCBI tutorials

NCBI Demo

## 2 Introduction to Pairwise Sequence Alignment

BLAST

Orthologs and Paralogs

Gene duplication and mutation

## 3 Next time

## 4 Reading for next class

# FASTA format is very common

Finishing up  
databases

Loose ends

RefSeq accession  
numbers

NCBI tutorials

NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST

Orthologs and  
Paralogs

Gene duplication  
and mutation

Next time

Reading for  
next class

- Line 1: ">" (greater than) followed by descriptive text
- Lines 2..n: sequence data (DNA or protein)

Examples:

>L19872\_AhR\_human mRNA

```
ATGAACAGCAGCAGCGCCAACATCACCTACGCCAGTCGCAAGCGGCGGAAGCCG
TGCAGAAAACAGTAAAGCCAATCCCAGCTGAAGGAATCAAGTCAAATCCTTCCA
GCGGCATAGAGACCGACTTAATACAGAGTTGGACCGTTTGGCTAGCCTGCTGCC
TTCCACAAGATGTTATTAATAAGTTGGACAAACT
```

>L19872\_AhR\_human\_1 amino acids

```
MNSSSANITYASRRRRKPVQKTVKPIPAEGIKSNPSKRHRDRLNTELDRLASLL
INKLDKLSVLRLSVSYLRAKSFFDVALKSSPTERNGGQDNCRAANFREGLNLQE
ALNGFVLVTTDALVFYASSTIQDYLGFFQQSDVIHQSVYELIHTEDEAEFQRQL
```

## Sample GenBank record:

```
LOCUS      YP_002302326 238 aa linear BCT 04-NOV-2008
DEFINITION green fluorescent protein
ACCESSION  YP_002302326
VERSION    YP_002302326.1  GI:211909965
AUTHORS    Srikhanta, Dowideit, ...
TITLE      Phasevarion mediated random switching ...
```

**Locus** Traditional identifier, not unique across databases, mainly historical

**Accession** Unique identifier, **stable forever**

**Version** An identifier that changes with each update to the sequence or annotation

**GI** Geninfo identifier

For full description:

<http://www.ncbi.nlm.nih.gov/Sitemap/samplerecord.html>

# Primary vs. Derivative Sequence Databases



Finishing up  
databases

Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

## 1 Finishing up databases

Loose ends

RefSeq accession numbers

NCBI tutorials

NCBI Demo

## 2 Introduction to Pairwise Sequence Alignment

BLAST

Orthologs and Paralogs

Gene duplication and mutation

## 3 Next time

## 4 Reading for next class

# NCBI RefSeq accession numbers

Finishing up  
databases

Loose ends

RefSeq accession  
numbers

NCBI tutorials

NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST

Orthologs and  
Paralogs

Gene duplication  
and mutation

Next time

Reading for  
next class

- Accession number format indicates RefSeq
- Format: *letter letter underscore number...*
  - Example: NM\_123456
- Two letters identify entry type. Important examples in **red**:

NG_	gene
NM_	mRNA sequence
XM_	computer predicted mRNA sequence
NP_	protein sequence
XP_	computer predicted protein sequence
NC_	full length chromosome (big!)
(others)	See RefSeq documentation

- Non-RefSeq GenBank accession numbers have no underscore.  
Example: AB123456



Finishing up  
databases

Loose ends

RefSeq accession  
numbers

**NCBI tutorials**

NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST

Orthologs and  
Paralogs

Gene duplication  
and mutation

Next time

Reading for  
next class

## 1 Finishing up databases

Loose ends

RefSeq accession numbers

**NCBI tutorials**

NCBI Demo

## 2 Introduction to Pairwise Sequence Alignment

BLAST

Orthologs and Paralogs

Gene duplication and mutation

## 3 Next time

## 4 Reading for next class

# There are lots of good tutorials at NCBI

Finishing up  
databases

Loose ends

RefSeq accession  
numbers

**NCBI tutorials**

NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST

Orthologs and  
Paralogs

Gene duplication  
and mutation

Next time

Reading for  
next class

## NCBI Home

### Resource List (A-Z)

All Resources

Chemicals & Bioassays

Data & Software

DNA & RNA

Domains & Structures

Genes & Expression

Genetics & Medicine

Genomes & Maps

Homology

Literature

Proteins

Sequence Analysis

Taxonomy

Training & Tutorials **1**

Variation

## Welcome to NCBI

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

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

## Get Started

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

**3**

### NCBI YouTube channel

Learn how to get the most out of NCBI tools and databases with video tutorials on the NCBI YouTube Channel.

**GO**

# Typical roadmap for simple NCBI queries

- 1 Search **All Databases** <http://www.ncbi.nlm.nih.gov/>
- 2 Read NCBI **Bookshelf** on the subject
- 3 If it is a human disease, look at **OMIM**  
(**O**nline **M**endelian **I**nheritance in **M**an)
- 4 Look at the papers in **PubMed**, start with **review papers**.
- 5 Find the relevant genes
- 6 Search All Databases for the relevant genes, refine the query by species, etc.
- 7 Look in **Gene** for information on gene  
<http://www.ncbi.nlm.nih.gov/gene>
- 8 Find **RefSeq**, **protein structure**, **homologues**, etc
- 9 More details throughout the semester

Finishing up  
databases

Loose ends  
RefSeq accession  
numbers

NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

Finishing up  
databases

Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
**NCBI Demo**

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

## 1 Finishing up databases

Loose ends

RefSeq accession numbers

NCBI tutorials

**NCBI Demo**

## 2 Introduction to Pairwise Sequence Alignment

BLAST

Orthologs and Paralogs

Gene duplication and mutation

## 3 Next time

## 4 Reading for next class

Finishing up  
databases  
Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment  
BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

Remember glycolysis? Let's check out **hexokinase**.

Let's checkout how **anemia** relates to **hexokinase**.

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

Finishing up  
databases

Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

## ① Finishing up databases

Loose ends

RefSeq accession numbers

NCBI tutorials

NCBI Demo

## ② Introduction to Pairwise Sequence Alignment

BLAST

Orthologs and Paralogs

Gene duplication and mutation

## ③ Next time

## ④ Reading for next class

# Introduction to Ch. 3 - Pairwise Sequence Alignment

Finishing up  
databases

Loose ends

RefSeq accession  
numbers

NCBI tutorials

NCBI Demo

**Introduction  
to Pairwise  
Sequence  
Alignment**

BLAST

Orthologs and  
Paralogs

Gene duplication  
and mutation

Next time

Reading for  
next class

Biological justification for much of bioinformatics

# Think back to Biol-360 day one...

Finishing up  
databases

Loose ends

RefSeq accession  
numbers

NCBI tutorials

NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST

Orthologs and  
Paralogs

Gene duplication  
and mutation

Next time

Reading for  
next class

## Three driving forces behind bioinformatics:

- Massive volumes of DNA sequence data
- Gene conservation between species
- Systems biology



Finishing up  
databases

Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

**BLAST**  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

## ① Finishing up databases

Loose ends  
RefSeq accession numbers  
NCBI tutorials  
NCBI Demo

## ② Introduction to Pairwise Sequence Alignment

**BLAST**  
Orthologs and Paralogs  
Gene duplication and mutation

## ③ Next time

## ④ Reading for next class

# Many of you have done BLAST at NCBI

Finishing up  
databases  
Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

## BLAST

Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

## ① Query NCBI for human insulin. Here's the RefSeq output:

```
LOCUS      NM_000207 450 bp mRNA linear PRI 20-DEC-2003
DEFINITION Homo sapiens insulin (INS), mRNA.
ACCESSION  NM_000207
VERSION    NM_000207.1  GI:4557670

ORIGIN
      1 gctgcatcag aagaggccat caagcacatc actgtccttc tgccatggcc ctgtggatgc
     61 gcctcctgcc cctgctggcg ctgctggccc tctggggacc tgaccagcc gcagccttgc
    121 tgaaccaaca cctgtgcggc tcacacctgg tggaagctct ctacctagtg tgcggggaac ...
```

## ② BLAST this mRNA sequence. Sample BLAST output:

83% identity

```
Human  catggccctgtggatgcgccctcctgccctgctggcgctgctggccc
      |||
Mouse  catggccctgtggatgcgcttcctgccctgctggccctgctcttcc
```

Two HUGE biological questions:

1. What's going on?
2. Why do we care?

Finishing up  
databases

Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

## ① Finishing up databases

Loose ends

RefSeq accession numbers

NCBI tutorials

NCBI Demo

## ② Introduction to Pairwise Sequence Alignment

BLAST

Orthologs and Paralogs

Gene duplication and mutation

## ③ Next time

## ④ Reading for next class

# Orthologs, Paralogs, and Homologues

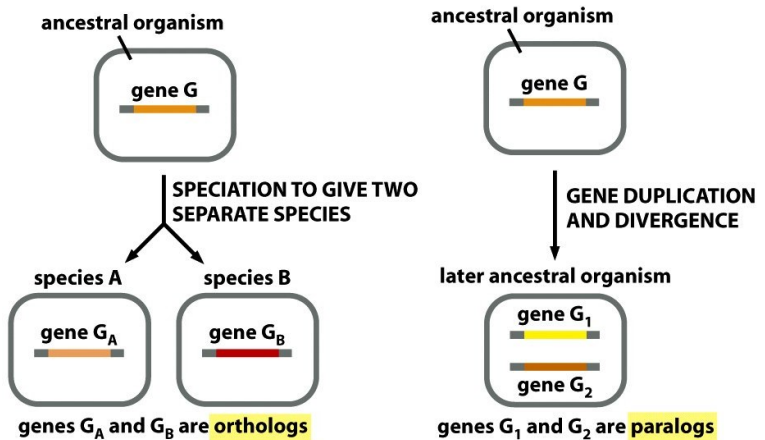


Figure 1-25 Molecular Biology of the Cell 5/e (© Garland Science 2008)

Note to self: This would be a great iClicker question.

Finishing up  
databases  
Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment  
BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

# Evolution of the globin family

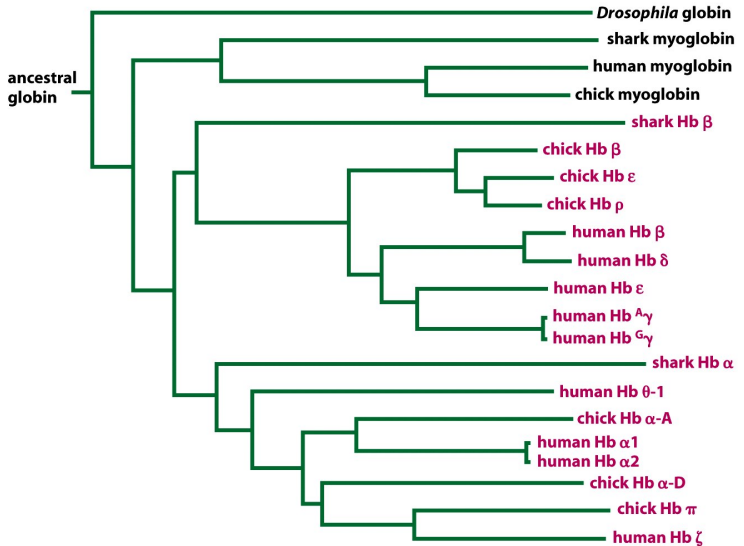


Figure 1-26 Molecular Biology of the Cell 5/e (© Garland Science 2008)

Finishing up  
databases  
Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment  
BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

# What happens to paralogs?

## Key Point:

Gene duplication produces new genes.

Fate of new genes:

- Copy is initially identical, so (simple model):
  - 2× amount of mRNA
  - 2× amount of protein product.
- Is this duplicated gene beneficial for the organism?
  - If beneficial, then natural selective pressure preserves both copies. Examples: tRNA and rRNA for protein synthesis exist in many copies.
  - If detrimental, then selective pressure tends to **deactivate** one copy. The inactive copy is then called a **pseudogene**, and is no longer under any selective pressure.
  - If neutral, then one copy is free to evolve, as long as one continues to function.

Finishing up  
databases  
Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment  
BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

Finishing up  
databases

Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

## ① Finishing up databases

Loose ends

RefSeq accession numbers

NCBI tutorials

NCBI Demo

## ② Introduction to Pairwise Sequence Alignment

BLAST

Orthologs and Paralogs

Gene duplication and mutation

## ③ Next time

## ④ Reading for next class

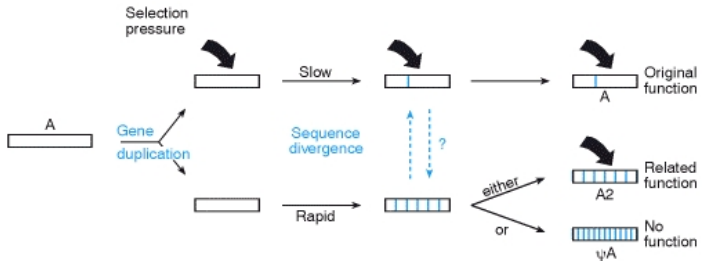
# Gene duplication & evolution

Finishing up  
databases  
Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment  
BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class





# How does all this relate to BLAST?

- BLAST looks for **sequence similarity**
  - Example: AAACCG is similar to TAACCG
- Homology + evolution creates **sequence similarity**

83% identity

Human	catggccctgtggatg	cgccctcctgcccctgctggcgctgctggccc
Mouse	catggccctgtggatg	cgcttcctgcccctgctggccctgctcttcc

Therefore, we look for **sequence similarity** as **evidence of homology**.

Therefore, we use BLAST to find **evidence of homology**.

Finishing up  
databases  
Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment  
BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

Finishing up  
databases

Loose ends

RefSeq accession  
numbers

NCBI tutorials

NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST

Orthologs and  
Paralogs

Gene duplication  
and mutation

Next time

Reading for  
next class

Note to self:

Previous slide will make another great iClicker question.

Prof keeps saying this is a key point of the course.

Finishing up  
databases

Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

## ① Finishing up databases

Loose ends  
RefSeq accession numbers  
NCBI tutorials  
NCBI Demo

## ② Introduction to Pairwise Sequence Alignment

BLAST  
Orthologs and Paralogs  
Gene duplication and mutation

## ③ Next time

## ④ Reading for next class

# Where we're going from here

- Ch. 3 & 4 cover sequence alignment & BLAST
- Alignment is finding the closest arrangement of two sequences:

--ACCTAGGA  
ACCTAGGA-- →

-ACCTAGGA  
ACCTAGGA- →

ACCTAGGA  
ACCTAGGA →

We have a winner  
100% identity

ACCTAGGA-  
-ACCTAGGA →

ACCTAGGA--  
--ACCTAGGA

- Alignment is the essence of BLAST and **many** other bioinformatic tools.

Finishing up  
databases  
Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment  
BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

Finishing up  
databases

Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

## ① Finishing up databases

Loose ends

RefSeq accession numbers

NCBI tutorials

NCBI Demo

## ② Introduction to Pairwise Sequence Alignment

BLAST

Orthologs and Paralogs

Gene duplication and mutation

## ③ Next time

## ④ Reading for next class

Finishing up  
databases

Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

## ① Finishing up databases

Loose ends  
RefSeq accession numbers  
NCBI tutorials  
NCBI Demo

## ② Introduction to Pairwise Sequence Alignment

BLAST  
Orthologs and Paralogs  
Gene duplication and mutation

## ③ Next time

## ④ Reading for next class

Finishing up  
databases  
Loose ends  
RefSeq accession  
numbers  
NCBI tutorials  
NCBI Demo

Introduction  
to Pairwise  
Sequence  
Alignment

BLAST  
Orthologs and  
Paralogs  
Gene duplication  
and mutation

Next time

Reading for  
next class

## Chapter 3

### Section “Scoring Matrices” to Section “Pairwise Alignment and Limits of Detection”

Pages	Notes
79–94	Read

Expect a few iClicker questions on these foundational ideas:

- What are homologues, orthologs, paralogs?
- What can happen to paralogs over evolutionary time?
- **Important:** What are the subtle distinctions between **sequence similarity** and **biological homology**?

(If you're not **1000% clear** on these, post anonymously to the Blackboard class-wide discussion group.)