

Sequence Analysis

BBSI 2006: Lecture #($\chi+1$)

Takis Benos (2006)



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Molecular Genetics 101



BBSI 2006 26-MAY-2006

© 2006 P. Benos

What is a “gene”?

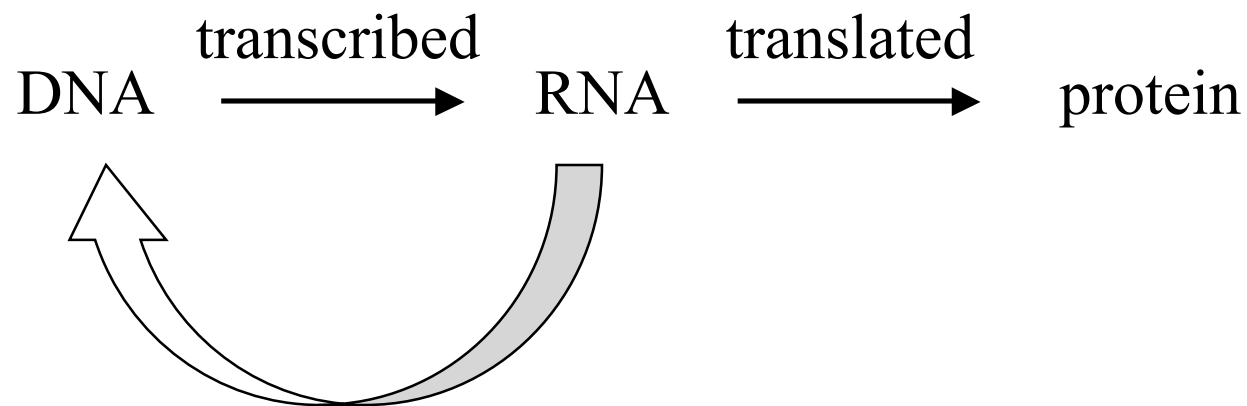
- We cannot define it (but we know it when we see it...)
- A loose definition:

“Gene” is a *DNA/RNA information unit* that is able to perform a function in a cellular environment



Protein coding genes

Central Dogma:



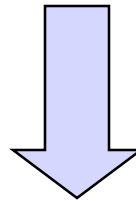
Open Reading Frames (ORFs)

aatagcgaat tttccaacga caaaagctaa atatcgcaaa aacctcagta aaaatcttgc 60
tggagctatt attgctaagt aacatttacc ccctgaagtt aatggatcaa tcaagagaga 120
tgtgggctgt aATGaatecgt cttattgaat taacagggtg gatcgttctt gtcgtttcag 180
tcattcttct tggcgtggcg agtcacattg acaactatca gccacctgaa cagagtgctt 240
cggtaacaaca caagTAAgct ctgcacttgt ggagcgacat gctgcccgtc cgggtgcatg 300
ttttcacttg tcggatatta aaccaggaat ttattatctt gttcgatggt gtaataaa 358



Open Reading Frames (ORFs)

aatagcgaat tttccaacga caaaagctaa atatcgcaaa aacctcagta aaaatcttgc 60
tggagctatt attgctaagt aacatttacc ccctgaagtt aatggatcaa tcaagagaga 120
tgtgggctgt a**ATGa**atcgt **ctt**attgaat **taac**agg**ttg** **gat**cg**ttc**tt **gtc**g**ttt**cag 180
tcatt**ctt**ct **tg**gc**gt**gg**cg** **ag**tcacattg **aca**actatca **gcc**ac**ct**gaa **cag**agt**g**ctt 240
cggtacaaca **caag**TAAgct ctgcacttgt ggagcgcacat gctgcccgtc cgggtgcatg 300
ttttcacttg **tc**gga**ta**ta **aac**caggaat ttattatctt gttcgatgtt gtaataaa 358



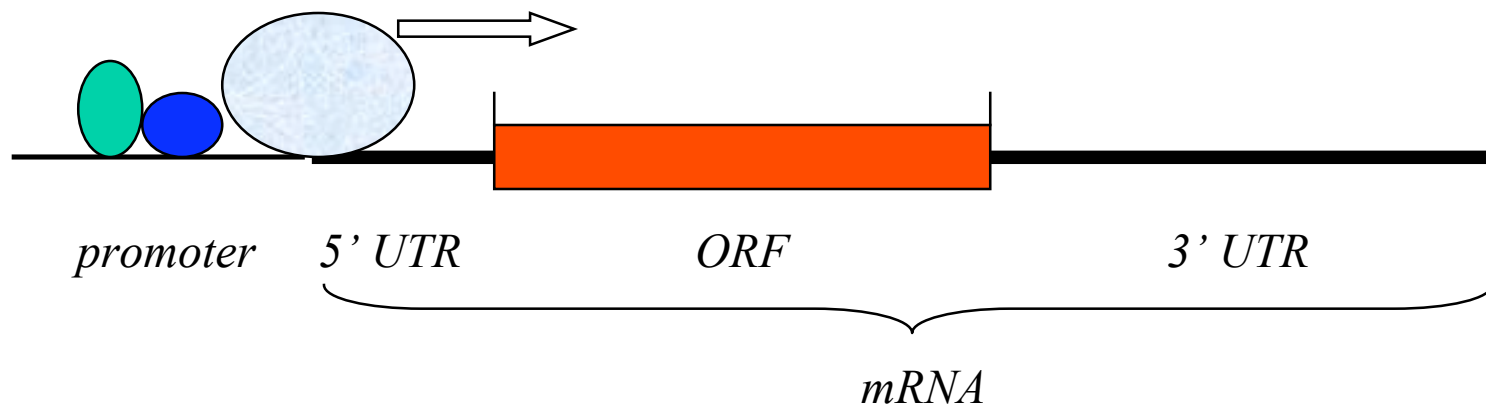
MNRLIELTGWIVLVVSVILLGVASHIDNYQPPEQSASVQHK



Gene's characteristics

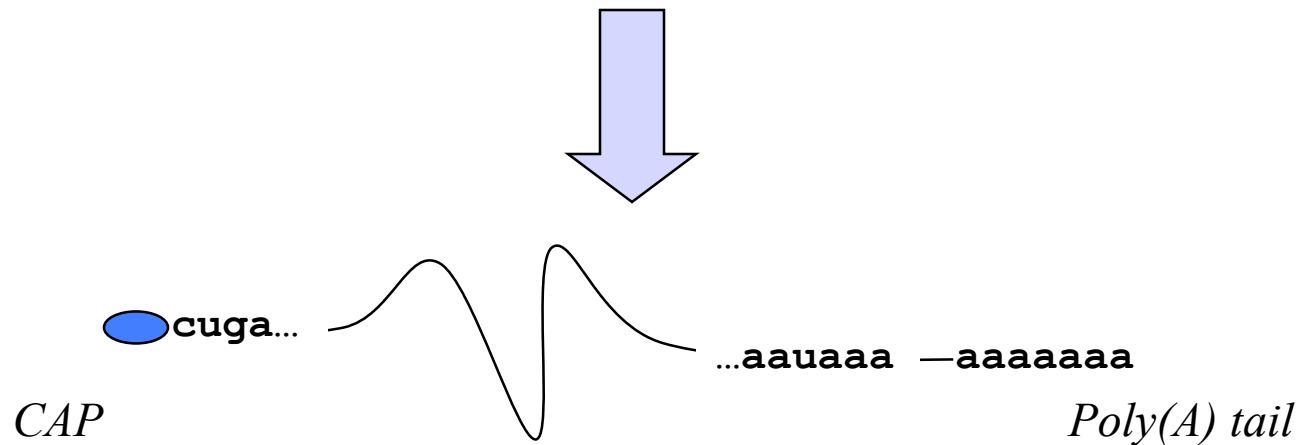
```

aatagcgaat tttccaacga caaaagctaa atatcgcaaa aacctcagta aaaatcttgc 60
tggagctatt attgctaagt aacatttacc ccctgaagtt aatggatcaa tcaagagaga 120
tgtgggctgt aATGaatcgt cttattgaat taacagggtt gatcgttctt gtcgtttcag 180
tcattcttct tggcgtggcg agtcacattg acaactatca gccacctgaa cagagtgctt 240
cggtacaaca caagTAAgct ctgcacttgt ggagcgacat gctgcccgtc cgggtgcatg 300
ttttcacttg tcggatatta aaccaggaat ttattatctt gttcgatgtt gtaataaa 358
  
```

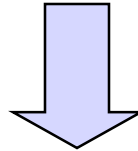
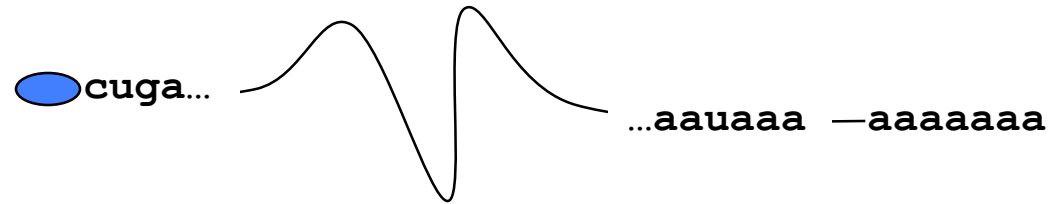


Transcription

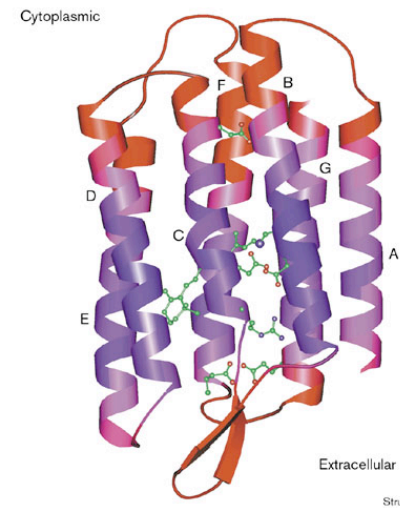
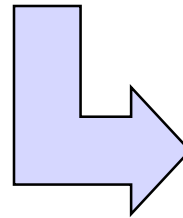
cugaaguu aauggaucaa ucaagagaga 120
ugugggcugu a**AUG**aaucgu cuuauugaau uaacagguug gaucguucuu gucguuucag 180
ucauucuucu uggcguggcg agucacauug acaacuauca gccaccugaa cagagugcuu 240
cgguacaaca caagUAAgcu cugcacuugu ggagcgacau gcugcccguc cgggugcaug 300
uuuucacuug ucggauauua aaccaggaau uuauuauuu guucgauguu guaauaaa 358



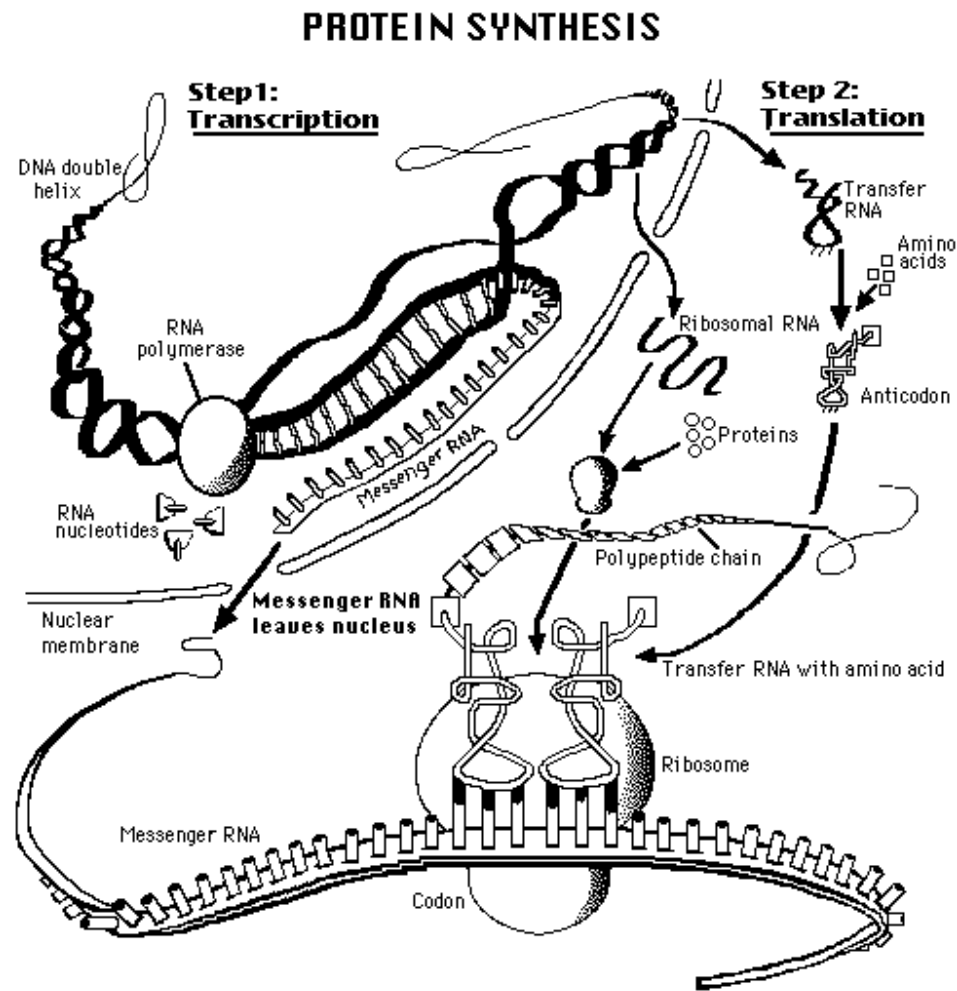
Translation



MNRLIELTGWIVLVVSVILLGVASHIDNYQPPEQSASVQHK



Protein coding genes (cntd)



Source:
<http://www.emc.maricopa.edu/faculty/farabee/BIOBK/BioBookPR>
OTSYN.html



BBSI 2006 26-MAY-2006

© 2006 P. Benos

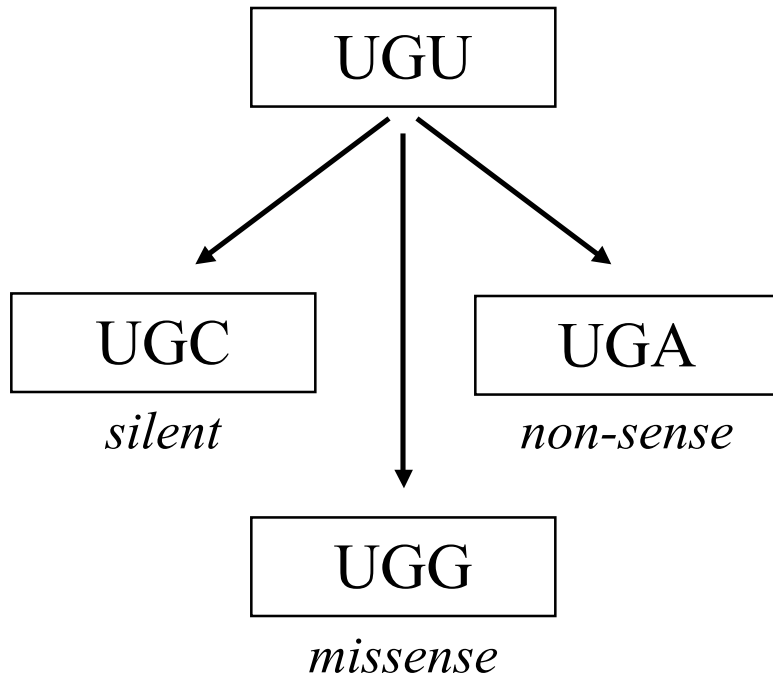
Alterations of the DNA

Base substitutions:

- silent: no a.a. replacement
- missense: a.a. replacement
- non-sense: a.a. → stop codon replacement



Alterations of the DNA (cntd)



		Second letter				
		U	C	A	G	
First letter	U	UUU Phenyl-alanine UUC UUA Leucine UUG	UCU Serine UCC UCA UCG	UAU Tyrosine UAC UAA Stop codon UAG Stop codon	UGU Cysteine UGC UGA Stop codon UGG Tryptophan	U C A G
	C	CUU Leucine CUC CUA CUG	CCU Proline CCC CCA CCG	CAU Histidine CAC CAA Glutamine CAG	CGU Arginine CGC CGA CGG	U C A G
	A	AUU Isoleucine AUC AUA AUG Methionine; initiation codon	ACU Threonine ACC ACA ACG	AAU Asparagine AAC AAA Lysine AAG	AGU Serine AGC AGA Arginine AGG	U C A G
	G	GUU Valine GUC GUA GUG	GCU Alanine GCC GCA GCG	GAU Aspartic acid GAC GAA Glutamic acid GAG	GGU Glycine GGC GGA GGG	U C A G

Source: <http://www.emc.maricopa.edu/faculty/farabee/BIOBK/BioBookPROTSYn.html>



Molecular evolution

Two species will acquire mutations proportionally to their divergence time. However:

- all proteins do not change in the same pace
- a given protein does not necessarily change in the same pace throughout time
- different parts of the same protein change at different paces



Molecular evolution (cntd)

Human (C11A_HUMAN; P05108) vs. Pig (C11A_PIG; P10612)

```
Query: 1   MLAKGLPPRSVLVKG YQTFLSAPREGLRRLRVPTGEGAGISTRSPRPFNEIPSPGDNGWL 60
          MLA+GL RSVLVKG O FLSAPRE G RV TGEGA IST++PRPF+EIPSPGDNGW+
Sbjct: 1   MLARGLALRSVLVKG CQPFLSAPRECPGHPRVGTGEGACISTKTPRPFSEIPSPGDNGWI 60

Query: 61  NLYHFWRETGTHKVLHHVQNFQKYGPIYREKLG NVESVYVIDPEDVALLFKSEGPNPER 120
          NLY FW+E GT K+H HHVQNFQKYGPIYREKLG N+ESVY+IDPEDVALLFK EGPNER
Sbjct: 61  NLYRFWKEKGTQKIHYHHVQNFQKYGPIYREKLG NLESVYIIDPEDVALLFKFEGPNPER 120

Query: 121  FLIPPWVAYHQYYQRP IGVLLKKSAAWKKDRVALNQEVMAPEATKNFLPLLDAVSRDFVS 180
          + IPPWVAYH O+Y O+P+GVLLKKS AWKKDR+ LN EVMAPEA KNF+PLLD VS+DFV
Sbjct: 121  YNIPPWVAYHQHYQKPVGVLLKKS GAWKKDR LVLNTEVMAPEAIKNFI PLLD TVSQDFVG 180

Query: 181  VLHRRIKKAGSGNYS GDISDDLFRFAFESITNVI FGERQGMLEEVVNPEAQRFIDAIYQM 240
          VLHRRIK+ GSG +SGDI +DLFRFAFESITNVI FGER GMLEE+V+PEAO+FIDA+YOM
Sbjct: 181  VLHRRIKQQGSGKFSGDIREDLFRFAFESITNVI FGERLGMLEEIVDPEAQKFIDAVYQM 240

Query: 241  FHTSVPMLNLPPDL FRLFRFTKTWKDHVAAWDVIFSKADIYTONFYWELRQKGSVHHDYRG 300
          FHTSVPMLNLPPDL FRLFRFTKTW+DHVAAWD IF+KA+ YTONFYW+LR+K ++Y G
Sbjct: 241  FHTSVPMLNLPPDL FRLFRFTKTWRDHVAAWDTIFNKAEKYTONFYWDLRRKRE-FNNYPG 299

Query: 301  MLYRLLGDSKMSFEDIKANVTEMLAGGVDTTSM TLQWHL YEMARNLKVQDMLRAEVLAAR 360
          +LYRLLG+ K+ ED+KANVTEMLAGGVDTTSM TLQWHL YEMAR+L VO+MLR EVL AR
Sbjct: 300  ILYRLLGNDKLLSE DVKANVTEMLAGGVDTTSM TLQWHL YEMARSLNVQEMLREEVLNAR 359

Query: 361  HQAQGD MATMLQLVPLLKASIKETLRLHPISVTL QRYLVNDLVL R DYMI PAKTLVQVAIY 420
          QAQGD + MLQLVPLLKASIKETLRLHPISVTL QRYLVNDLVL R DYMI PAKTLVQVA+Y
Sbjct: 360  RQAQGDTSKMLQLVPLLKASIKETLRLHPISVTL QRYLVNDLVL R DYMI PAKTLVQVAVY 419

Query: 421  ALGREPTFFFDPENFD PTRWLSKDKNITYFRNLGFGWGVRQCLGRRIAELEMTIFLINML 480
          A+GR+P FF +P FDPTRWL K++++ +FRNLGFGWGVRQC+GRRIAELEMT+FLI++L
Sbjct: 420  AMGRDPAFFSNPGQFDPTRWLGKERDLIHFRNLGFGWGVRQC VGRRIAELEMTLFLIHIL 479

Query: 481  ENFRVEIQHLSDVGTTFNLILMPEKPI SFTFWPFNQEATQ 520
          ENF+VE+OH SDV T FNLILMP+KPI F PFNO+ O
Sbjct: 480  ENFKVELQHFSVDVTIFNLILMPDKPIFLVFRPFNQDPLQ 519
```



Molecular evolution (cntd)

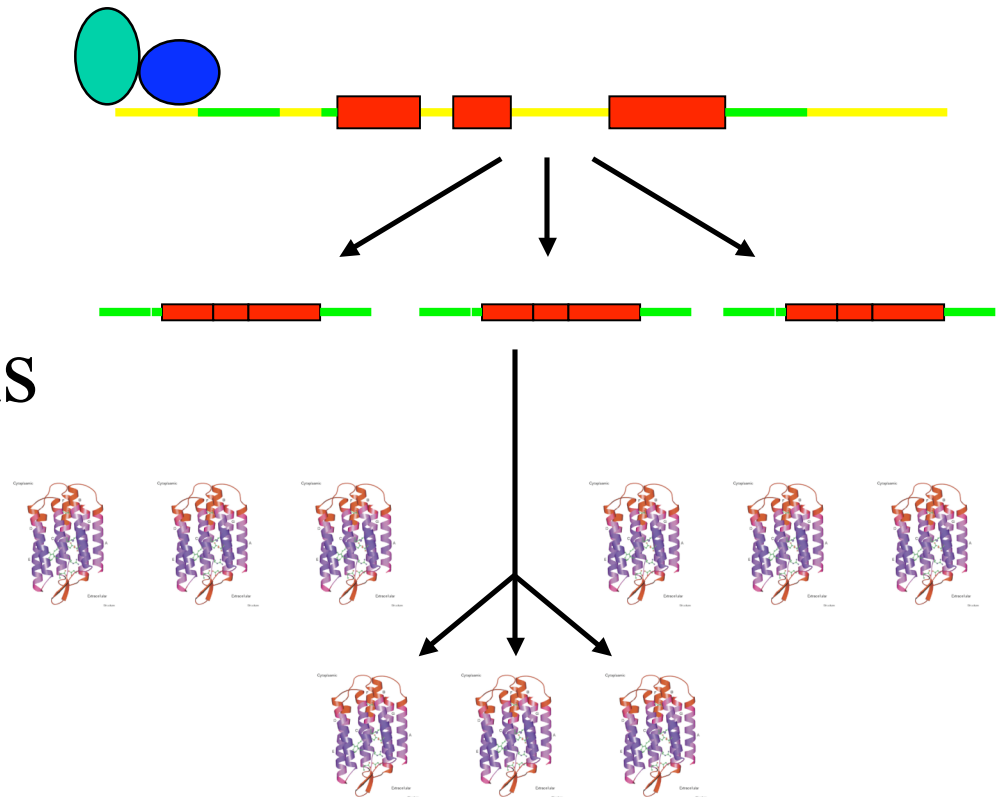
Human (C11A_HUMAN; P05108) vs. zebrafish (Cyp11a1; Q8JH93)

Query: 34 T GEGAGISTRSPRPFNEIPSPGDNGWLNLYHFWRETGTHKVLHVVQNFQKYGPIYREKL 93
Sbjct: 27 T G GRAPQNSTVQPFNKIPGRWRNSLLSVLAFTKMGGLRNVHRIMVHNFKTFGPIYREKV 86
Query: 94 G NVESVYVIDPEDVALLFKSEGPNPERFLIPPWVAYHQYYQRPIGVLLKKSAAWKKDRVA 153
Sbjct: 87 G IYDSVYIIKPEDGAILFKAEGHPNRINVDAWTAYRDYRNQKYGVLLKEGKAWKTDRMI 146
Query: 154 L NQEVMAPEATKNFLPLLDVSRDFVSVLHRRIKKAGSGNYSGDISDDLFRFAFESITNV 213
Sbjct: 147 L NKELLLPKLQGTFFVPLLDEVGQDFVARVVKQIERSGQKQWTTDLTHDLFRFSLESVSAV 206
Query: 214 I FGERQGMLEEVVNPEAQRFIDAIYQMFHTSVPMLNLPDDLFRFRFRTKTWKDHVAAWDVI 273
Sbjct: 207 I YGERLGLLLDNIDPEFQHFIDCVSVMFKTTSPLYLPPGLLRSIGSNIWKNHVEAWDGI 266
Query: 274 F SKADIYTONFYWELRQKGSVHHDYRGMLYRLLGDSKMSFEDIKANVTEMLAGGVDTTSM 333
Sbjct: 267 F NQADRCIQNIFKQWKENPEGNGKYPGVLAILLMQDKLSIEDIKASVTELMAGGVDSVTF 326
Query: 334 T LQWHLYEMARNLKVQDMLRAEVLAARHQAGDMATMLQLVPLLKASIKETLRLHPISVT 393
Sbjct: 327 T LLWTLYELARQPDLQDELRAEISAARIAFKGDMVQMVKMIPLLKAALKETLRLHPVAMS 386
Query: 394 L QRYLVNDLVLRDYMIPAKTLVQVAIYALGREPTFFFDPENFDPTRWLSKDKNITYFRNL 453
Sbjct: 387 L PRYITEDTVIQNYHIPAGTLVQLGVYAMGRDHQFFPKPEQYCPSRWISSNRQ--YFKSL 444
Query: 454 G FGGWVRQCLGRRIAEMTIFLINMLENFRVEIQHLSDVGTTFNLILMPEKPIISFTFWP 513
Sbjct: 445 G FGGFGPRQCLGRRIAETEMQIFLIHMLENFRIEKQKQIEVRSKFELLMLMPEKPIILTIP 504
Query: 514 F N 515
Sbjct: 505 L N 506

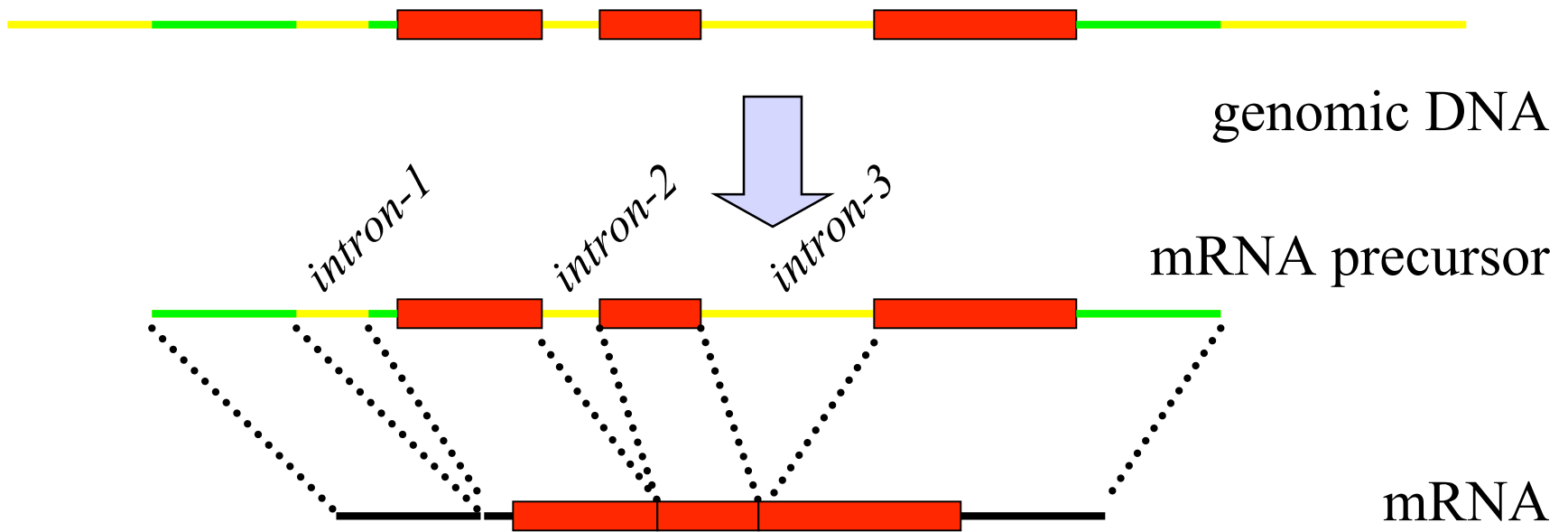


Gene expression regulation

- promoter region
- expression rates
- degradation
- post modifications

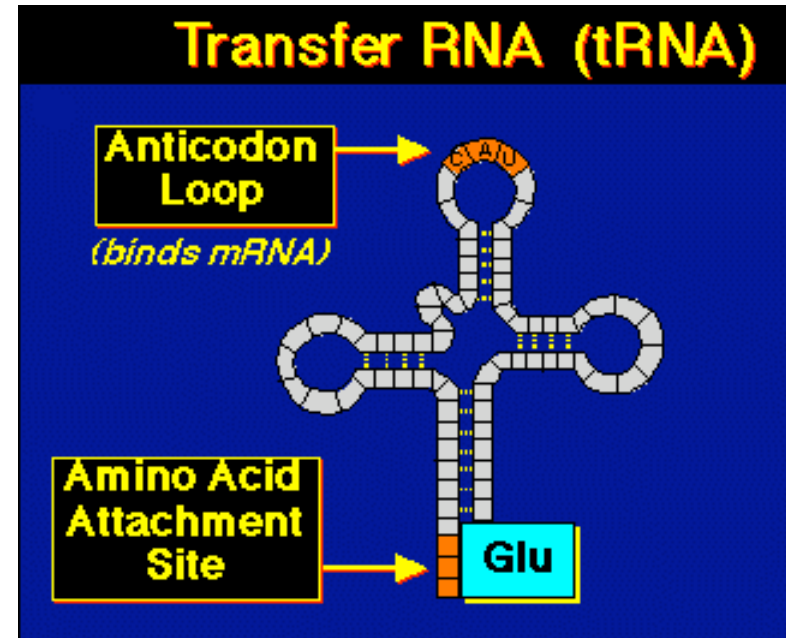


Splicing



Non-coding genes

- tRNA
- ribosomal RNA
- snoRNA
- microRNA
- etc



Source: <http://www.emc.maricopa.edu/faculty/farabee/BIOBK/BioBookPROTSYn.html>



Elements of Probability Theory (with examples)



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Outline

- Conditional Probabilities
- Markov Chains
- Hidden Markov Models
- Information measures



Probabilities

Definition:

$$P(x) = \frac{\# \text{ favourable outcomes } (x)}{\text{total \# possible outcomes}}$$

Conditional probabilities:

$$P(x|A) = \frac{\# \text{ favourable outcomes } (x) \text{ given } A}{\text{total \# possible outcomes given } A}$$



Conditional Probabilities

- *Joint probability:* $P(X, Y) = P(X|Y) P(Y)$
- If $P(X|Y) = P(X)$ then X, Y *independent*

$$P(X, Y) = P(X) P(Y)$$

- *Marginal probability:*

$$P(X) = \sum_Y P(X, Y) = \sum_Y P(X|Y)P(Y)$$



Conditional Probabilities (cntd)

- *Bayes' theorem*

$$P(X | Y) = \frac{P(Y | X)P(X)}{P(Y)} = \frac{P(Y | X)P(X)}{\sum_x P(Y | X)P(X)}$$

- *Posterior probabilities are the compromise between data and prior information.*



Bayes: Application-1

- Problem (from Durbin *et al.*, 1998):

A rare genetic disease is discovered with population frequency one in 1 million. An extremely good genetic test is 100% sensitive (always correct if you have the disease) and 99.99% specific (false positive rate 0.01%). Will you be willing to take such a test?

- Hint: What is the probability that you have the disease, if the test is positive?



Bayes: Application-1 (cntd)

- Answer:

$$\begin{aligned} P(D | +) &= P(+ | D) P(D) / P(+) = \\ &= 1.0 * 10^{-6} / [1.0 * 10^{-6} + 10^{-4} * (1 - 10^{-6})] = \\ &= 0.0099 \end{aligned}$$



Application of Bayes-2

- Problem:

Given a set of transmembrane proteins with specified membrane domains of length L (*training set*), can you develop a probabilistic model that predicts which parts of a new transmembrane protein are likely to be membrane domains?



Application of Bayes-2 (cntd)

- Solution:
 - Suppose that we suspect that the amino acid frequencies differ between membrane and non-membrane regions.
 - Using the *training* set, calculate the probabilities, $P(a_i|D)$, that each amino acid a_i is part of a membrane domain (D). Also, using the non-membrane parts, calculate the corresponding probabilities, $P(a_i|\text{not } D)$.



Application of Bayes-2 (cntd)

- Solution (*cntd*):
 - Divide the new protein into segments.
 - Using Bayes theorem, calculate the posterior probability of each segment being a membrane domain using the $P(a_i|D)$.

$$P(X | M); M := \arg \max_M \frac{P(M | D)P(D)}{\sum_d P(M | d)P(d)}$$



Markov chains

- What is a Markov chain?
- Markov chain of order n is a stochastic process of a series of outcomes, in which the probability of outcome x depends on the state of the previous n outcomes.



Markov chains (cntd)

- Markov chain (of *first* order):

$$\begin{aligned} P(x) &= P(X_L, X_{L-1}, \dots, X_1) = \\ &= P(X_L | X_{L-1}, \dots, X_1) P(X_{L-1} | X_{L-2}, \dots, X_1) \dots P(X_1) = \\ &= P(X_L | X_{L-1}) P(X_{L-1} | X_{L-2}) \dots P(X_2 | X_1) P(X_1) = \\ &= P(X_1) \prod_{i=2}^L P(X_i | X_{i-1}) \end{aligned}$$

- *Transition probabilities*: $P(X_i | X_{i-1})$



Application of Markov chains

- Problem (from Durbin *et al.*): CpG islands

Given two sets of sequences from the human genome, one with CpG islands and one without, can you calculate a model that can predict the CpG islands?




Application of Markov chains (cntd)

- Solution:

+	A	C	G	T
A	0.180	0.274	0.426	0.120
C	0.171	0.368	0.274	0.188
G	0.161	0.339	0.375	0.125
T	0.079	0.355	0.384	0.182

-	A	C	G	T
A	0.300	0.205	0.285	0.210
C	0.322	0.298	0.078	0.302
G	0.248	0.246	0.298	0.208
T	0.177	0.239	0.292	0.292

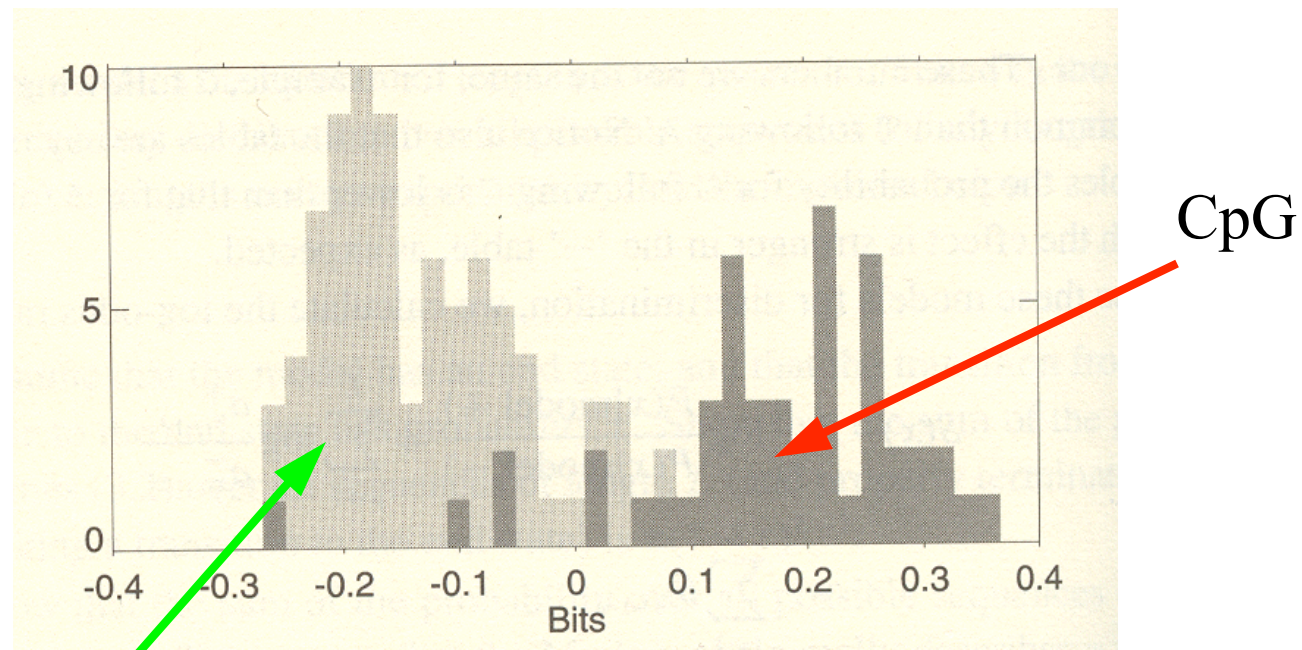


	A	C	G	T
A	-0.740	0.419	0.580	-0.803
C	-0.913	0.302	1.812	-0.685
G	-0.624	0.461	0.331	-0.730
T	-1.169	0.573	0.393	-0.679



Application of Markov chains (cntd)

- Histogram of scores (CpG islands):



other

Hidden Markov Models

- What is a Hidden Markov Model?
- A Markov process in which the probability of an outcome depends also in a (hidden) random variable (state).
- *Transition* probability: the probability of reaching a state given the previous state.
- *Emission* probability: the probability of an outcome given the state.

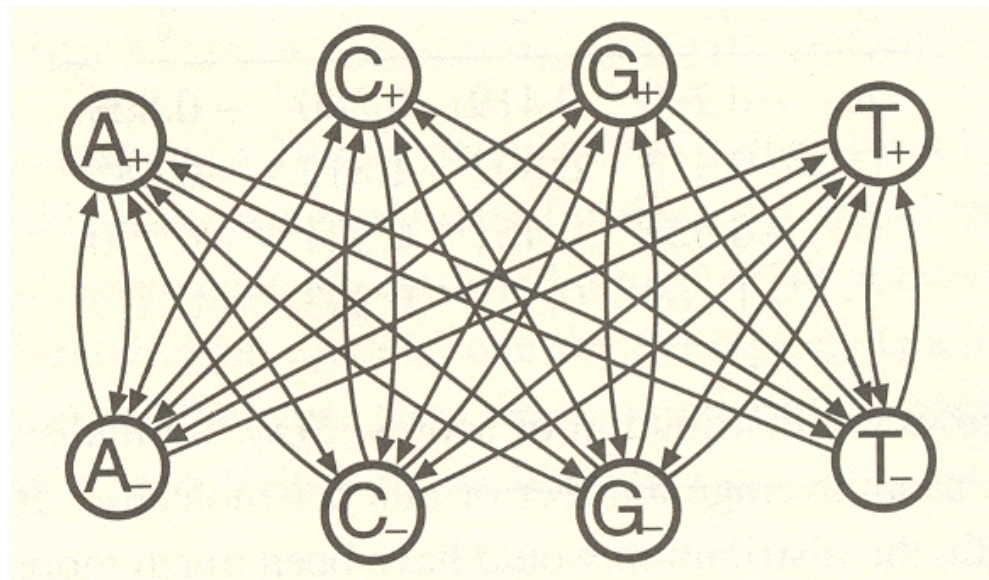


Hidden Markov Models (cntd)

- Graphical representation of the HMM:

CpG islands

(transition probabilities)

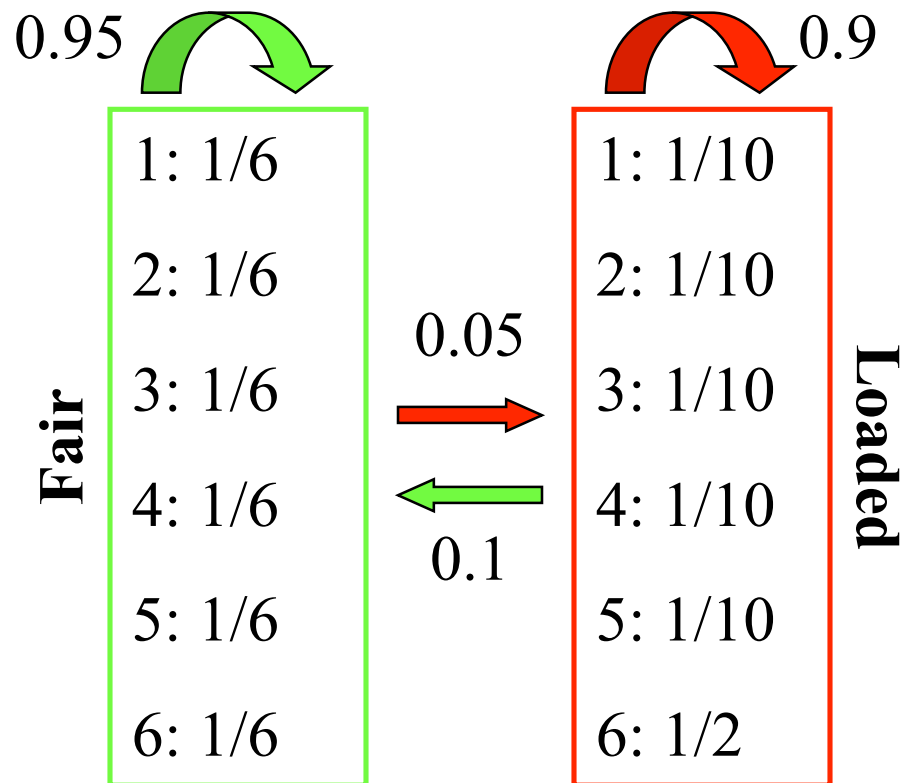


- Question: Where is the Markov process here?



Application of HMMs

- Problem (from Durbin *et al.*): dishonest casino



Application of HMMs (cntd)

- Problem (from Durbin *et al.*): dishonest casino

Given (1) the previous model and (2) a series of die rolls $(x_i, i=1, \dots, L)$, can we predict which of the rolls are coming from the fair and which from the loaded die?

- Question: What is “hidden” here?



Application of HMMs (cntd)

- Answer: YES
 - Viterbi algorithm (best path)
 - Forward-backward algorithm (probability of state k in outcome x_i)



HMMs: Viterbi algorithm

- Viterbi predictions: 300 rolls of die

```
Rolls 315116246446644245311321631164152133625144543631656626566666
Die   FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFL
Viterbi FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFL
```

```
Rolls 651166453132651245636664631636663162326455236266666625151631
Die   LLLLLLFFFFFFFFFFFFFFFFLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFL
Viterbi LLLLLLFFFFFFFFFFFFFFFFLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFL
```

```
Rolls 222555441666566563564324364131513465146353411126414626253356
Die   FFFFFFFFFLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL
Viterbi FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFL
```

```
Rolls 366163666466232534413661661163252562462255265252266435353336
Die   LLLLLLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFL
Viterbi LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL
```

```
Rolls 233121625364414432335163243633665562466662632666612355245242
Die   FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFLFFFFFFFFFFFFFFFFL
Viterbi FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFLFFFFFFFFFFFFFFFFL
```



HMMs in biology

- General comments:
 - Usually the structure of the model is unknown
 - The transition and emission probabilities are calculated based on trusted training set(s) and the postulated model
 - Predictions are based on the Viterbi or the forward-backward algorithm, depending on the question asked



Information measures

- Definitions:

- Entropy:

$$H(P) = \mathbf{E}(-\log P) = -\sum_{i=1}^n p_i \log p_i$$

- Relative Entropy:

$$H(P, Q) = \sum_{i=1}^n p_i \log \frac{p_i}{q_i}$$

- Mutual Information:

$$M(X, Y) = \sum_{i,j} P(x_i, y_j) \log \frac{P(x_i, y_j)}{P(x_i)P(y_j)}$$

