

# CISC 889 Bioinformatics (Spring 2004)

## Gene Finding

CISC889, S04, Lec10, Liao

1

### Gene prediction strategies

- Content-based
  - Codon usage
  - Periodicity of repeats
  - Compositional complexity
- Site-based
  - Binding sites for transcription factors
  - polyA tracts,
  - Donor and acceptor splice sites
  - Start and stop codons
- Comparative
  - BLAST

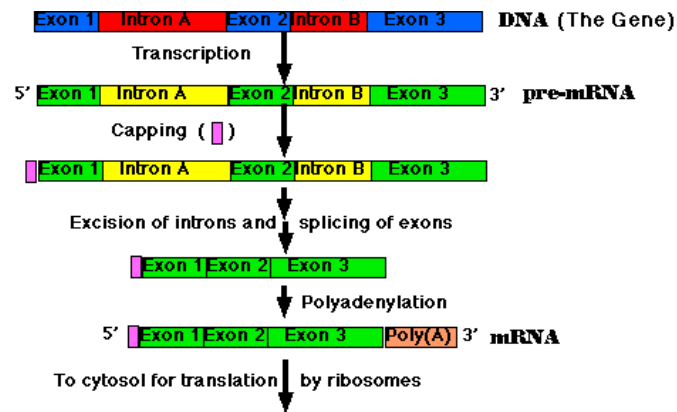
CISC889, S04, Lec10, Liao

2

## Gene expression

Transcription: DNA  $\rightarrow$  mRNA

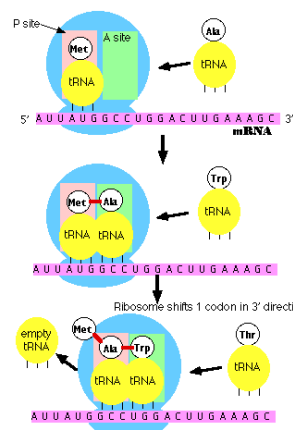
Translation: mRNA  $\rightarrow$  Protein



Kimball's Biology page

CISC889, S04, Lec10, Liao

3



Kimball's Biology page

CISC889, S04, Lec10, Liao

4

## Universal Genetic Code

5'	2 <sup>nd</sup> Position				3'
	T	C	A	G	
T	Phe	Ser	Tyr	Cys	T
	Phe	Ser	Tyr	Cys	C
	Leu	Ser	Stop	Stop	A
	Leu	Ser	Stop	Trp	G
C	Leu	Pro	His	Arg	T
	Leu	Pro	His	Arg	C
	Leu	Pro	Gln	Arg	A
	Leu	Pro	Gln	Arg	G
A	Ile	Thr	Asn	Ser	T
	Ile	Thr	Asn	Ser	C
	Ile	Thr	Lys	Arg	A
	Met	Thr	Lys	Arg	G
G	Val	Ala	Asp	Gly	T
	Val	Ala	Glu	Gly	A
	Val	Ala	Glu	Gly	G

ACCUUAGCGUA } Reading frame 1  
Thr Leu Ala

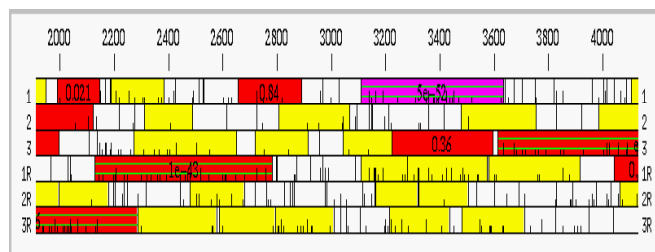
ACCUUAGCGUA } Reading frame 2  
Pro Stop Arg

ACCUUAGCGUA } Reading frame 3  
Leu Ser Val

CISC889, S04, Lec10, Liao

5

## Open Reading Frame (ORF)



CISC889, S04, Lec10, Liao

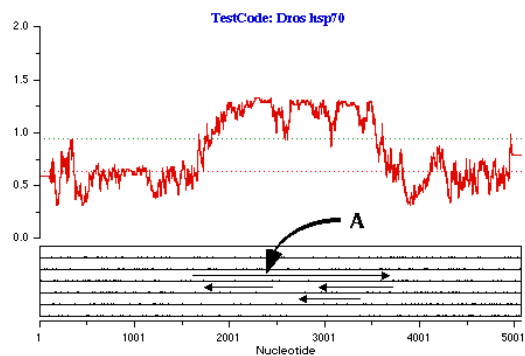
6

- Prokaryotic
  - Most regions of DNA are coding regions
  - No introns
- Eukaryotic
  - Introns and Exons

CISC889, S04, Lec10, Liao

7

- Fickett's rule (1982)
  - In ORFs, every third base tends to be the same one more often than by chance alone.
    - Regardless species,
    - No knowledge of codon preference is required.



CISC889, S04, Lec10, Liao

8

- Codon Usage Index

- There are 64 codons but 20 amino acids to code, therefore some AAs are coded by multiple codons.
  - For example, 6 codons for Leu, and 4 for Ala, but only one for Try.
  - For random DNA sequences, the frequency of having these three AAs would be 6/4/1 for Lue:Ala:Trp. In real protein sequences, ratio was found to be 6.9/6.5/1, which implies coding DNA sequence is not random;
  - some codons are preferred (depending on species.)

CISC889, S04, Lec10, Liao

9

Codon usage database: [www.kazusa.or.jp/codon](http://www.kazusa.or.jp/codon)

Address [http://www.kazusa.or.jp/codon/cgi-bin/showcodon.cgi?species=Escherichia+coli+K12+\[gbbct\]](http://www.kazusa.or.jp/codon/cgi-bin/showcodon.cgi?species=Escherichia+coli+K12+[gbbct])

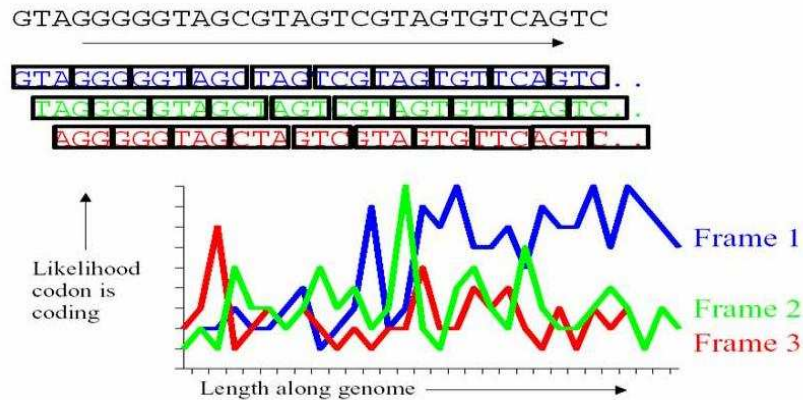
<i>Escherichia coli</i> K12 [gbbct]: 5089 CDS's (1608122 codons)			
fields: [triplet] [frequency: per thousand] ([number])			
UUU 22.4 ( 35982)	UCU 8.5 ( 13687)	UAU 16.3 ( 26266)	UGU 5.2 ( 8340)
UUC 16.6 ( 26678)	UCC 8.6 ( 13849)	UAC 12.3 ( 19728)	UGC 6.4 ( 10347)
UUA 13.9 ( 22376)	UCA 7.2 ( 11511)	UAA 2.0 ( 3246)	UGA 0.9 ( 1468)
UUG 13.7 ( 22070)	UCG 8.9 ( 14379)	UAG 0.2 ( 378)	UGG 15.3 ( 24615)
CUU 11.0 ( 17754)	CCU 7.1 ( 11340)	CAU 12.9 ( 20728)	CGU 21.0 ( 33694)
CUC 11.0 ( 17723)	CCC 5.5 ( 8915)	CAC 9.7 ( 15595)	CGC 22.0 ( 35306)
CUA 3.9 ( 6212)	CCA 8.5 ( 13707)	CAA 15.4 ( 24835)	CGA 3.6 ( 5716)
CUG 52.7 ( 84673)	CCG 23.2 ( 37328)	CAG 28.8 ( 46319)	CGG 5.4 ( 8684)
AUU 30.4 ( 48818)	ACU 9.0 ( 14397)	AAU 17.7 ( 28465)	AGU 8.8 ( 14092)
AUC 25.0 ( 40176)	ACC 23.4 ( 37624)	AAC 21.7 ( 34912)	AGC 16.1 ( 25843)
AUA 4.3 ( 6962)	ACA 7.1 ( 11366)	AAA 33.6 ( 54097)	AGA 2.1 ( 3337)
AUG 27.7 ( 44614)	ACG 14.4 ( 23124)	AAG 10.2 ( 16401)	AGG 1.2 ( 1987)
GUU 18.4 ( 29569)	GCU 15.4 ( 24719)	GAU 32.2 ( 51852)	GGU 24.9 ( 40019)
GUC 15.2 ( 24477)	GCC 25.5 ( 40993)	GAC 19.0 ( 30627)	GGC 29.4 ( 47309)
GUA 10.9 ( 17508)	GCA 20.3 ( 32666)	GAA 39.5 ( 63517)	GGA 7.9 ( 12776)
GUG 26.2 ( 42212)	GCG 33.6 ( 53988)	GAG 17.7 ( 28522)	GGG 11.0 ( 17704)

Coding GC 51.80% 1st letter GC 58.89% 2nd letter GC 40.72% 3rd letter GC 55.79%

CISC889, S04, Lec10, Liao

10

## Using codon usage for prediction ORFs



Can identify regions where there is a high frequency of 'preferred' codons

CISC889, S04, Lec10, Liao

11

## ORFs as Markov chains

- Glimmer: Interpolated Markov models (IMM)
  - 1<sup>st</sup> order model:  $p(a|a), p(a|c), \dots, p(t|t)$ . Probability of having a amino acid given its previous neighbor.
  - 2<sup>nd</sup> order model:  $p(a|xx)$
  - Up to 8<sup>th</sup> order model (0<sup>th</sup> for random, 1<sup>st</sup> to 6<sup>th</sup> for 6 reading frames, why higher order? Why stop at 8<sup>th</sup> ?
    - E.g, 5<sup>th</sup> order model, need  $4^6$  conditional probabilities  $p(a|xxxxx)$ . In a genome of 1.8Mb, for each 6mer, we can observe about  $1.8\text{Mb}/4096$  samples. But the higher  $k$ , the less number of samples for  $k$ mers .
  - Interpolation (linear combination of models of different orders)
 
$$P(S|M) = \sum_{x=1}^n \text{IMM}_g(S_x)$$
 where  $S_x$  is the oligomer ending at position  $x$  and  $n$  is the sequence length. The interpolated Markov model score is
 
$$\text{IMM}_k(S_x) = \lambda_k(S_{x-1}) P_k(S_x) + [1 - \lambda_k(S_{x-1})] \text{IMM}_{k-1}(S_x)$$
 where  $\lambda_k(S_{x-1})$  is the numerical weight associated with the  $k$ mer ending at position  $x-1$ , and  $P_k(S_x)$  is the probability of having  $S_x$ , predicted by  $k$ -th order model.

CISC889, S04, Lec10, Liao

12

- Glimmer (cont'd)

Results for *H. Influenzae*:

model	Found	Missed	New
Glimmer	1680	37	209
5 <sup>th</sup> order	1574	143	104

## Self-identification (Audic and Claverie '98)

- Probability of sequence  $W$  of length  $L$  is generated by a  $k$ -th order Markov chain

$$P(W|M) = P(S_0) \prod_{i=k}^{L-1} P(n_i | S_{i-k}) \quad \text{eq(1)}$$

where  $S_i$  is a  $k$ mer starting at position  $i$  in  $W$ . The model contains all the probabilities for any possible kmers to be followed by one of nucleotides A, C, G, or T.  $k = 5$  is used.

- Which model is better?

$$P(M_j | W) = P(W|M_j)P(M_j) / \sum_{r=1 \text{ to } N} P(W|M_r)P(M_r) \quad \text{eq(2)}$$

where *a priori* probability  $P(M_j)$  is assumed to be equally probable for  $N$  models, i.e., is  $1/N$ .

If we have three models corresponding to coding, reverse coding, and non-coding, then the posterior probability tells what sequence  $W$  is more likely to be.

- Model building (no training data is required)

- How to build the three models?

- If we have regions labeled as coding, reverse coding, and non-coding, then we can count the frequencies to train the transition matrices.

- Self consistent

- Randomly cut into non-overlapping pieces of  $w$  bases long, and assign them randomly into three distinct subsets, and build three Markov models  $M_1$ ,  $M_2$ , and  $M_3$  respectively.
    - Scan genomic sequence using size  $w$  window. For each window segment, determine its class using eq(2). Slide the window by 5 bases, and repeat the process.
    - If a region is covered by  $n$  (for  $5n \geq w$ ) successive windows of  $M_j$  type, it is qualified to be assigned into the  $j$  data set.
    - After finishing one scan of the genomic sequence, the 3 subsets are updated, and new Markov models are built for each of the 3 subsets.
    - Repeat the whole process until convergence is reached.

CISC889, S04, Lec10, Liao

15

- Results ( $k = 5$ ,  $w = 100$ )

- Convergence is reached after 50 iterations

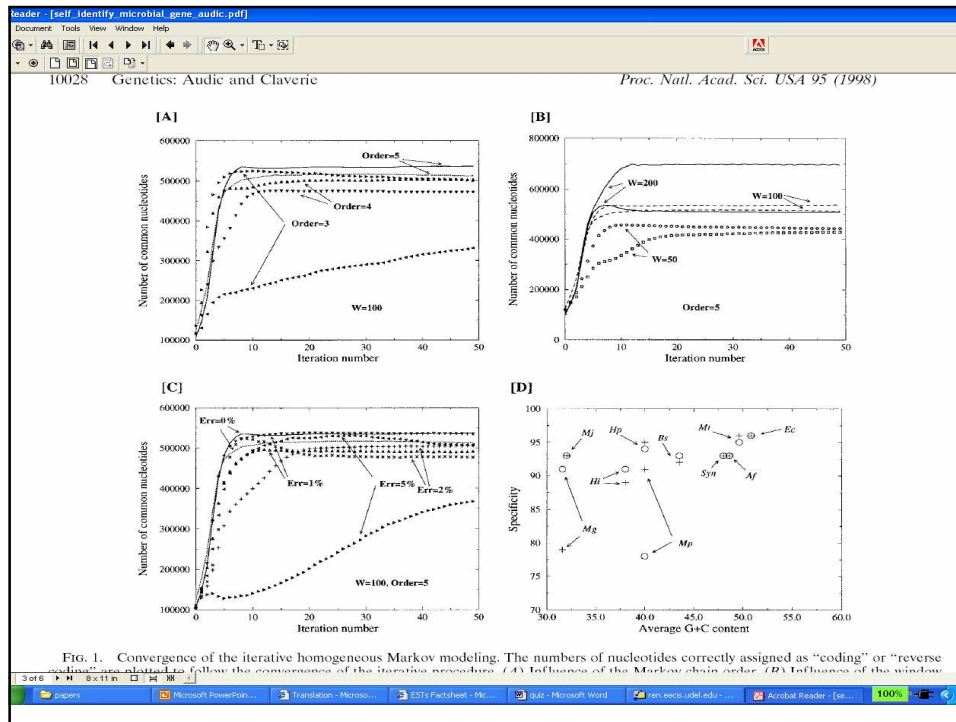
- *H. pylori*

- Correct rate: 95%, 94%, 93.8% for coding, reverse coding and non-coding respectively.

CISC889, S04, Lec10, Liao

16





## More markov based tools

- **HMMgene**  
<http://www.cbs.dtu.dk/services/HMMgene/>
  - Hidden Markov model
    - whole genes
    - partial genes
    - Cosmids or even longer sequences.
- **GeneScan**
- **Genemark**
  - <http://opal.biology.gatech.edu/GeneMark/>

- **blastx** compares a nucleotide query sequence translated in all reading frames against a protein sequence database
- **tblastn** compares a protein query sequence against a nucleotide sequence database dynamically translated in all reading frames.
- **tblastx** compares the six-frame translations of a nucleotide query sequence against the six-frame translations of a nucleotide sequence database. Please note that tblastx program cannot be used with the nr database on the BLAST Web page.

CISC889, S04, Lec10, Liao

19

## Prediction Assessment

- (Baxevanis & Ouellete, 2<sup>nd</sup> ed, page 246)
  - TP (true positive), FP (false positive), TN (true negative), FN (false negative)
  - Actual positive, negative; Predicted positive, negative
  - $TP + TN + FP + FN = N$
  - Sensitivity =  $TP / (TP + FN)$
  - Specificity =  $TP / (TP + FP)$
  - Correlation Coefficient =  $TP \cdot TN - FP \cdot FN / \sqrt{(PP \cdot PN \cdot AP \cdot AN)}$
  - $CC = [-1, 1]$

CISC889, S04, Lec10, Liao

20

## **Strategies for Gene finding**

- Challenges remain:
  - partial genes, non-coding RNA genes, etc.
- MZEF best for single exons
- GENSCAN best for whole genes
- shall try one more method for cross reference
- Shall resort to comparative method, e.g., run BLAST against dbEST and/or protein databases.