# DNA as Biological Information

Rasmus Wernersson

Henrik Nielsen
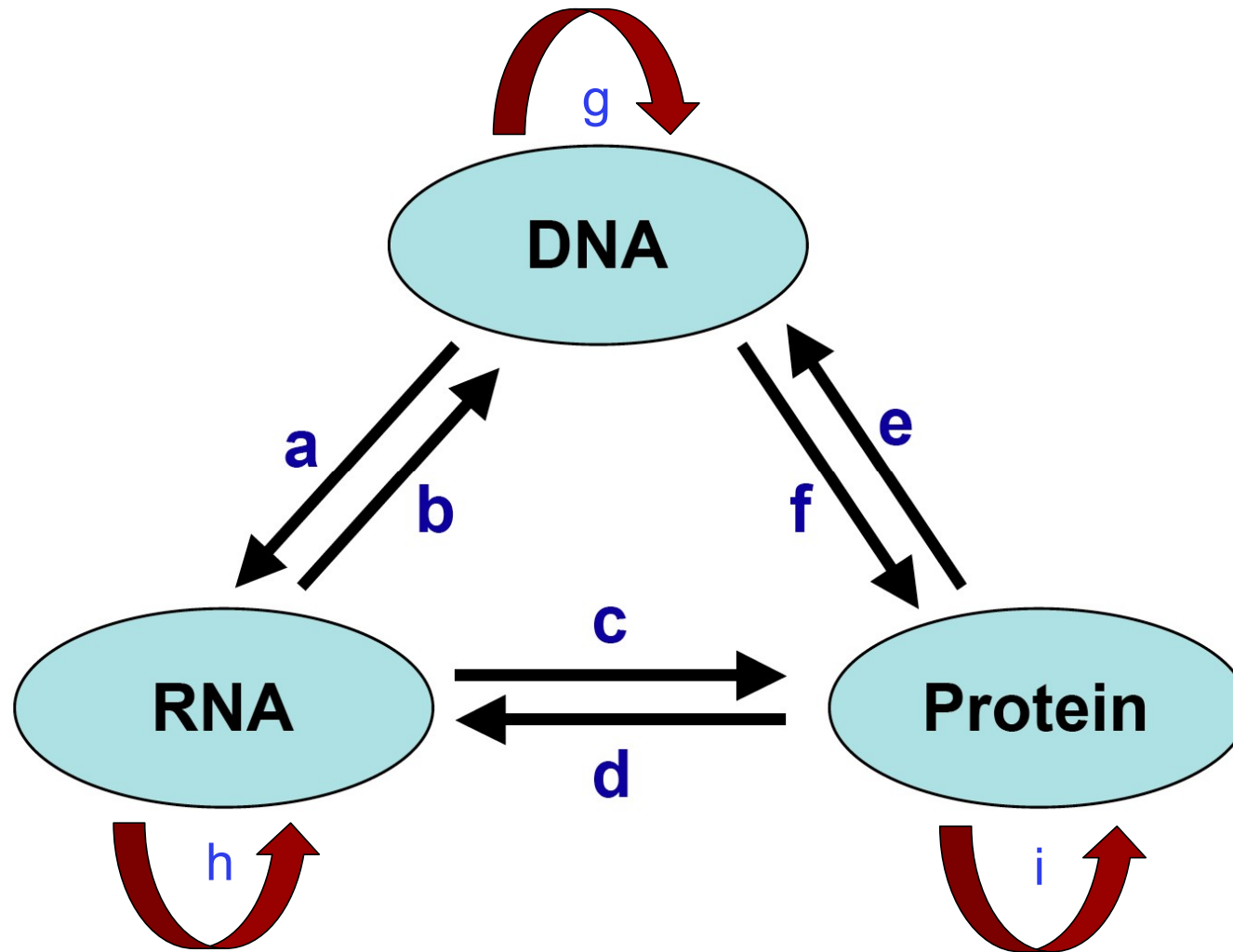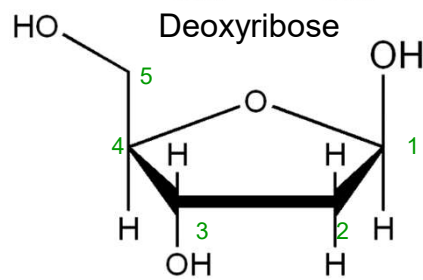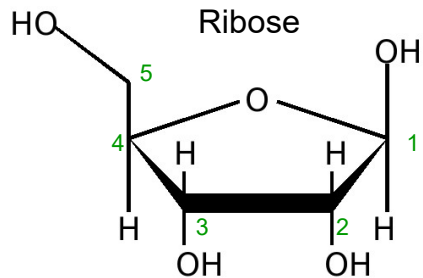
Carolina Barra

# Learning objectives

- Interpret DNA as Biological information

- Describe DNA sequencing techniques and DNA data

- Identify file formats used to store DNA data

- Recognize how information is stored in GenBank database

# DNA as *information*
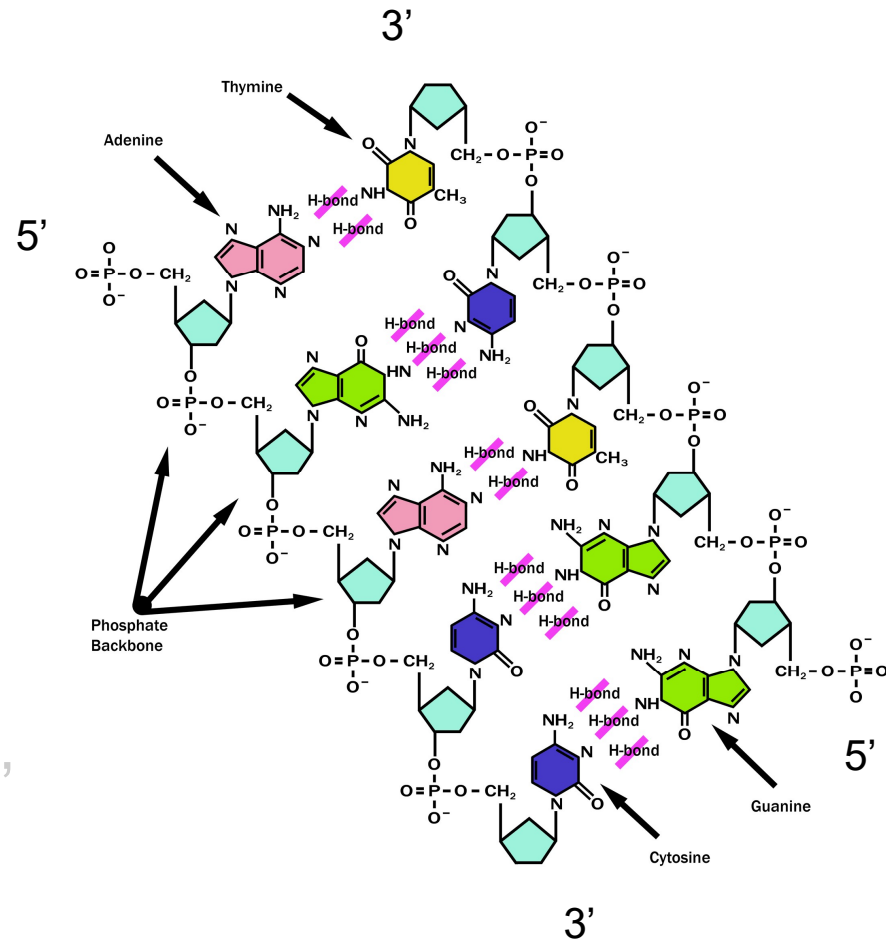
# Information flow in biological systems

# DNA sequences = summary of information



Ribose

Deoxyribose

5' AGCC 3'

3' TCGG 5'

5' ATGGCCAGGTAA 3'



DNA backbone: http://en.wikipedia.org/wiki/DNA
(Deoxy)ribose: http://en.wikipedia.org/

# From molecules to computer files

# Reminder: PCR
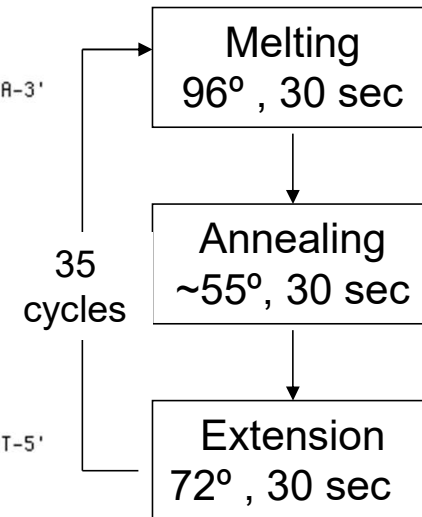
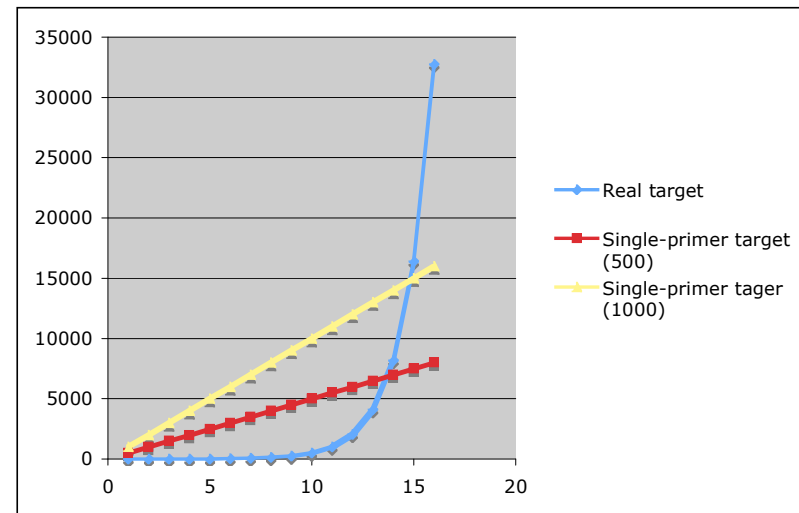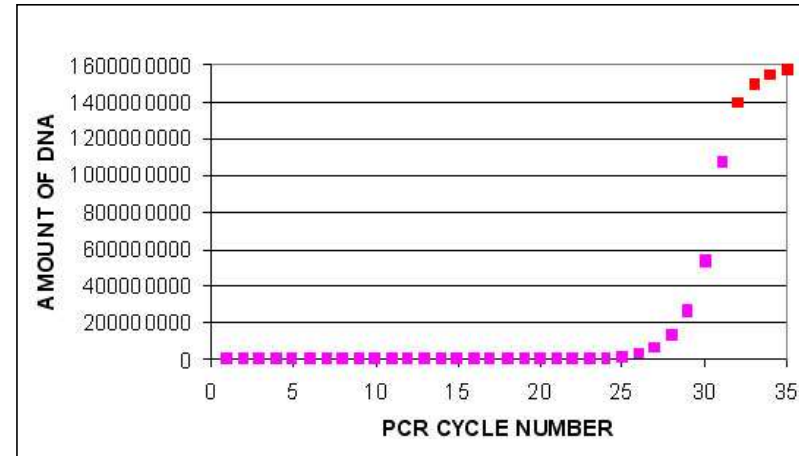Cycle 1

5'-CTAGAATATGAAACCTATAGGTACGGTGGCCATTCTATGTCTGATCCCGGTACTACCTACAGAA-3'
                                                  |||||||||||||
                                         3'-GGGCCATGATGG-5'

5'-ATGAAACCTATAG-3'
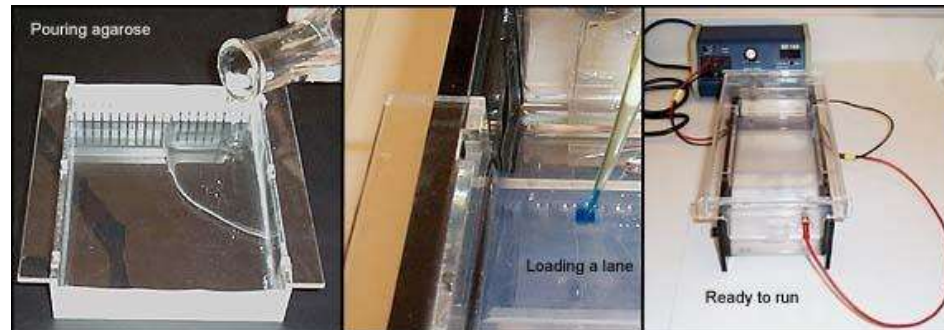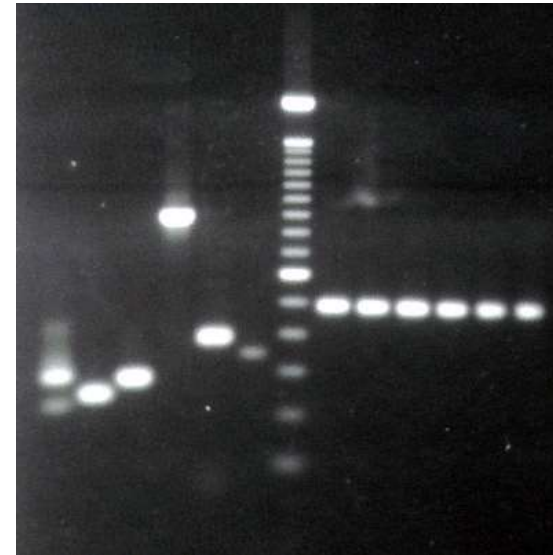|||||||||||||
3'-GATCTTATACTTTGGATATCCATGCCACCGGTAAGATACAGACTAGGGCCATGATGGATGTCTT-5'

```
        ┌──────────────→ ┌─────────────────┐
        │                │   Melting        │
        │                │  96º , 30 sec    │
        │                └─────────────────┘
        │                         │
  35    │                         ▼
 cycles │                ┌─────────────────┐
        │                │   Annealing      │
        │                │  ~55º, 30 sec    │
        │                └─────────────────┘
        │                         │
        │                         ▼
        └──────────────→ ┌─────────────────┐
                         │   Extension      │
                         │  72º , 30 sec    │
                         └─────────────────┘
```

Animation: http://depts.washington.edu/~genetics/courses/genet371b-aut99/PCR_contents.html

# Reminder: PCR





Animation: http://www.people.virginia.edu/~rjh9u/pcranim.html
PCR graph: http://pathmicro.med.sc.edu/pcr/realtime-home.htm
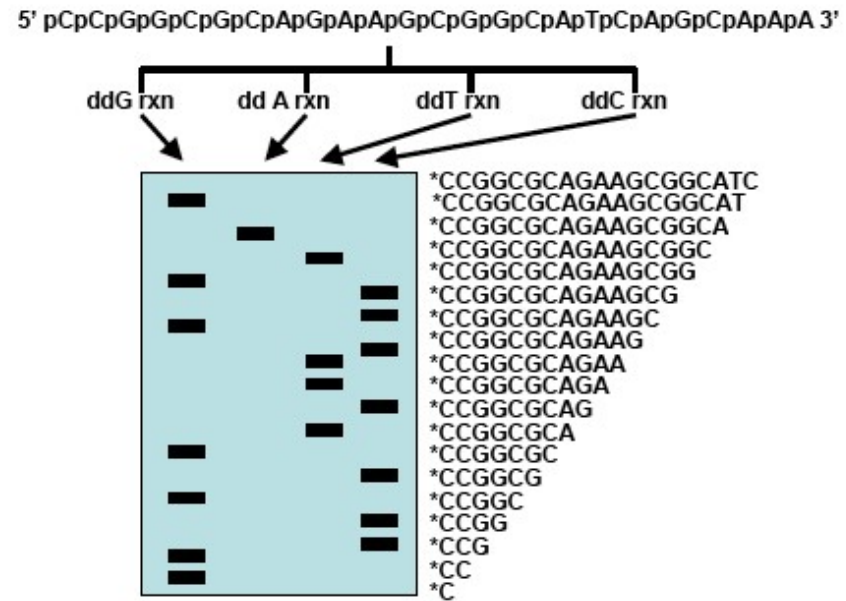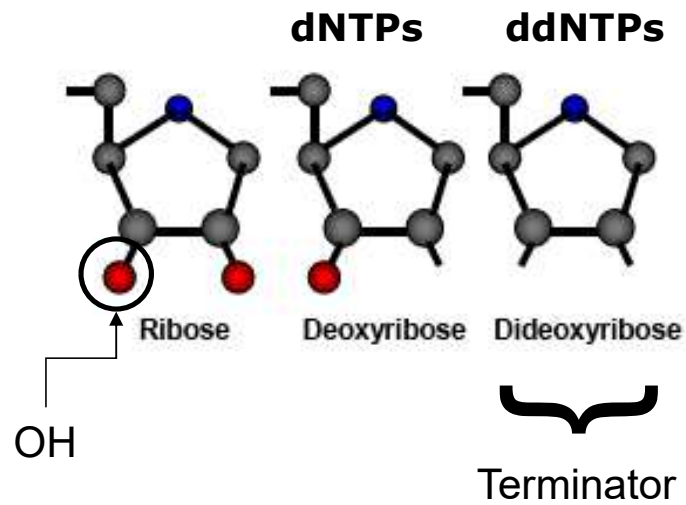
# Gel electrophoresis

- DNA fragments are separated using gel electrophoresis
  - Typically 1% agarose
  - Colored with EtBr or ZybrGreen (glows in UV light).
  - A DNA "ladder" is used for identification of known DNA lengths.

−

+





Pouring agarose

Loading a lane

Ready to run

Gel picture: http://www.pharmaceutical-technology.com/projects/roche/images/roche3.jpg
PCR setup: http://arbl.cvmbs.colostate.edu/hbooks/genetics/biotech/gels/agardna.html

# The Sanger method of DNA sequencing

**dNTPs**  **ddNTPs**

Ribose  Deoxyribose  Dideoxyribose

OH

Terminator

5' pCpCpGpGpCpGpCpApGpApApGpCpGpGpGpCpApTpCpApGpCpApApA 3'

ddG rxn  dd A rxn  ddT rxn  ddC rxn

*CCGGCGCAGAAGCGGCATC
*CCGGCGCAGAAGCGGCAT
*CCGGCGCAGAAGCGGCA
*CCGGCGCAGAAGCGGC
*CCGGCGCAGAAGCGG
*CCGGCGCAGAAGCG
*CCGGCGCAGAAGC
*CCGGCGCAGAAG
*CCGGCGCAGAA
*CCGGCGCAGA
*CCGGCGCAG
*CCGGCGCA
*CCGGCGC
*CCGGCG
*CCGGC
*CCGG
*CCG
*CC
*C

X-ray sequencing gel

video: https://www.youtube.com/watch?v=FvHRio1yyhQ

Images: http://www.idtdna.com/support/technical/TechnicalBulletinPDF/DNA_Sequencing.pdf

# Automated sequencing

- The major break-through of sequencing has happened through *automation.*

- Fluorescent dyes.

- Laser based scanning.

- Capillary electrophoresis

- Computer based base-calling and assembly.



Images: http://www.idtdna.com/support/technical/TechnicalBulletinPDF/DNA_Sequencing.pdf

# Handout exercise: "base-calling"

- Handout: Chromatogram

- Groups of 2-3.

- Tasks:
  - Identify "difficult" regions
  - Identify likely errors
  - Try to estimate the best interval to use

# Automatic assignment of quality



Figure 1. An example of a DNA sequence tracing and the Phred score (grey bars) corresponding to each colored peak. The colored peaks on the trace correspond to each DNA letter. For example 'T' bases are represented in red, and this sequence has four 'T' bases on a row, as viewed by the four red peaks in the sequence. The aqua horizontal line placed across the grey bars represents a Phred score of 20 which is considered an acceptable level of accuracy.

# DNA sequencing - history

**1972** Recombinant DNA technology [Paul Berg].

**1976** The first sequenced genome, the bacteriophage MS2 [Walter Fiers *et al*.]

**1977** DNA sequencing by chemical cleavage [Allan Maxam & Walter Gilbert]
DNA sequencing by enzymatic synthesis [Fred Sanger].

**1982** *GenBank* (public database of DNA sequences).

**1987** The first automatic sequencer, *Prism 373* [Applied Biosystems].

**1990** *Human Genome Project* is launched.

**1995** The first genome of a free-living organism, the bacterium *Haemophilus influenzae* (1.8 Mb) [The Institute for Genomic Research (TIGR)].

**1996** The first genome of a eukaryote, Baker's Yeast, *Saccharomyces cerevisiae* (12.1 Mb) [International consortium].

**1998** The first genome of an animal, the round worm *Caenorhabditis elegans* (97Mb) [Sanger Center and collaborators].

**2001** The first "drafts" of the human genome (3Gb) [Human Genome Project Consortium (Nature, 15 Feb) + Celera (Science, 16 Feb)].

**2000**→ Development of several generations of "Next Generation Sequencing" (NGS)

**October 2022** *GenBank release 252* contains 20.35 trillion bases and 3.10 billion records (including *Whole Genome Shotgun* (WGS) sequences).

Frederick Sanger

Two Nobel Prizes – one for protein sequencing (1958) and one for DNA sequencing (1980)

# "Shotgun sequencing"

# Sequence read mapping

# NGS read mapping

# Cost of sequencing



Cost per Raw Megabase of DNA Sequence

# Cost of sequencing the human genome



Cost per Human Genome

Moore's Law

genome.gov/sequencingcosts

# DNA data bases and data formats

# Background - Nucleotide databases

- **GenBank**, http://www.ncbi.nlm.nih.gov/Genbank/
- National Center for Biotechnology Information (NCBI), National Library of Medicine (NLM), National Institutes of Health (NIH), USA
- Established in 1982.
- **EMBL**, http://www.ebi.ac.uk/embl/
- European Bioinformatics Institute (EBI), England
- Established in 1980 by the European Molecular Biology Laboratory, Heidelberg, Germany
- Now part of **ENA**, the European Nucleotide Archive, http://www.ebi.ac.uk/ena/
- **DDBJ**, http://www.ddbj.nig.ac.jp/
- National Institute of Genetics, Japan

- *Together they form*
- International Nucleotide Sequence Database Collaboration, http://www.insdc.org/

# Nucleotide database growth



NCBI: Growth in public available sequence databases

# FASTA format

**>alpha-D**
ATGCTGACCGACTCTGACAAGAAGCTGGTCCTGCAGGTGTGGGAGAAGGTGATCCGCCAC
CCAGACTGTGGAGCCGAGGCCCTGGAGAGGTGCGGGCTGAGCTTGGGGAAACCATGGGCA
AGGGGGGCGACTGGGTGGGAGCCCTACAGGGCTGCTGGGGGTTGTTCGGCTGGGGGTCAG
CACTGACCATCCCGCTCCCGCAGCTGTTCACCACCTACCCCCAGACCAAGACCTACTTCC
CCCACTTCGACTTGCACCATGGCTCCGACCAGGTCCGCAACCACGGCAAGAAGGTGTTGG
CCGCCTTGGGCAACGCTGTCAAGAGCCTGGGCAACCTCAGCCAAGCCCTGTCTGACCTCA
GCGACCTGCATGCCTACAACCTGCGTGTCGACCCTGTCAACTTCAAGGCAGGCGGGGGAC
GGGGGTCAGGGGCCGGGGAGTTGGGGGCCAGGGACCTGGTTGGGGATCCGGGGCCATGCC
GGCGGTACTGAGCCCTGTTTTGCCTTGCAGCTGCTGGCGCAGTGCTTCCACGTGGTGCTG
GCCACACACCTGGGCAACGACTACACCCCGGAGGCACATGCTGCCTTCGACAAGTTCCTG
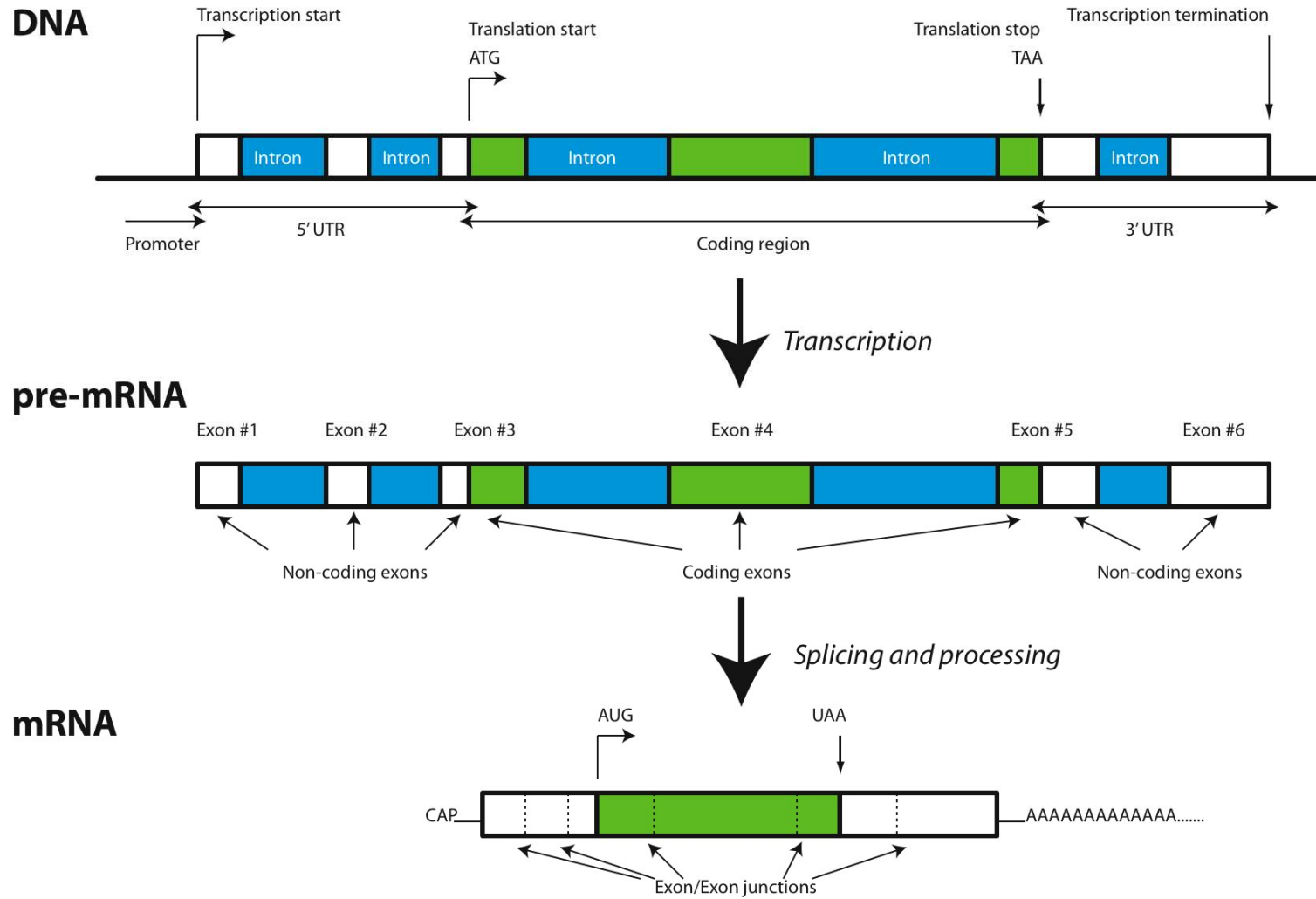TCGGCTGTGTGCACCGTGCTGGCCGAGAAGTACAGATAA
**>alpha-A**
ATGGTGCTGTCTGCCAACGACAAGAGCAACGTGAAGGCCGTCTTCGGCAAAATCGGCGGC
CAGGCCGGTGACTTGGGTGGTGAAGCCCTGGAGAGGTATGTGGTCATCCGTCATTACCCC
ATCTCTTGTCTGTCTGTGACTCCATCCCATCTGCCCCCATACTCTCCCCATCCATAACTG
TCCCTGTTCTATGTGGCCCTGGCTCTGTCTCATCTGTCCCCAACTGTCCCTGATTGCCTC
TGTCCCCCAGGTTGTTCATCACCTACCCCCAGACCAAGACCTACTTCCCCCACTTCGACC
TGTCACATGGCTCCGCTCAGATCAAGGGGCACGGCAAGAAGGTGGCGGAGGCACTGGTTG
AGGCTGCCAACCACATCGATGACATCGCTGGTGCCCTCTCCAAGCTGAGCGACCTCCACG
CCCAAAAGCTCCGTGTGGACCCCGTCAACTTCAAAGTGAGCATCTGGGAAGGGGTGACCA
GTCTGGCTCCCCTCCTGCACACACCTCTGGCTACCCCCTCACCTCACCCCCTTGCTCACC
ATCTCCTTTTGCCTTTCAGCTGCTGGGTCACTGCTTCCTGGTGGTCGTGGCCGTCCACTT
CCCCTCTCTCCTGACCCCGGAGGTCCATGCTTCCCTGGACAAGTTCGTGTGTGCCGTGGG
CACCGTCCTTACTGCCAAGTACCGTTAA

# So we got the DNA sequence – now what?

```
CGATCAGGTTACATTTACTGCCCATGCCTGTCTCAGAGGAATTCTGACACGAAAAGGTGGGCACAAATTC
TTAAGCACACTCTGATGGTACAACGTGAGCTGGCACTACAAGCTGTGTTCCTCATCCCGTTTACAAAATT
TTGAGACTGTGTTTGGGCAAGGGGGAGAGAGACAGTGCAGAAGCTCTGAAGCCACTGAATTTCTCTAAAT
GTGTTTAGAGAAGCTTTTCAAACATGCTACATTTGTGGGTCTCAAGTAACCCGAACATTTAAATCCACAG
CTGAGGTGGTAGGTCAGCAACTGGTGTATCCCTAACTCCAAGTTTGTACCAAAGACTTTAACAGAGGCTG
AGTGAAGGATGTGACAGTCACCAGCCACCATATGCCCTGCCAAATGTCCCAGTGTTTACAGAGGGATAGC
AACACACACTGGGGTGGGAAGAAGGAAAGAAGACCAGGCCTGACAAGCATCACAGGATGGATTTTGGGAA
GACTATGGACCTGAAAAGGAGATTCTTCCCCACTCAGGTCTCTCCAGGATGCTGGGGAGATGCTGTTTCC
TGTGGTAGATCCCCAGCATGAACCAGGAGGGCATGTCCGTGGCTCCTGCTCTGAGGCTCACAGTGTCTTT
GGGTGGAGAGGGGATGGATGCACTGGGGTCTGGAGGACATGAGGGACTGGGGGCGCTCGTGGGATCTACT
CTGACACCTGCAGAGACAGGGAGGACCCTGGCCTGGCCAGAAGGGAATGGTGGATCCCAACAGGAAGCTT
GAGGATATGCAGGTTTGTGAGGCCGAGGCTGTGGCACCCGTGGGACATGCCGATGGCTGCTGTTGACCAT
GGGGCAGCTCAGCCAAGTGCTGCCCCCAGCCCCCAGCCCAGCGTGGGGCTGGTGCAGTGCGGCACATCAG
GGCAGGGCAGCCGCCCCATTGGGGCCCCCTCGGGGCTGGGCCTCCCAGGGCAGTCGGGGCCCCCTGAGGC
AGTGGCCCCCCACCCCTTGGTGCCGATAAGATAACGCTGGGGCGGAGGTGCCGACCACTATAAGAGGATG
TCCTGGTGGGCCCTGCTACCACTGAGCCCTGACCGCCACCCCCAGCCGCCACCATGCTGACCGACTCTGA
CAAGAAGCTGGTCCTGCAGGTGTGGGAGAAGGTGATCCGCCACCCAGACTGTGGAGCCGAGGCCCTGGAG
AGGTGCGGGCTGAGCTTGGGGAAACCATGGGCAAGGGGGGCGACTGGGTGGGAGCCCTACAGGGCTGCTG
GGGGTTGTTCGGCTGGGGGTCAGCACTGACCATCCCGCTCCCGCAGCTGTTCACCACCTACCCCCAGACC
AAGACCTACTTCCCCCACTTCGACTTGCACCATGGCTCCGACCAGGTCCGCAACCACGGCAAGAAGGTGT
TGGCCGCCTTGGGCAACGCTGTCAAGAGCCTGGGCAACCTCAGCCAAGCCCTGTCTGACCTCAGCGACCT
GCATGCCTACAACCTGCGTGTCGACCCTGTCAACTTCAAGGCAGGCGGGGGACGGGGGTCAGGGGCCGGG
GAGTTGGGGGCCAGGGACCTGGTTGGGGATCCGGGGCCATGCCGGCGGTACTGAGCCCTGTTTTGCCTTG
CAGCTGCTGGCGCAGTGCTTCCACGTGGTGCTGGCCACACACCTGGGCAACGACTACACCCCGGAGGCAC
ATGCTGCCTTCGACAAGTTCCTGTCGGCTGTGTGCACCGTGCTGGCCGAGAAGTACAGATAAGCCATCGC
```

# Reminder: Eukaryotic gene structure

```
 901 ggcacatcag ggcagggcag ccgcccccatt ggggcccccct cggggctggg cctcccaggg
 961 cagtcggggc cccctgaggc agtggcccccc caccccttgg tgccgataag ataacgctgg
1021 ggcggaggtg ccgaccacta taagaggatg tcctggtggg ccctgctacc actgagccct
1081 gaccgccacc cccagccgcc accatgctga ccgactctga caagaagctg gtcctgcagg
1141 tgtgggagaa ggtgatccgc cacccagact gtggagccga ggccctggag aggtgcgggc
1201 tgagcttggg gaaaccatgg gcaaggggggg cgactgggtg ggagccctac agggctgctg
1261 ggggttgttc ggctgggggt cagcactgac catcccgctc ccgcagctgt tcaccaccta
1321 cccccagacc aagacctact tccccccactt cgacttgcac catggctccg accaggtccg
1381 caaccacggc aagaaggtgt tggccgcctt gggcaacgct gtcaagagcc tgggcaacct
1441 cagccaagcc ctgtctgacc tcagcgacct gcatgcctac aacctgcgtg tcgaccctgt
1501 caacttcaag gcaggcgggg gacgggggtc agggggccggg gagttggggg ccagggacct
1561 ggttggggat ccggggccat gccggcggta ctgagccctg ttttgccttg cagctgctgg
1621 cgcagtgctt ccacgtggtg ctggccacac acctgggcaa cgactacacc ccggaggcac
1681 atgctgcctt cgacaagttc ctgtcggctg tgtgcaccgt gctggccgag aagtacagat
1741 aagccatcgc tcgtgccgaa gtgccgtcaa taaagacacc tttgctgcag catcgtgtcc
1801 gtctgtgctg gggccaggga cctgggtggg ctgtgctcct gtgggaggga gggaggccgt
1861 ggaacagggg gcagcactgg ccatgggttg cctgggtgcc ccaccaagag ccttgcccac
1921 ttccacaccc ccctctccag cttgggatgt ggctgatggt ggtagcaggg ccagagcgat
1981 ggagctcagc ctgtcacctc gccatgcctg caccctctgg ggagcgggag ctaaagatga
2041 agacaagcgg ttcttgcatc cccttgaagg actctccagg agggtagaat taatccttcc
2101 tggctccatc tatcatgact gttgtttagg actggggaag ctgttggagc tgagtgcctg
2161 catgagctca aaggcagtca gaaaagttta tggagtgaaa aacatccacc acagactatt
2221 aagcacaaac agtgttcttg ggctcaagcc cctgagccac aaatcccagc aactgccagg
2281 gaagtgcccc catacctgcc ctgcgtgctg gtttggctga acgaaggtaa cttttccatga
2341 ggttagttaa ttttcttatt gatagttagc atgggccttt gtttttgaat gtagtttgag
2401 aacaataata tagcatcttg ggacaaagtt aatgttttct tccagtaact ctccagtgtt
```

# GenBank format



- Originates from the GenBank database.

- Contains both a DNA sequence and annotations of features (*e.g.* location of genes).

# GenBank format - HEADER

```
LOCUS       CMGLOAD                   1185 bp    DNA     linear   VRT 18-APR-2005
DEFINITION  Cairina moschata (duck) gene for alpha-D globin.
ACCESSION   X01831
VERSION     X01831.1  GI:62724
KEYWORDS    alpha-globin; globin.
SOURCE      Cairina moschata (Muscovy duck)
  ORGANISM  Cairina moschata
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Archosauria; Aves; Neognathae; Anseriformes; Anatidae; Cairina.
REFERENCE   1  (bases 1 to 1185)
  AUTHORS   Erbil,C. and Niessing,J.
  TITLE     The primary structure of the duck alpha D-globin gene: an unusual
            5' splice junction sequence
  JOURNAL   EMBO J. 2 (8), 1339-1343 (1983)
   PUBMED   10872328
COMMENT     Data kindly reviewed (13-NOV-1985) by J. Niessing.
```

# GenBank format - ORIGIN section

```
ORIGIN
        1 ctgcgtggcc tcagcccctc cacccctcca cgctgataag ataaggccag ggcgggagcg
       61 cagggtgcta taagagctcg gccccgcggg tgtctccacc acagaaaccc gtcagttgcc
      121 agcctgccac gccgctgccg ccatgctgac cgccgaggac aagaagctca tcgtgcaggt
      181 gtgggagaag gtggctggcc accaggagga attcggaagt gaagctctgc agaggtgtgg
      241 gctgggccca gggggcactc acagggtggg cagcagggag caggagccct gcagcgggtg
      301 tgggctggga cccagagcgc cacggggtgc gggctgagat gggcaaagca gcagggcacc
      361 aaaactgact ggcctcgctc cggcaggatg ttcctcgcct accccagac caagacctac
      421 ttcccccact tcgacctgca tcccggctct gaacaggtcc gtggccatgg caagaaagtg
      481 gcggctgccc tgggcaatgc cgtgaagagc ctggacaacc tcagccaggc cctgtctgag
      541 ctcagcaacc tgcatgccta caacctgcgt gttgaccctg tcaacttcaa ggcaagcggg
      601 gactagggtc cttgggtctg ggggtctgag ggtgtggggt gcagggtctg ggggtccagg
      661 ggtctgagtt tcctggggtc tggcagtcct gggggctgag ggccagggtc ctgtggtctt
      721 gggtaccagg gtcctggggg ccagcagcca gacagcaggg gctgggattg catctgggat
      781 gtgggccaga ggctgggatt gtgtttggaa tgggagctgg gcaggggcta gggccagggt
      841 gggggactca gggcctcagg gggactcggg ggggactga gggagactca gggccatctg
      901 tccggagcag gggtactaag ccctggtttg ccttgcagct gctggcacag tgcttccagg
      961 tggtgctggc cgcacacctg ggcaaagact acagccccga gatgcatgct gcctttgaca
     1021 agttcttgtc cgccgtggct gccgtgctgg ctgaaaagta cagatgagcc actgcctgca
     1081 cccttgcacc ttcaataaag acaccattac cacagctctg tgtctgtgtg tgctgggact
     1141 gggcatcggg ggtcccaggg agggctgggt tgcttccaca catcc
//
```

# GenBank format - FEATURE section

```
FEATURES             Location/Qualifiers
     source          1..1185
                     /organism="Cairina moschata"
                     /mol_type="genomic DNA"
                     /db_xref="taxon:8855"
     CAAT_signal     20..24
     TATA_signal     69..73
     precursor_RNA   101..1114
                     /note="primary transcript"
     exon            101..234
                     /number=1
     CDS             join(143..234,387..591,939..1067)
                     /codon_start=1
                     /product="alpha D-globin"
                     /protein_id="CAA25966.2"
                     /db_xref="GI:4455876"
                     /db_xref="GOA:P02003"
                     /db_xref="InterPro:IPR000971"
                     /db_xref="InterPro:IPR002338"
                     /db_xref="InterPro:IPR002340"
                     /db_xref="InterPro:IPR009050"
                     /db_xref="UniProt/Swiss-Prot:P02003"
                     /translation="MLTAEDKKLIVQVWEKVAGHQEEFGSEALQRMFLAYPQTKTYFP
                     HFDLHPGSEQVRGHGKKVAAALGNAVKSLDNLSQALSELSNLHAYNLRVDPVNFKLLA
                     QCFQVVLAAHLGKDYSPEMHAAFDKFLSAVAAVLAEKYR"
     repeat_region   227..246
                     /note="direct repeat 1"
     intron          235..386
                     /number=1
     repeat_region   289..309
                     /note="direct repeat 1"
     exon            387..591
                     /number=2
     intron          592..939
                     /number=2
     exon            940..1114
                     /number=3
     polyA_signal    1095..1100
     polyA_signal    1114
```
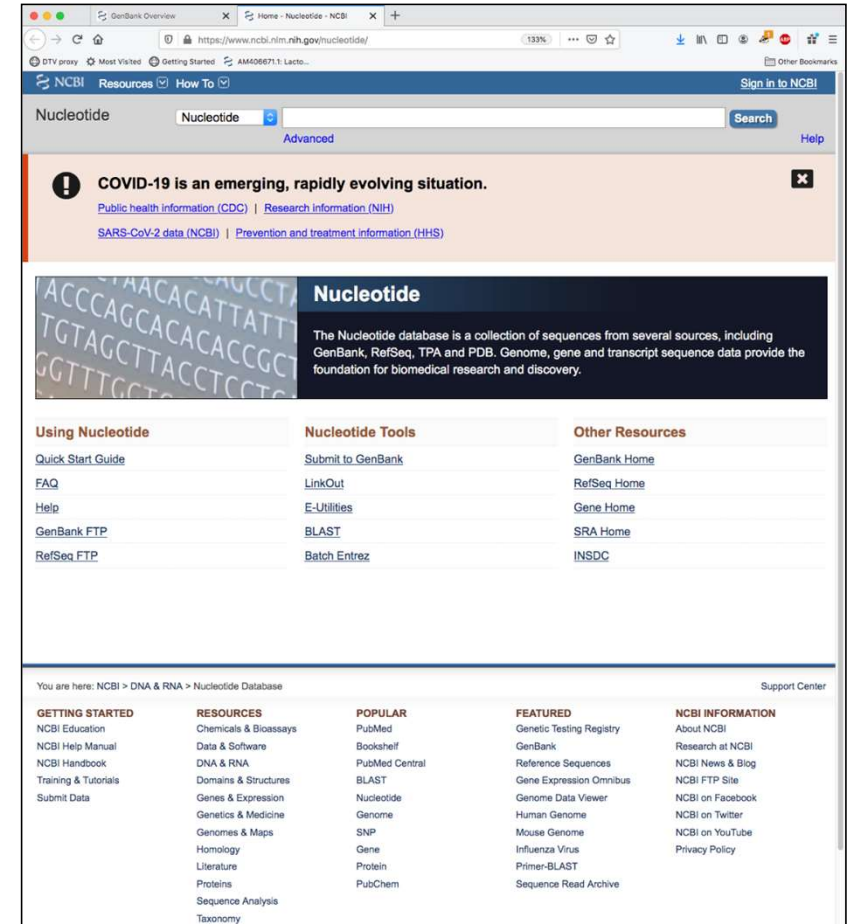
# Exercise: GenBank

- The exercise guide is linked from the course programme.

- Read the guide carefully - it contains a lot of information about GenBank.

- Remember your handouts:
  - GenBank & FASTA format
  - Eukaryotic gene structure