# Sequences, Structures, and Gene Regulatory Networks
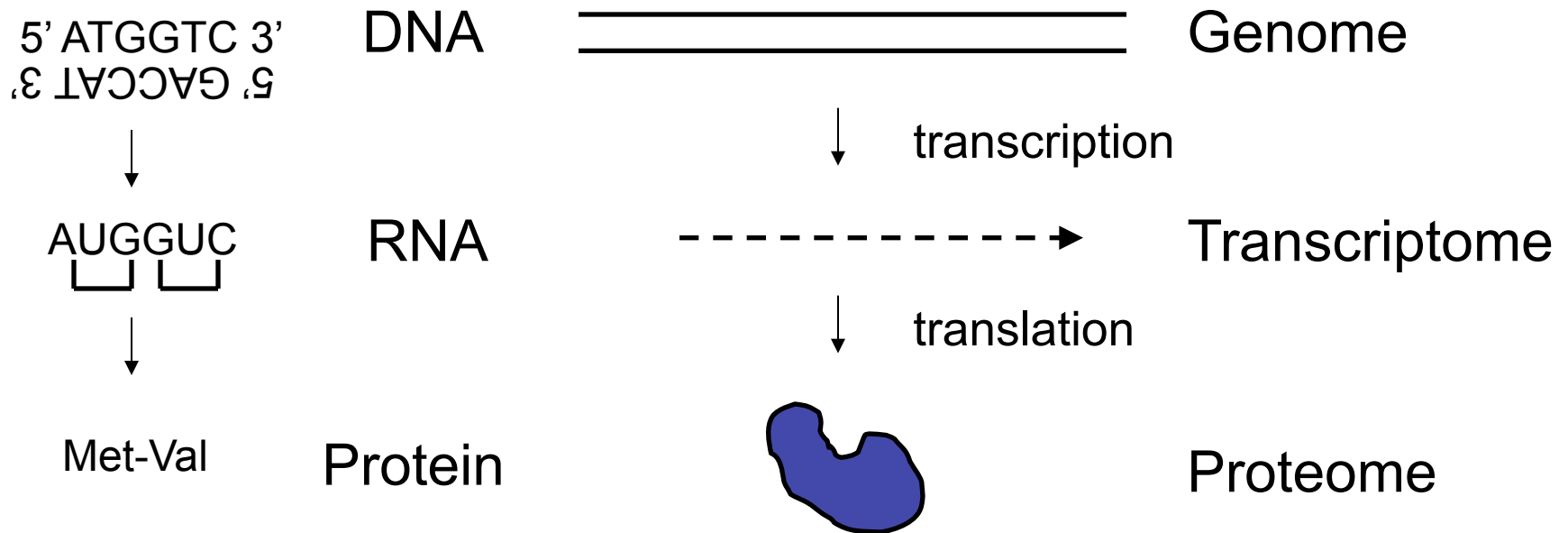
# Learning Outcomes

After this class, you will

- Understand gene expression and protein structure in more detail

- Appreciate why biologists like to align sequences, and have a general idea of how the most commonly used algorithm, BLAST, works

- Be able to use your knowledge of biology to help you critique visual representations of alignments and gene regulatory networks

# Outline

- Sequences and Structures: The Central Dogma in a little more detail

- Alignment – why is this so important? What are important features to visualize?

- Gene regulatory networks

- Appendix: Representation of sequences in databases at the NCBI

# The Central Dogma of Molecular Biology: Genes Encode Proteins

5' ATGGTC 3'
3' TACCAG 5'

DNA

Genome

transcription

AUGGUC

RNA

Transcriptome

translation

Met-Val

Protein

Proteome

# Prokaryotic and eukaryotic transcription and translation



Figure 8-11
*Introduction to Genetic Analysis, Ninth Edition*
© 2008 W. H. Freeman and Company

Griffeths, Introduction to Genetic Analysis, 2008

# Transcription is mediated by RNA Polymerase



© 2012 Pearson Education, Inc.

# Translation is mediated by ribosomes

# The Genetic Code

Second letter

| | U | C | A | G | |
|---|---|---|---|---|---|
| **U** | UUU, UUC } Phe; UUA, UUG } Leu | UCU, UCC, UCA, UCG } Ser | UAU, UAC } Tyr; UAA Stop; UAG Stop | UGU, UGC } Cys; UGA Stop; UGG Trp | U C A G |
| **C** | CUU, CUC, CUA, CUG } Leu | CCU, CCC, CCA, CCG } Pro | CAU, CAC } His; CAA, CAG } Gln | CGU, CGC, CGA, CGG } Arg | U C A G |
| **A** | AUU, AUC } Ile; AUA; AUG Met | ACU, ACC, ACA, ACG } Thr | AAU, AAC } Asn; AAA, AAG } Lys | AGU, AGC } Ser; AGA, AGG } Arg | U C A G |
| **G** | GUU, GUC, GUA, GUG } Val | GCU, GCC, GCA, GCG } Ala | GAU, GAC } Asp; GAA, GAG } Glu | GGU, GGC, GGA, GGG } Gly | U C A G |

First letter

Third letter

Some evolutionary thoughts

http://courses.bio.indiana.edu/L104-Bonner/F09/imagesF09/L23/Ribosomes.html

# Transcription vs. Translation: Lac operon control region

Promoter

Operator
**O**

Coding region
**Z**

Met-Thr-Met…

GAAUUGUGAGCGGAUAACAAUUUCACAC AGGAAAC AGC AUG ACCAUG

→ mRNA

CTTCCGGCTCG TATGTT GTGTGGAATTGTGAGCGGATAACAATTTCACAC AGGAAAC AGCT ATG ACCATG
GAAGGCCGAGC ATACAA CACACCTTAACACTCGCCTATTGTTAAAGTGTG TCCTTTG TCGATACTGGTAC 5′

Operator

Signals are not perfect matches to consensus sequences

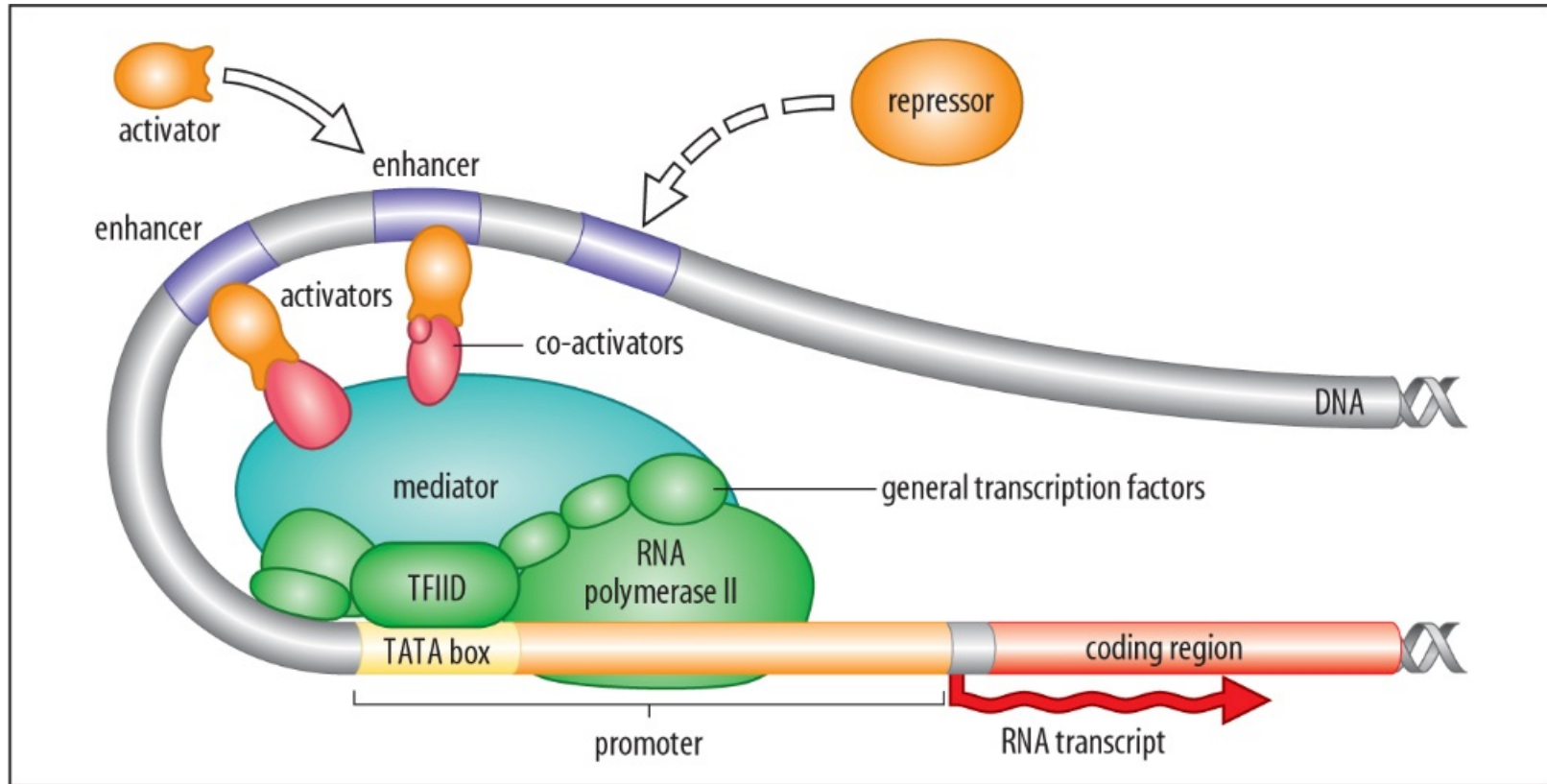RNAP recognizes sequence in promoter DNA

Ribosome recognizes signal in RNA

Seen by computer in DNA

RNAP does not care about codons; it cannot distinguish coding and non-coding DNA.

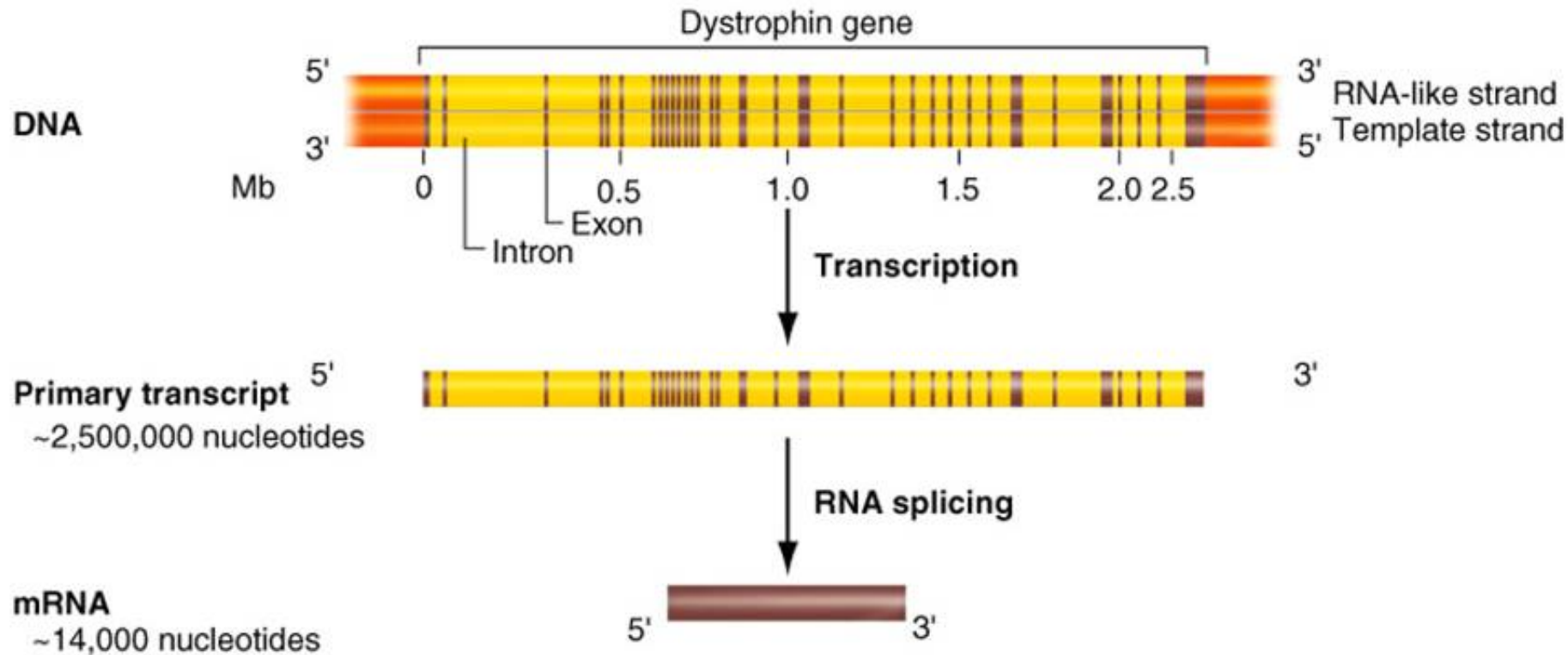1$^o$ transcript == mRNA in bacteria; no splicing, capping, polyA

# Eukaryotic transcription is complex!



- Basal transcriptional regulators
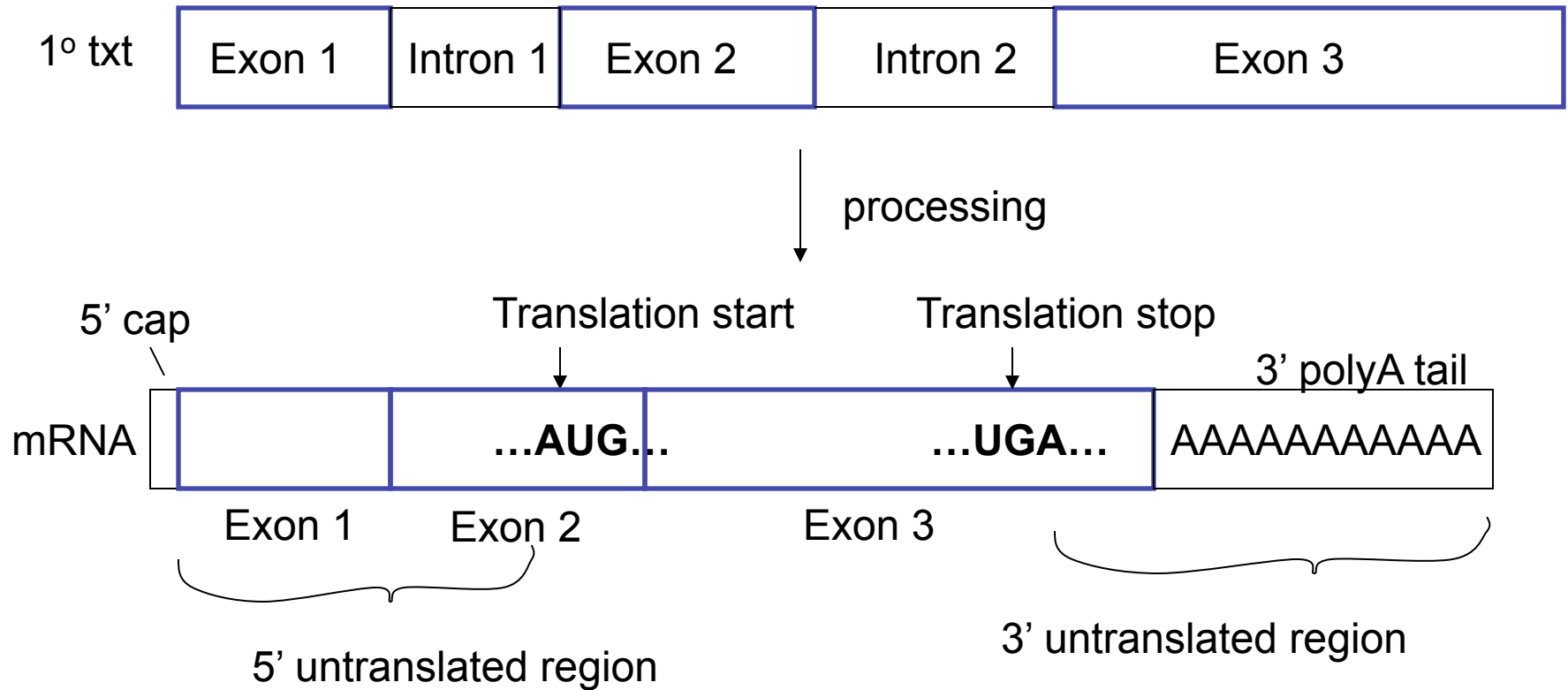- Cell type specific enhancers and repressors

http://www.mun.ca/biology/desmid/brian/BIOL3530/DEVO_10/devo_10.html

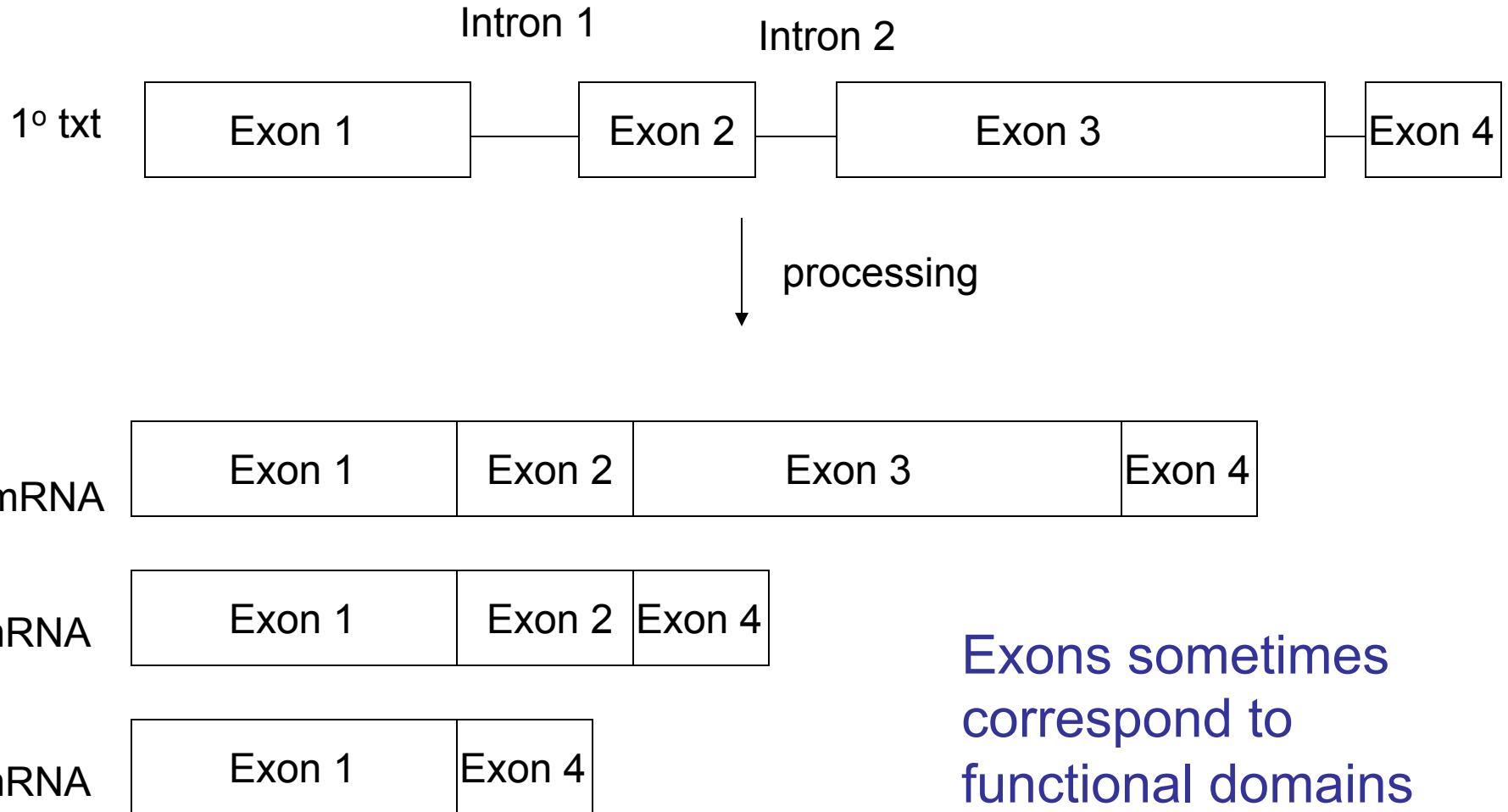# RNA processing in more complex eukaryotes



Hartwell et al. Genetics: From Genes to Genomes

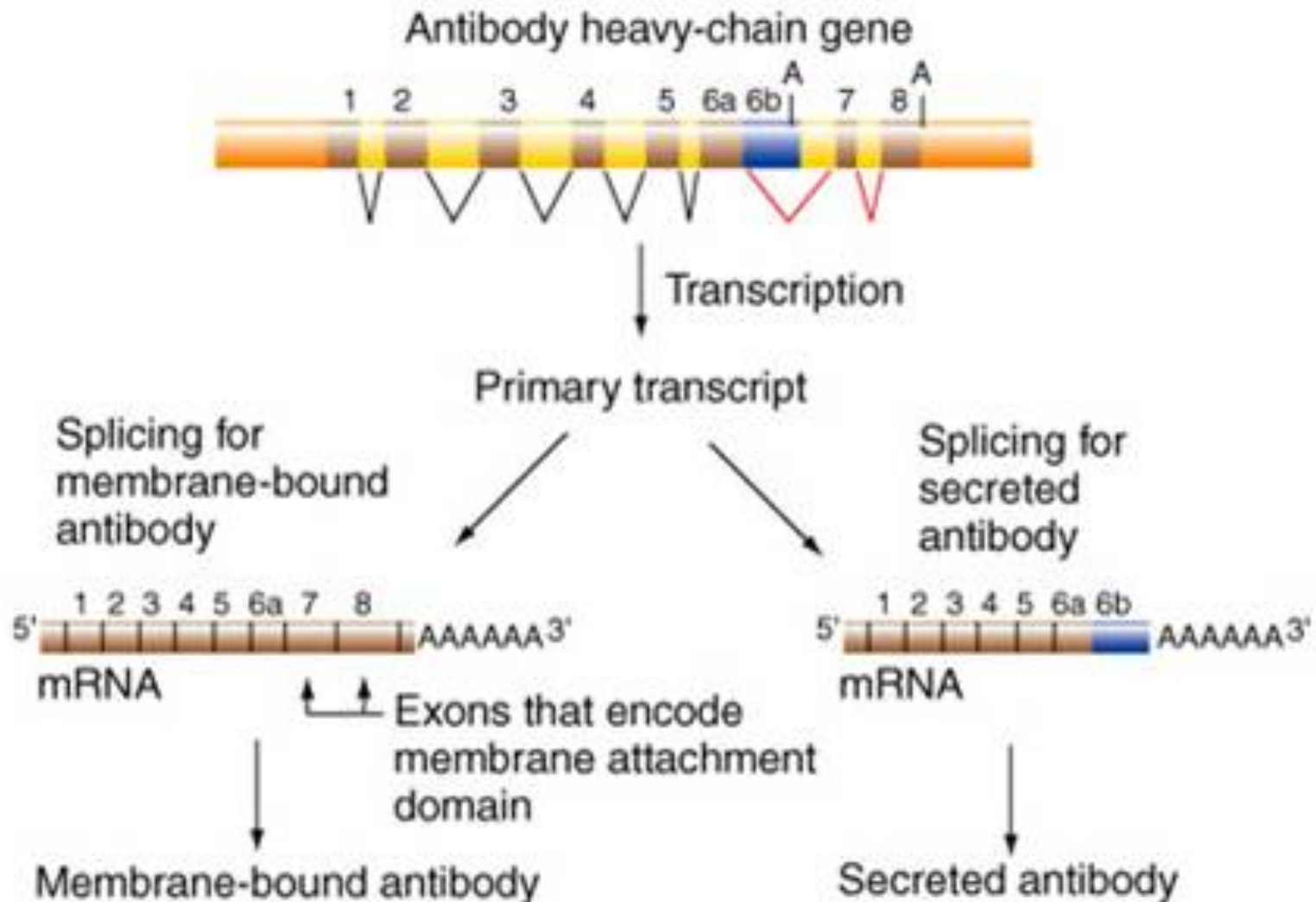# Schematic view of processing of mRNA in eukaryotes

1° txt | Exon 1 | Intron 1 | Exon 2 | Intron 2 | Exon 3

processing

5' cap

Translation start    Translation stop

3' polyA tail

mRNA |  | **...AUG...** | **...UGA...** | AAAAAAAAAAA

Exon 1    Exon 2    Exon 3

5' untranslated region

3' untranslated region

Is the entire transcript translated?

# Alternative splicing

Intron 1  Intron 2

1° txt  | Exon 1 | Exon 2 | Exon 3 | Exon 4 |

↓ processing

mRNA | Exon 1 | Exon 2 | Exon 3 | Exon 4 |

mRNA | Exon 1 | Exon 2 | Exon 4 |

mRNA | Exon 1 | Exon 4 |

Exons sometimes correspond to functional domains of proteins

# Different splice forms can function very differently in the cell



Hartwell et al. Genetics: From Genes to Genomes

# The Genetic Code

Second letter

|  | U | C | A | G |  |
|---|---|---|---|---|---|
| **U** | UUU ⎫ Phe<br>UUC ⎭<br>UUA ⎫ Leu<br>UUG ⎭ | UCU ⎫<br>UCC ⎪ Ser<br>UCA ⎪<br>UCG ⎭ | UAU ⎫ Tyr<br>UAC ⎭<br>UAA  Stop<br>UAG  Stop | UGU ⎫ Cys<br>UGC ⎭<br>UGA  Stop<br>UGG  Trp | U<br>C<br>A<br>G |
| **C** | CUU ⎫<br>CUC ⎪ Leu<br>CUA ⎪<br>CUG ⎭ | CCU ⎫<br>CCC ⎪ Pro<br>CCA ⎪<br>CCG ⎭ | CAU ⎫ His<br>CAC ⎭<br>CAA ⎫ Gln<br>CAG ⎭ | CGU ⎫<br>CGC ⎪ Arg<br>CGA ⎪<br>CGG ⎭ | U<br>C<br>A<br>G |
| **A** | AUU ⎫<br>AUC ⎬ Ile<br>AUA ⎭<br>AUG  Met | ACU ⎫<br>ACC ⎪ Thr<br>ACA ⎪<br>ACG ⎭ | AAU ⎫ Asn<br>AAC ⎭<br>AAA ⎫ Lys<br>AAG ⎭ | AGU ⎫ Ser<br>AGC ⎭<br>AGA ⎫ Arg<br>AGG ⎭ | U<br>C<br>A<br>G |
| **G** | GUU ⎫<br>GUC ⎪ Val<br>GUA ⎪<br>GUG ⎭ | GCU ⎫<br>GCC ⎪ Ala<br>GCA ⎪<br>GCG ⎭ | GAU ⎫ Asp<br>GAC ⎭<br>GAA ⎫ Glu<br>GAG ⎭ | GGU ⎫<br>GGC ⎪ Gly<br>GGA ⎪<br>GGG ⎭ | U<br>C<br>A<br>G |

First letter (left axis) · Third letter (right axis)

# Amino acids fall into different classes

Hydrophobic:
Nonpolar side chains often found in protein core, transmembrane regions

Hydrophilic:
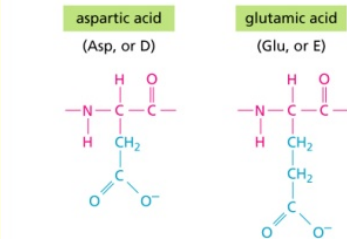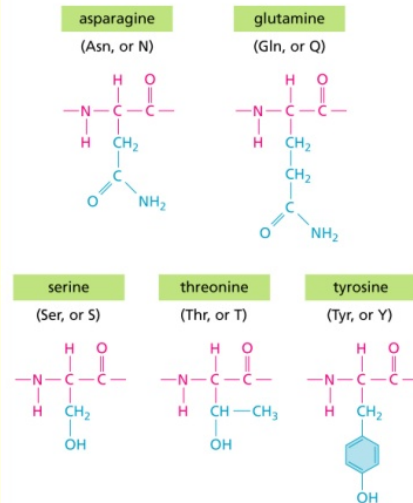Polar and charged side chains often found nearer protein surface, interact with water

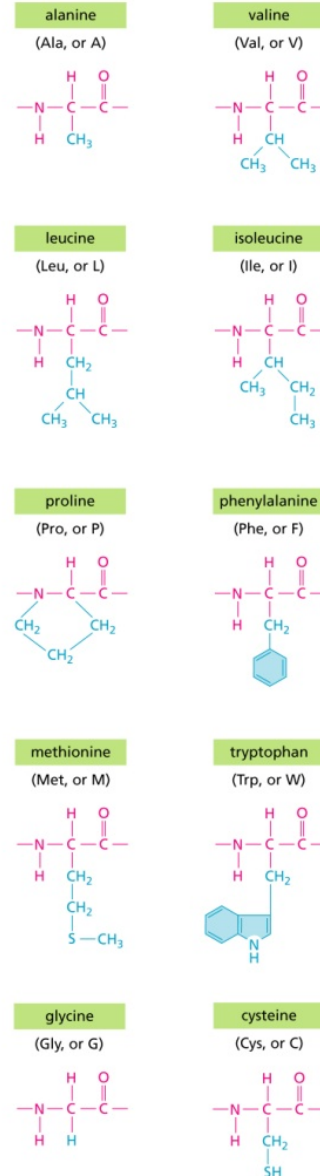# Four levels of protein structure



PRIMARY

N terminus–...MYCATISEATINGFISHANDMEATANDWATER...–C terminus

SECONDARY

TERTIARY

QUATERNARY

# Making a polypeptide chain: primary structure

# Secondary structure: $\alpha$-helices

(A)

(B)

amino acid side chain

oxygen

H-bond

carbon

hydrogen

nitrogen

0.54 nm

carbon

nitrogen

Ala, Glu, Leu, Met 'like' $\alpha$-helices

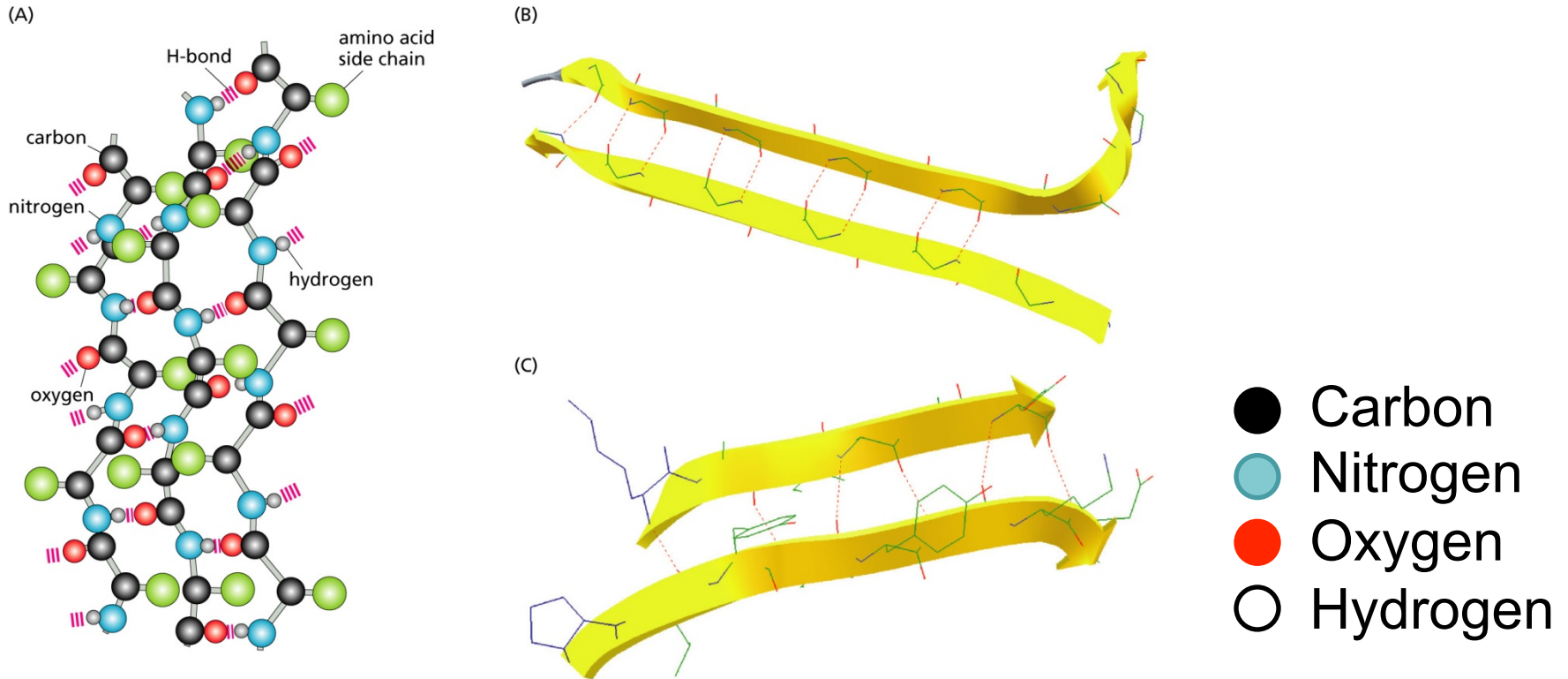Pro rarely found in helices

Gly, Tyr, also poor helix formers

● Carbon
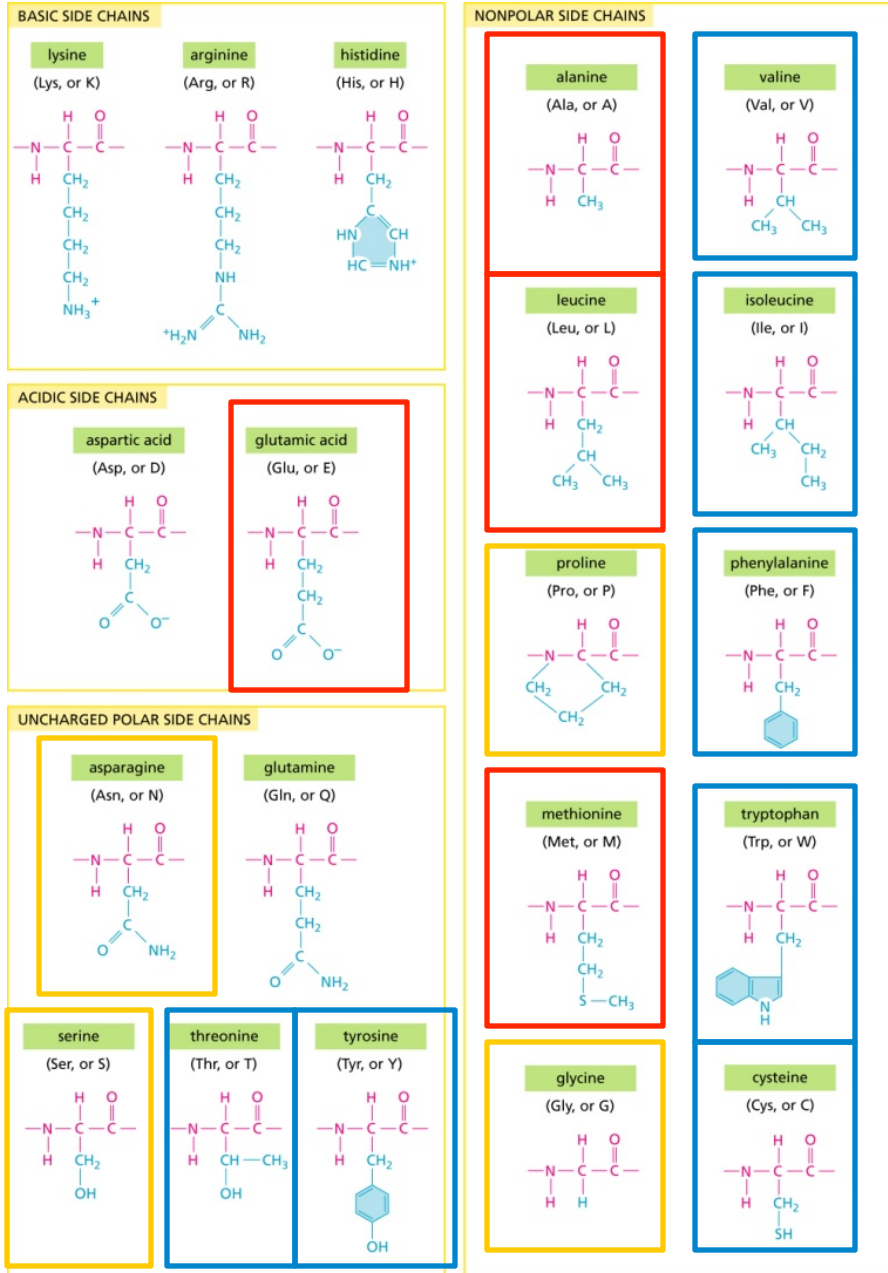● Nitrogen
● Oxygen
○ Hydrogen

# Secondary structure: β-strands



● Carbon
● Nitrogen
● Oxygen
○ Hydrogen

Val, Ile, Tyr, Cys, Trp, Phe, Thr 'like' β-strands
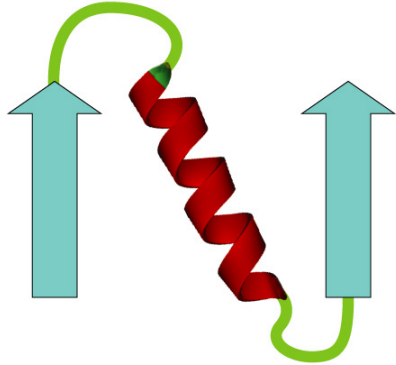
# Amino acids fall into different classes

$\alpha$-helix formers

$\beta$-strand formers

Turn segments

Gly, Tyr, and especially Pro are poor $\alpha$-helix formers

Some evolutionary thoughts

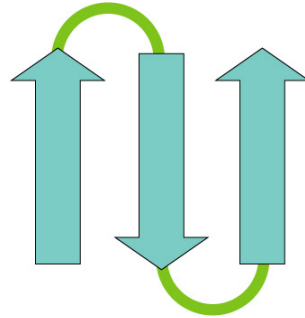# Supersecondary structure


(A)


(B)


(C)


(D)
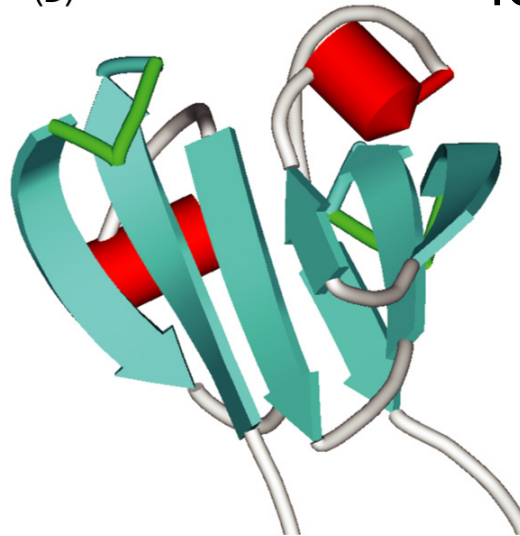
A) βαβ repeat
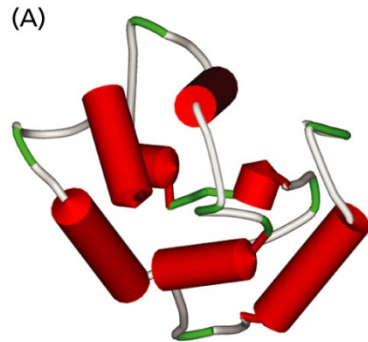B) β meander
C) Greek Key
D) Greek Key, β-crystallin

'only' 2,000 fold families
for 35,000 structures

# Prediction Examples: Known Structures



(A)

1B8C

(B)

1BKB

(C)

1CJW

(D)

1CT5

# α helix / β sheet

Knowledge based | Neural net

X-ray:
GOR IV:
GOR V:
PredS:
PredM:
Zpred:
PROF:
NNSSP:
PHD:
PSIPRED:
Jnet:

# α + β and α / β fold

X-ray:
GOR IV:
GOR V:
PredS:
PredM:
Zpred:
PROF:
NNSSP:
PHD:
PSIPRED:
Jnet:

X-ray:
GOR IV:
GOR V:
PredS:
PredM:
Zpred:
PROF:
NNSSP:
PHD:
PSIPRED:
Jnet:

# Some evolutionary thoughts

- Mutations occur at random in _____
- Are all mutations bad?  The Genetic Code
- Are some more likely to affect protein function than others?  Amino acids
- Which ones might be selected against?
- Which ones might be selected for?
- How does this relate to sequence alignment?

# Alignments

# Outline

- Why align sequences?
- Principles of alignments
- Performing alignments
- Scoring alignments: substitution matrices

# Why align sequences?

- To determine whether sequences are **homologous**: this is, they are derived from a common ancestor sequence

- Are two sequences so similar that we can conclude they are homologous; or is the similarity just due to chance?

- Databases are so large now that lots of surprising things may occur by chance

# Homologs, orthologs, paralogs

Homologs

Orthologs     Paralogs

- Homologs: Two genes with a common ancester
- Orthologs: Homologous genes arising through **speciation**
- Paralogs: Homologous genes arising through **duplication**

# Orthologs vs. paralogs



(A)

All 5 genes are homologs.
Which are orthologs, and which are paralogs?
Which are most likely to function similarly?

# Why align sequences?
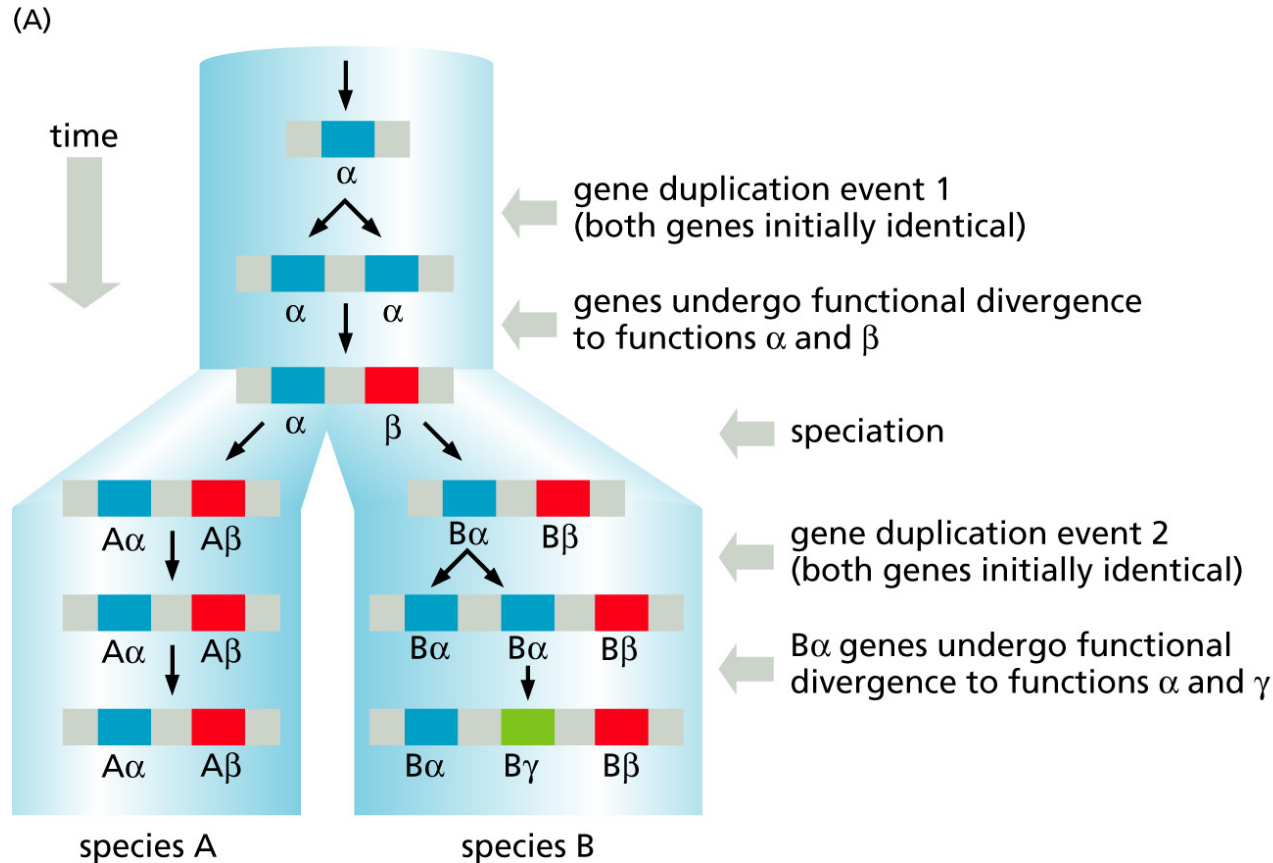
- To determine whether sequences are **homologous**: this is, they are derived from a common ancestor.

- We hope that **orthologs** will have similar functions.

- Therefore, we make alignments to understand more about how proteins function.

- Alignments also allow us to infer how closely related proteins are

  – Important application: evolution / transmission of disease (e.g. flu, malaria)

# Principles of Alignment

- There are almost always multiple ways to align two sequences

- Which way is best? Is the alignment due to more than just chance?

- Need a way to score

# Principles of Alignment

**THISISASEQUENCE**

**THATSEQUENCE**

**THISISASEQUENCE**
**THAT---SEQUENCE**

**THISISASEQUENCE**
**TH---ATSEQUENCE**

**THISISA-SEQUENCE**
**TH----ATSEQUENCE**

Which one is best?

# Scoring alignments

- Simplest score: % identity
- Number of matches/length of match

**THISISASEQUENCE**   **THISISASEQUENCE**

**THAT---SEQUENCE**   **TH---ATSEQUENCE**

    10/15 = 66.7%          10/15 = 66.7%

**THISISA-SEQUENCE**

**TH----ATSEQUENCE**

    11/16 = 68.8%

Is the third alignment the best?

# Inserting gaps in alignments

- Gaps in alignment presumably due to insertion or deletion in one sequence during evolution
- Too many gaps => meaningless alignment
- Gap penalty
- Gap extension penalty (lower)
- Penalties should be adjusted based on whether you are looking for highly related (high penalty) or more distant sequences (low penalty)
- Gap penalty may depend on residue opposite gap (e.g. tryptophan – large penalty), but typically this is ignored

# Identity is not always satisfactory: Substitution matrices

- Use real data to derive scoring matrix

- Genuine matches need not be identical

- How likely is it that an amino acid at a particular position substituted for another amino acid at that position during evolution?

- Substitutions of amino acids with similar physicochemical properties (e.g. size, charge, hydrophobicity) are more likely to conserve function
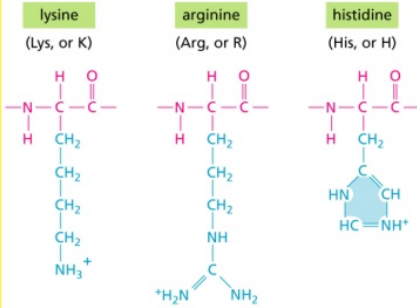
**BASIC SIDE CHAINS**
lysine (Lys, or K); arginine (Arg, or R); histidine (His, or H)

**ACIDIC SIDE CHAINS**
aspartic acid (Asp, or D); glutamic acid (Glu, or E)

**UNCHARGED POLAR SIDE CHAINS**
asparagine (Asn, or N); glutamine (Gln, or Q); serine (Ser, or S); threonine (Thr, or T); tyrosine (Tyr, or Y)

**NONPOLAR SIDE CHAINS**
alanine (Ala, or A); valine (Val, or V); leucine (Leu, or L); isoleucine (Ile, or I); proline (Pro, or P); phenylalanine (Phe, or F); methionine (Met, or M); tryptophan (Trp, or W); glycine (Gly, or G); cysteine (Cys, or C)

# Amino acids fall into different classes

☐ $\alpha$-helix formers

☐ $\beta$-strand formers

☐ Turn segments

Gly, Tyr, and especially Pro are poor $\alpha$-helix formers
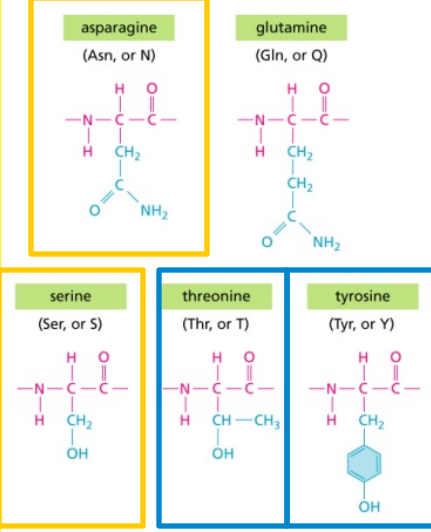
# PAM120 Substitution Matrix

(B)

|   | C | S | T | P | A | G | N | D | E | Q | H | R | K | M | I | L | V | F | Y | W |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| C | 9 | | | | | | | | | | | | | | | | | | | |
| S | -1 | 3 | | | | | | | | | | | | | | | | | | |
| T | -3 | 2 | 4 | | | | | | | | | | | | | | | | | |
| P | -3 | 1 | -1 | 6 | | | | | | | | | | | | | | | | |
| A | -3 | 1 | 1 | 1 | 3 | | | | | | | | | | | | | | | |
| G | -5 | 1 | -1 | -2 | 1 | 5 | | | | | | | | | | | | | | |
| N | -5 | 1 | 0 | -2 | 0 | 0 | 4 | | | | | | | | | | | | | |
| D | -7 | 0 | -1 | -2 | 0 | 0 | 2 | 5 | | | | | | | | | | | | |
| E | -7 | -1 | -2 | -1 | 0 | -1 | 1 | 3 | 5 | | | | | | | | | | | |
| Q | -7 | -2 | -2 | 0 | -1 | -3 | 0 | 1 | 2 | 6 | | | | | | | | | | |
| H | -4 | -2 | -3 | -1 | -3 | -4 | 2 | 0 | -1 | 3 | 7 | | | | | | | | | |
| R | -4 | -1 | -2 | -1 | -3 | -4 | -1 | -3 | -3 | 1 | 1 | 6 | | | | | | | | |
| K | -7 | -1 | -1 | -2 | -2 | -3 | 1 | -1 | -1 | 0 | -2 | 2 | 5 | | | | | | | |
| M | -6 | -2 | -1 | -3 | -2 | -4 | -3 | -4 | -4 | -1 | -4 | -1 | 0 | 8 | | | | | | |
| I | -3 | -2 | 0 | -3 | -1 | -4 | -2 | -3 | -3 | -3 | -4 | -2 | -2 | 1 | 6 | | | | | |
| L | -7 | -4 | -3 | -3 | -3 | -5 | -4 | -5 | -4 | -2 | -3 | -4 | -4 | 3 | 1 | 5 | | | | |
| V | -2 | -2 | 0 | -2 | 0 | -2 | -3 | -3 | -3 | -3 | -3 | -3 | -4 | 1 | 3 | 1 | 5 | | | |
| F | -6 | -3 | -4 | -5 | -4 | -5 | -4 | -7 | -6 | -6 | -2 | -4 | -6 | -1 | 0 | 0 | -3 | 8 | | |
| Y | -1 | -3 | -3 | -6 | -4 | -6 | -2 | -5 | -4 | -5 | -1 | -6 | -6 | -4 | -2 | -3 | -3 | 4 | 8 | |
| W | -8 | -2 | -6 | -7 | -7 | -8 | -5 | -8 | -8 | -6 | -5 | 1 | -5 | -7 | -7 | -5 | -8 | -1 | -1 | 12 |

Yellow: small and polar
White: small and nonpolar/hydrophobic
Red: polar or acidic
Blue: basic
Green: larger nonpolar/hydrophobic
Orange: large, aromatic

You will see different types of amino acid groupings!

Score the alignment:

**CSTPEDWLV**
**CTNCDEWDI**

# BLAST

- Basic Local Alignment Search Tool
- Most widely used local alignment algorithm

# BLAST basics

- Starts with short 'words' in the query sequence (default length 3 for proteins, 11 for nucleotides)

- Finds matches in target sequence (using a substitution matrix for proteins, score of match must be above a threshold; for nucleotides, exact match)

- When match is found (two nearby words for proteins), BLAST tries to extend forward and backward to make alignment

- Continues extension until negative scores make the score drop by a critical amount

# How do we know an alignment is significant?

- Expect value (E value): 'the number of times that an alignment as good or better than that found by BLAST would be expected to occur by chance, given the **size of the database searched**'      --From BLAST QuickStart tutorial

$$E = Kmne^{-\lambda S}$$

- S is the score
- Sometimes better to search smaller databases
  - m.n are the lengths of the sequences being compared
  - When comparing to a large database, consider m=query length, n= total length of all sequences in database
- Default E value = 10
- Typically don't consider matches with E > 0.001
- Often see E values like $10^{-36}$

# How do we know an alignment is significant?

- Expect value (E value): 'the number of times that an alignment as good or better than that found by BLAST would be expected to occur by chance, given the **size of the database searched**'      --From BLAST QuickStart tutorial

$$E = Kmne^{-\lambda S}$$

- The parameters *K* and *lambda* can be thought of simply as natural scales for the search space size and the scoring system respectively

- For details: http://www.ncbi.nlm.nih.gov/BLAST/tutorial/Altschul-1.html#head2

# A word about nucleotides

- Why do we typically align proteins and not nucleotide sequences?

- Possible to align nucleotides, but more difficult
    - Only 4 bases
    - Matches more likely to occur by chance
    - Amino acids more conserved over evolution (genetic code is degenerate; DNA less conserved)

- Use protein sequence when available

- Scoring matrices much simpler – e.g. BLASTn uses +2 for a match, -3 for a mismatch

- When would we have to align nucleotides?

# Alignment example

- We'll try aligning the C. elegans protein MIG-10 to the Refseq database

- See if we can decide whether there are homologs of MIG-10 in other species

- Initially, we'll judge by % identity, % coverage, and E value

- DM0C32V6014 (Expires on 02-11 20:02PM)

# BLAST Output

# An individual alignment

# Multiple alignment:
# COBALT Output (NCBI)

```
NP_001038978    152   SLDDITAQLEQASLSMDEAAQQ-      SLVEDPKPLVTNQHRRTASAGTVSDAEARSISNSSRSSITSA-ASSMDSLDIDK    226

NP_001021248    281   ----DSLNTPSPTQVSPRNGELNAEEAKAQKIRQALEKMKEAKVTKIFVKFFVEDGEALQLLIDERWTVADTLKQLAEKN    356
XP_002641939    152   ----DSLNTPSPTQVSPRTGELNAEEAKSLKIRQALEKMKEAKIIKMLVKFFVEDGQPLQMLIDERWTVADTMKQLAEKN    227
XP_006610888    163   -hKPPQTAMHTGPQQQSHQLMDAASRVKAEKIRLALEKMREASVQKLFIKAFTLDGSGKSLLVDEGMSVAHVCRLLADKN    241
NP_001121716    244   --RGQENETQSQNQSQTSTEEEQAAKAKAEKIRVALEKIKEAQVKKLVIRVHMSDESSKTMMVDERQTVRQVLDSLLDKS    321
XP_004033136    227   vtRPQELDLT--HQGQPITEEEQAAKLKAEKIRVALEKIKEAQVKKLVIRVHMSDDSSKTMMVDERQTVRQVLDNLMDKS    304
NP_001006357    139   --------PPPPPPPEPLSQEEQEARAKADKIKLALEKLKEAKIKKLVVKVHMYDNSTKSLMVDERQVTRDVLDNLFEKT    210
NP_001038978    227   vtRPQELDLTT-HQGQPITEEEQAAKLKAEKIRVALEKIKEAQVKKLVIRVHMSDDSSKTMMVDERQTVRQVLDNLMDKS    305

NP_001021248    357   HIALMEDHCIVEEYPELYIKRVYEDHEKVVENIQMWVQDSPNK-LYFMRRPDKYAFISRPELYLLT---PKTSDHMEIPS    432
XP_002641939    228   HIALMEDHCIVEEYPELYIKRVYEDHEKVVENITMWVQDSPNK-LYFMRRPDKYTFISRPELYLLT---PKTSDHMEIPP    303
XP_006610888    242   HVPMDPKWTVVEHLPDLFMERVYEDHELLVENLLLWTRDSKNK-LLFVERPEKTQLFLTPERFLLG---------LSDRS    311
NP_001121716    322   HCGYSPDWALVETIPELQMERIFEDHENLVENLLNWTRDSQNK-LMFIERIEKYALFKNPQNYLL---GRKETSEMADRN    397
XP_004033136    305   HCGYSLDWSLVETVSELQMERIFEDHENLVENLLNWTRDSQNK-LIFMERIEKYALFKNPQNYLL---GKKETAEMADRN    380
NP_001006357    211   HCDCSVDWCLYEVYPELQIERFFEDHENVVEVLSDWTRDSENKvLFL-EKKEKYALFKNPQNFYLAnkGKNESKEMNDKS    289
NP_001038978    306   HCGYSLDWSLVETISELQMERIFEDHENLVENLLNWTRDSQNK-LIFMERIEKYALFKNPQNYLL---GKKETAEMADRN    381

NP_001021248    433   [8]KQKFVSEYFHREPVVPPEMEGFLYLKSDGRKSWKKHYFVLRPSGLYYAPKSKKPTTKDLTCLMNLHSNQVYTGIGWE    517
XP_002641939    304   [8]KQKFVHDYFNREPVVPPEMEGFLYLKSDGRKSWKKHYFVLRPSGLYYAPKSKKPTTKDLTCLMNLHSNQVYTGIGWE    388
XP_006610888    312   [8]RNILLEEFFSSSNVGVPEVEGPLYLKSDSKKGWKRYHFILRASGLYYWPKEKARTARDLVCLATFDVNQIYYGIGWK    396
NP_001121716    398      KEALLEECFCGSSVSVPEIEGVLWLKEDGKKSWKRRYFLLRASGIYFVPKGKAKASRDLVCFLQLDHVNVYYGQDYR    474
XP_004033136    381      KEVLLEECFCGSSVTVPEIEGVLWLKDDGKKSWKRRYFLLRASGIYYVPKGKAKVSRDLVCFLQLDHVNVYYGQDYR    457
NP_001006357    290      KEALLEESFCGASVIVPELEGALYLKEDGKKSWKRRYFLLRASGIYYVPKGKTKTSRDLMCFIQFENMNVYYGSQHK    366
NP_001038978    382      KEVLLEECFCGSSVTVPEIEGVLWLKDDGKKSWKRRYFLLRASGIYYVPKGKAKVSRDLVCFLQLDHVNVYYGQDYR    458
```

# MIG-10 in Jalview (from EBI)

# Jalview colored by conservation

# Gene regulatory networks

How is gene expression controlled?

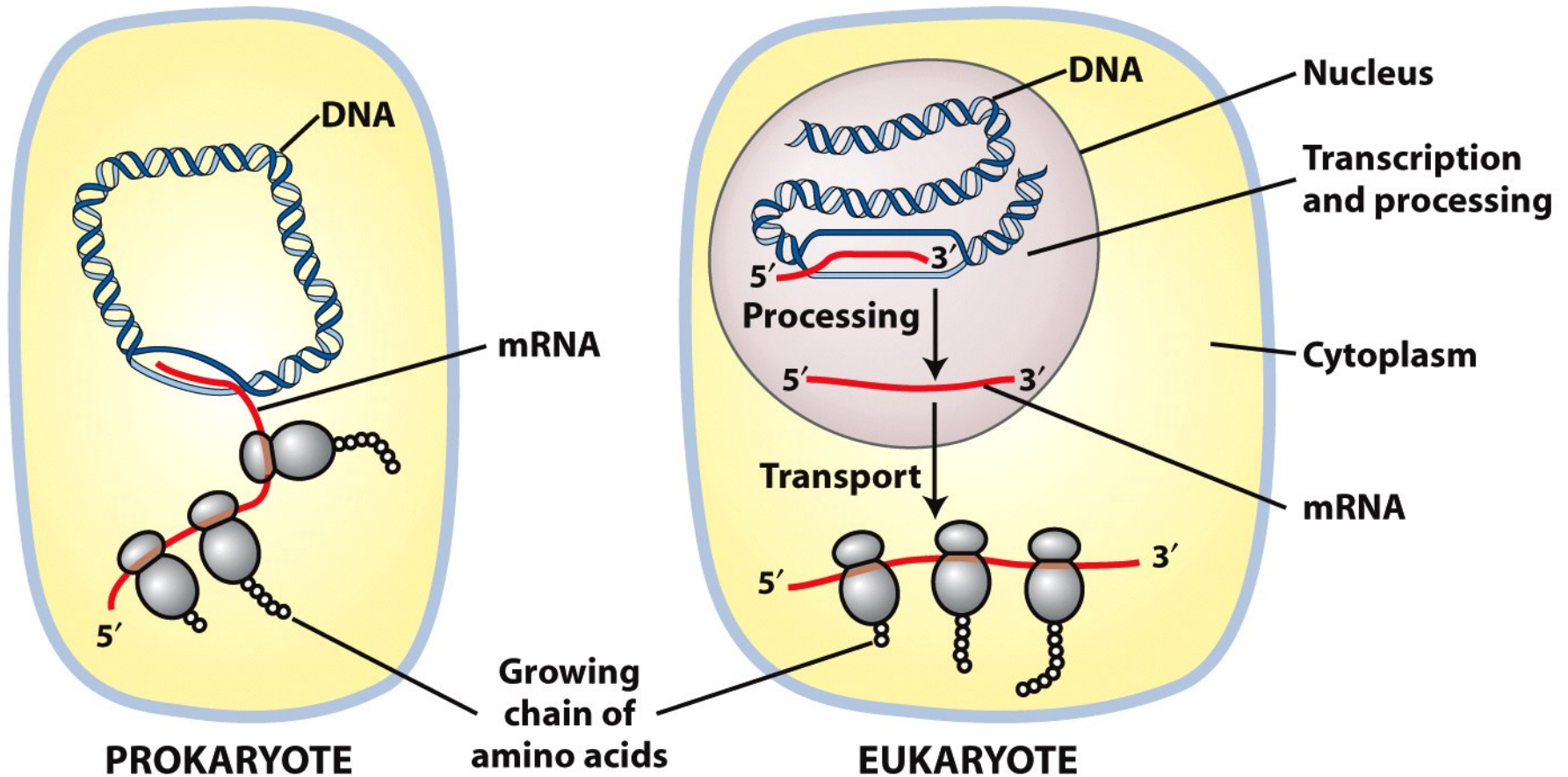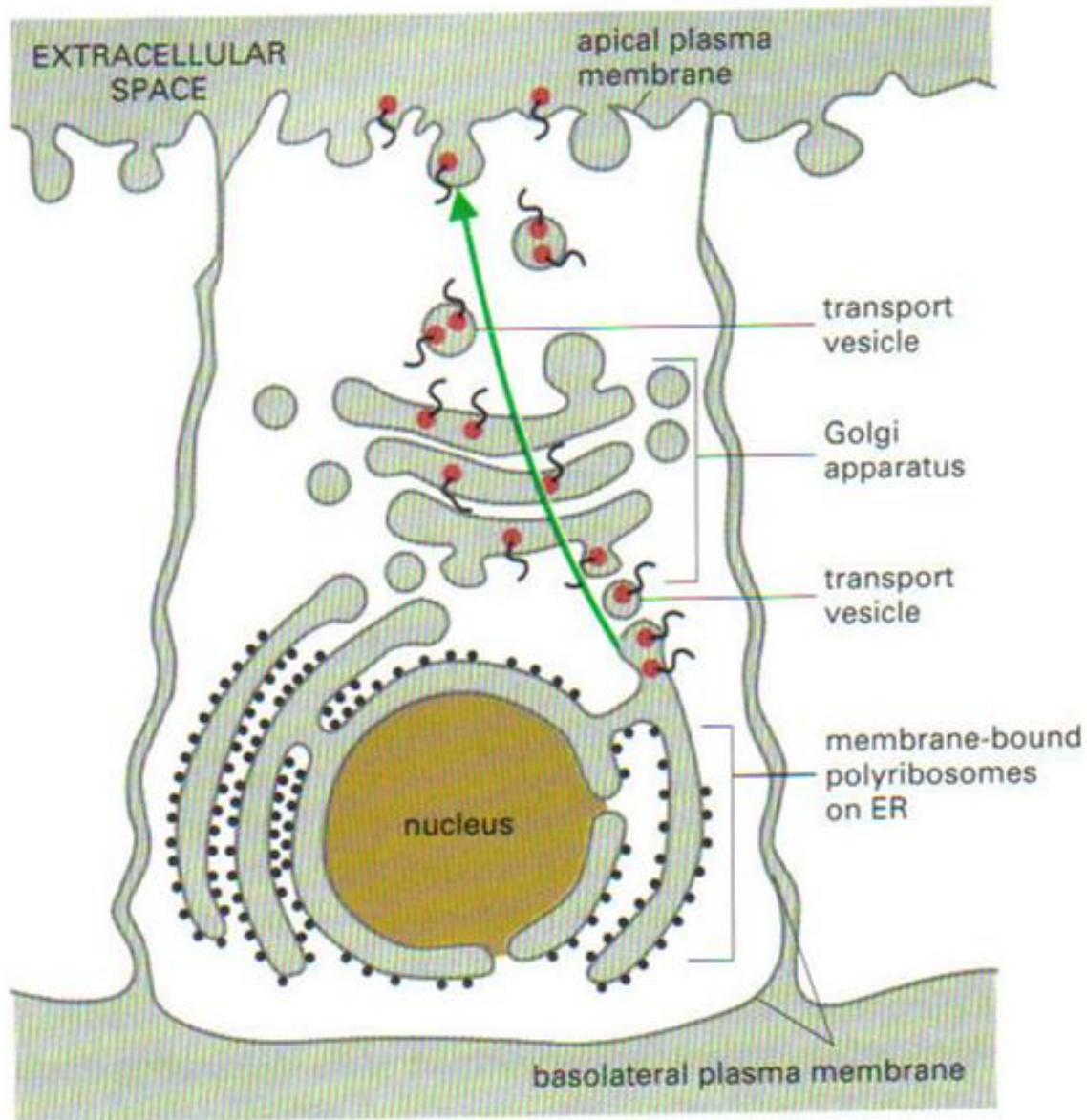# Prokaryotic and eukaryotic transcription and translation



**Figure 8-11**
*Introduction to Genetic Analysis, Ninth Edition*
© 2008 W. H. Freeman and Company

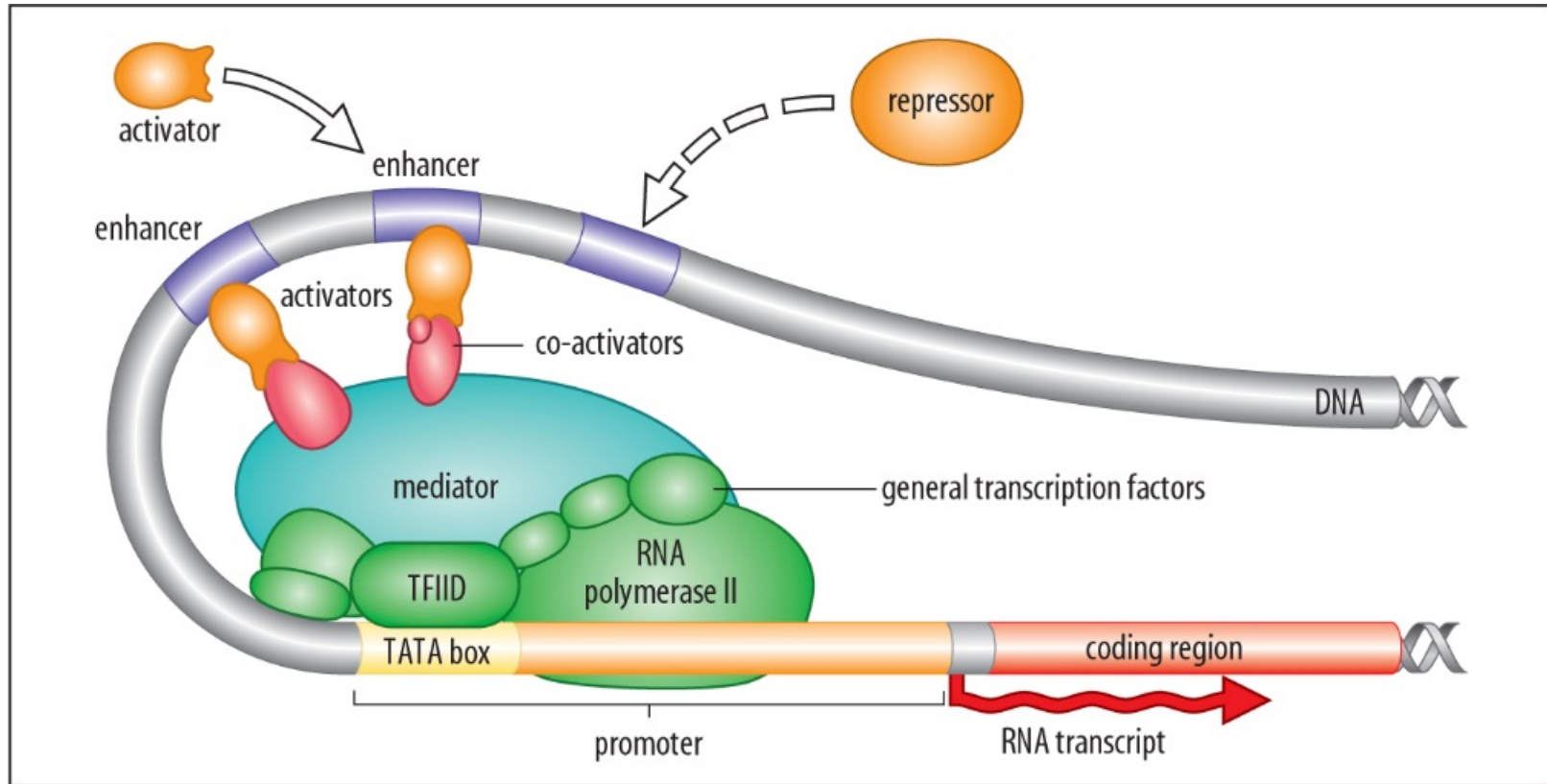Griffeths, Introduction to Genetic Analysis, 2008

# Topology of glycosylation



$\zeta$ membrane protein

● N-linked oligosaccharide added in ER

Glycosylation and phosphorylation occur in different cellular compartments…
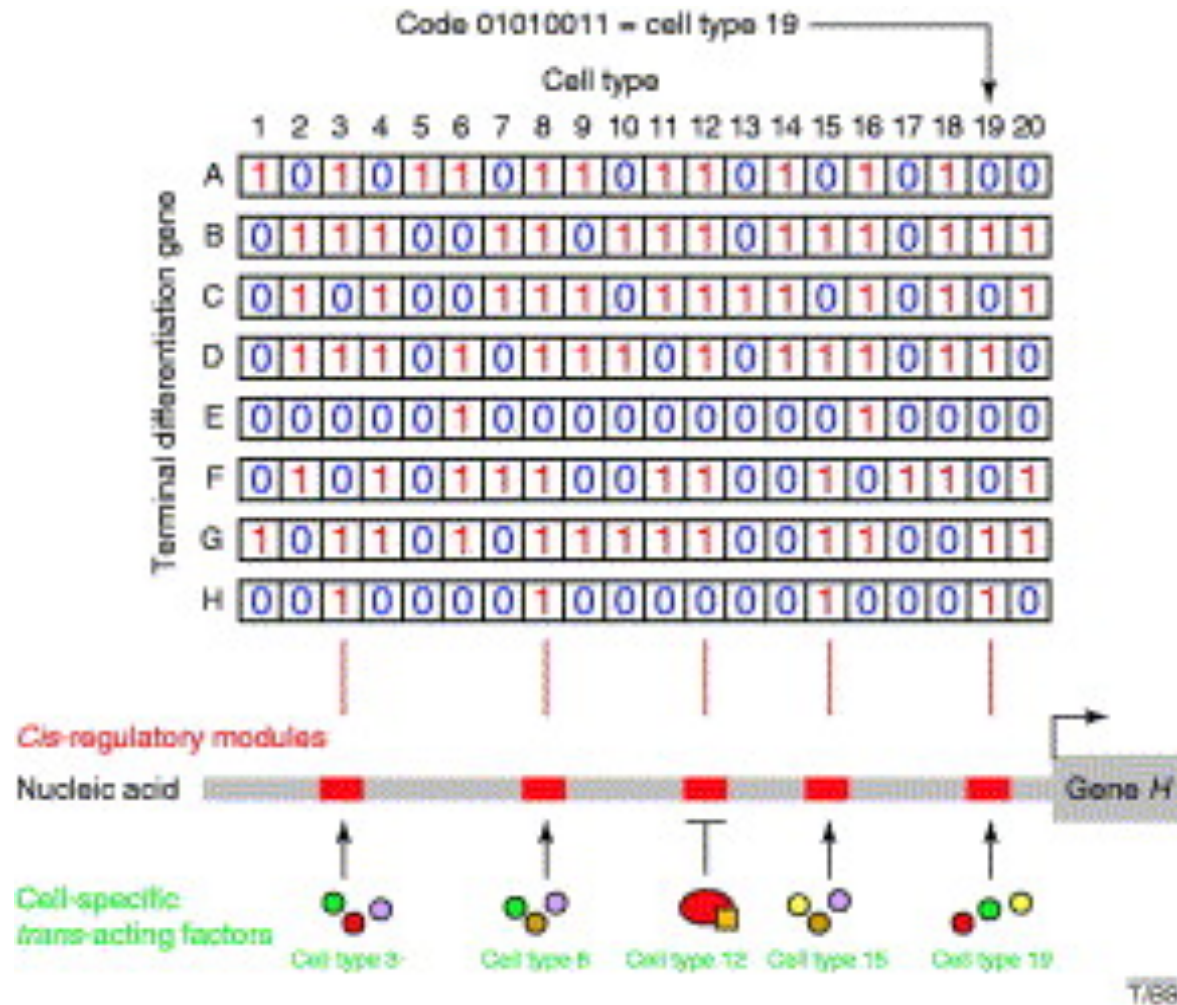
Alberts et al., Molecular Biology of the Cell

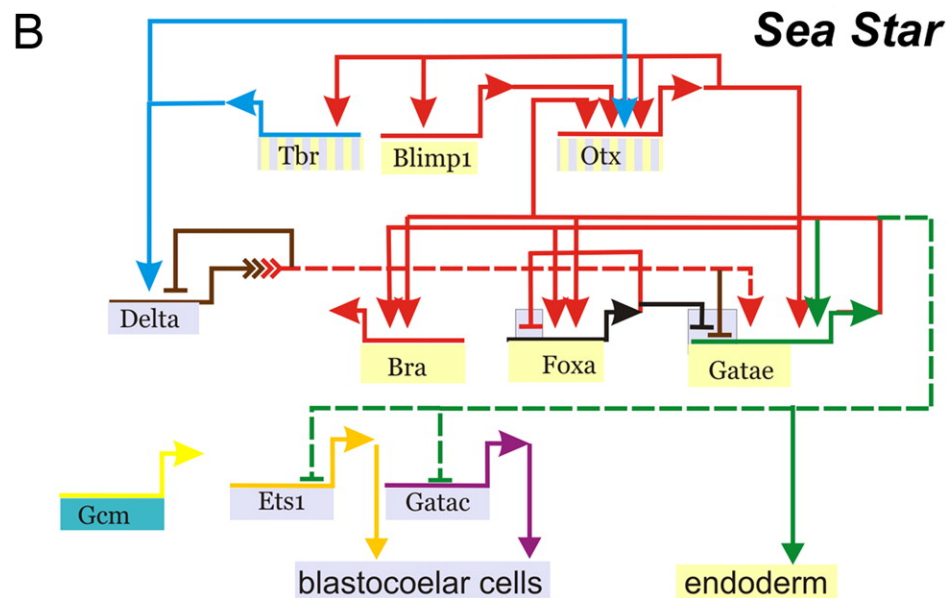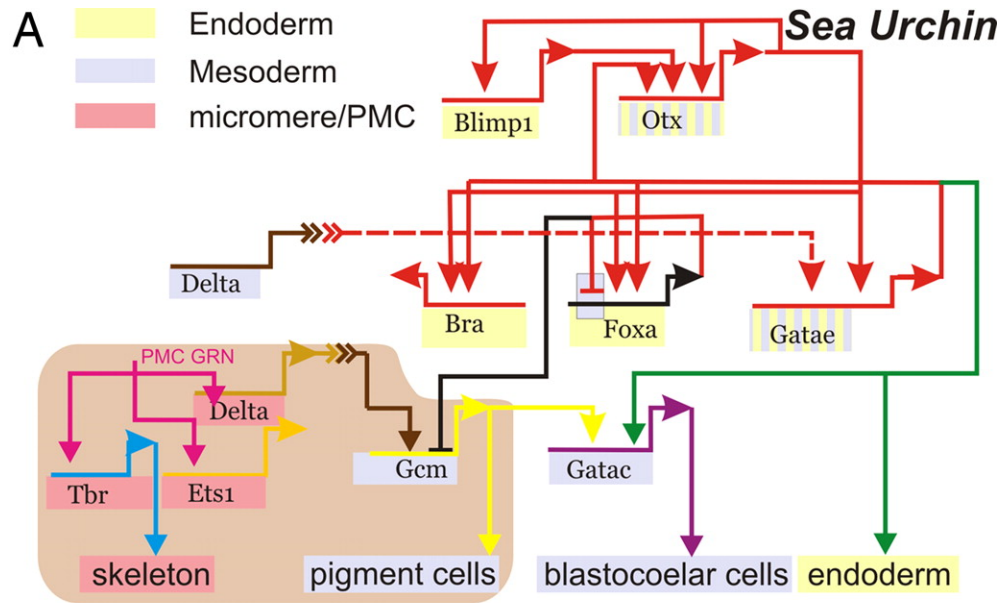# Eukaryotic transcription is complex!



- Basal transcriptional regulators
- Cell type specific enhancers and repressors

http://www.mun.ca/biology/desmid/brian/BIOL3530/DEVO_10/devo_10.html

# 'Gene batteries'



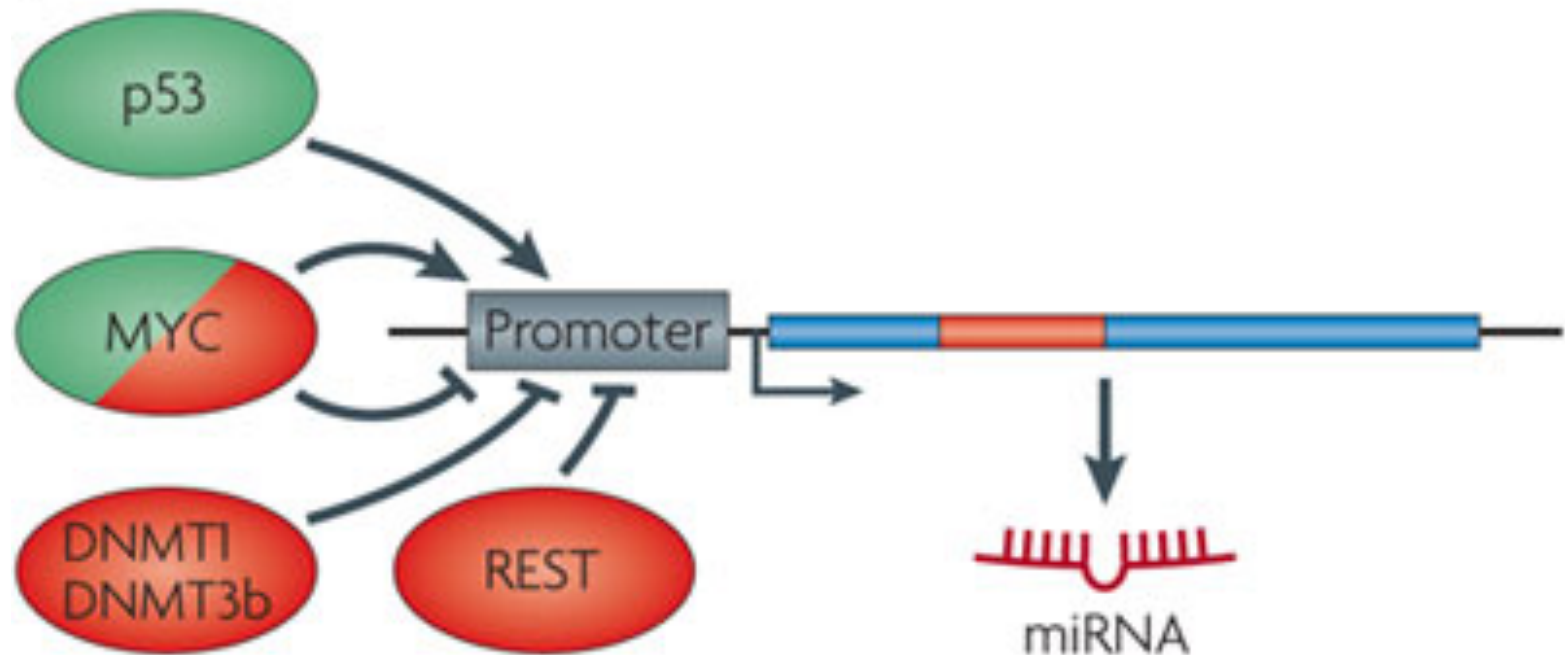Terminal differentiation genes expressed by different cell types

# GRN: Sea urchin vs. Sea star

- All genes except Delta are transcription factors
- Arrows: + regulation
- T's: - regulation
- Colors: ???

http://www.pnas.org/content/104/49/19404/F5.expansion.html

# Micro RNA (miRNA) genes are regulated similarly to protein-encoding genes

# miRNA regulates mRNA



Biosynthetic pathway of microRNA

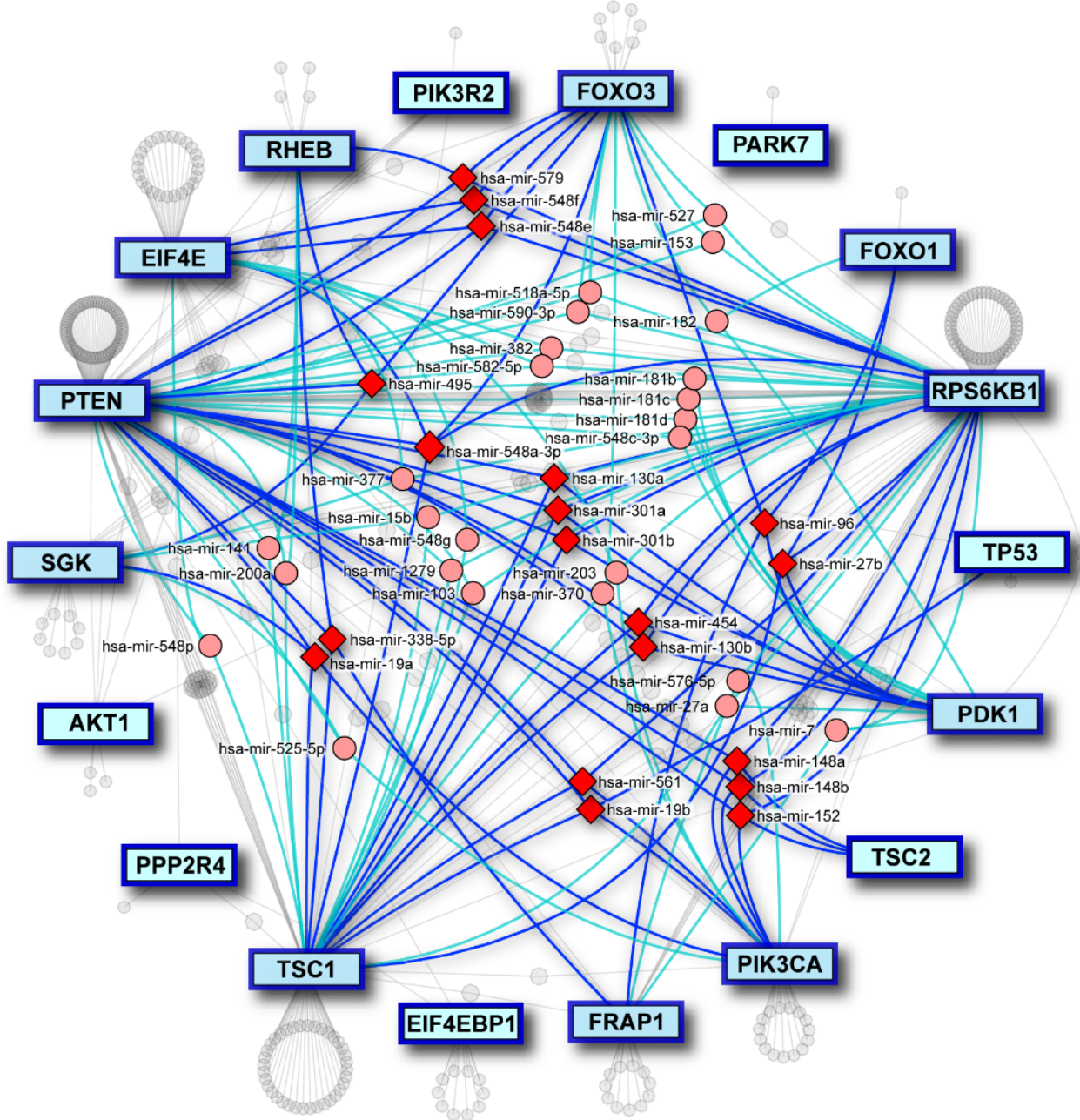# Common logic:
# transcription factors and miRNA

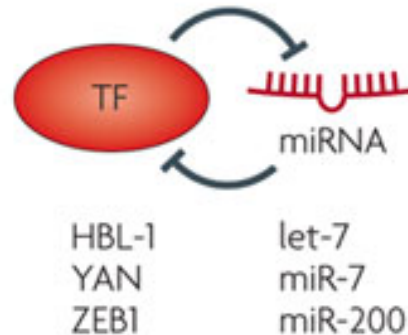Regulation of oncogenes and tumor suppressors by miRNA

http://www.cs.toronto.edu/~juris/home.htm

# Regulatory circuits

**b**

**Unilateral negative feedback loops**

TF → miRNA

| | |
|---|---|
| PITX3 | miR-133b |
| RUNX1 | miR-27a |
| MYB | miR-15a |

**Reciprocal negative feedback loops**

TF ⇄ miRNA

| | |
|---|---|
| HBL-1 | let-7 |
| YAN | miR-7 |
| ZEB1 | miR-200 |

**Double-negative feedback loop**

DIE-1 → lys-6 → COG-1 → miR-273 → DIE-1

- Transcription factors can be + or –
- miRNA typically negative regulator
- Regulatory circuits typically contain both

http://deepbase.sysu.edu.cn/chipbase/tfmiRtargetNetworks

# A microarray experiment



**Data are variable!**

- http://azcc.arizona.edu/research/shared-resources/gsr/services

Samples

Genes (mRNAs)

What are the trees telling you?

http://www.cbioc.com/en/services/bioinformatics-services/

# Appendix: Databases and web sites

- DNA repositories
  - Genbank; NCBI (National Center for Biotechnology Information)
  - EMBL-bank; EMBL-EBI (European Molecular Biology Laboratory-European Bioinformatics Institute)
  - DDBJ (DNA Data Bank of Japan)

# A few important websites

- NCBI
  - http://www.ncbi.nlm.nih.gov/
  - Multiple databases, tools
- EMBL-EBI
  - http://www.ebi.ac.uk/
  - Alignment tools in particular
- ExPASy (Expert Protein Analysis System – Swiss)
  - http://expasy.org/
  - Multiple tools, especially useful for secondary structure prediction

# A sampling of NCBI Databases

- Nucleotides (Genbank entries)
- Gene (Refseq; annotated model organisms)
- Unigene (expression data)
- Homologene
- OMIM (Online Mendelian Inheritance in Man)
- Cn3D (3D structure info)
- CD (conserved domains)
- dbEST (expressed sequence tags)

# Accession numbers vs. gi numbers

- Accession numbers
  - Unique, stable
  - AB123456.2     ← Version 2
  - New version made whenever any change is made to a sequence; old version info is retained
- gi (GenInfo) numbers
  - Assigned sequentially to each nucleotide sequence processed by Genbank
  - 12345678
  - New number whenever any change is made to a sequence; no relationship to old number
  - old number info is retained
- gi and Accession numbers are incremented in parallel

# Refseq database

- Highly annotated entries from subset of organisms
- From the NCBI Glossary
  - RefSeq is the NCBI database of reference sequences; a curated, non-redundant set including genomic DNA contigs, mRNAs and proteins for known genes, and entire chromosomes
- Accession numbers start with two letters and an underscore
  - NM_123456.2 mRNA (version 2)
  - NP_123456.3 protein
  - NT_123456.4 contig
  - NC_000003.7 chromosome
  - XM_ or XP_ entries are 'models' predicted from sequence; no experimental evidence of existence

# Some useful links

- NCBI Education pages
  - Glossary    http://www.ncbi.nlm.nih.gov/books/NBK21106/
  - Tutorials    http://www.ncbi.nlm.nih.gov/education/
- Tour of various NCBI databases and tools
- GenBank sample record explaining info in all fields
  - http://www.ncbi.nlm.nih.gov/Sitemap/samplerecord.html#ModificationsDateB
- Entrez nucleotide and protein FAQs
  - http://www.ncbi.nlm.nih.gov/books/NBK49541/
- The NCBI handbook
  - http://www.ncbi.nlm.nih.gov/books/NBK21101/