



# Методы предсказания структуры генов эукариот

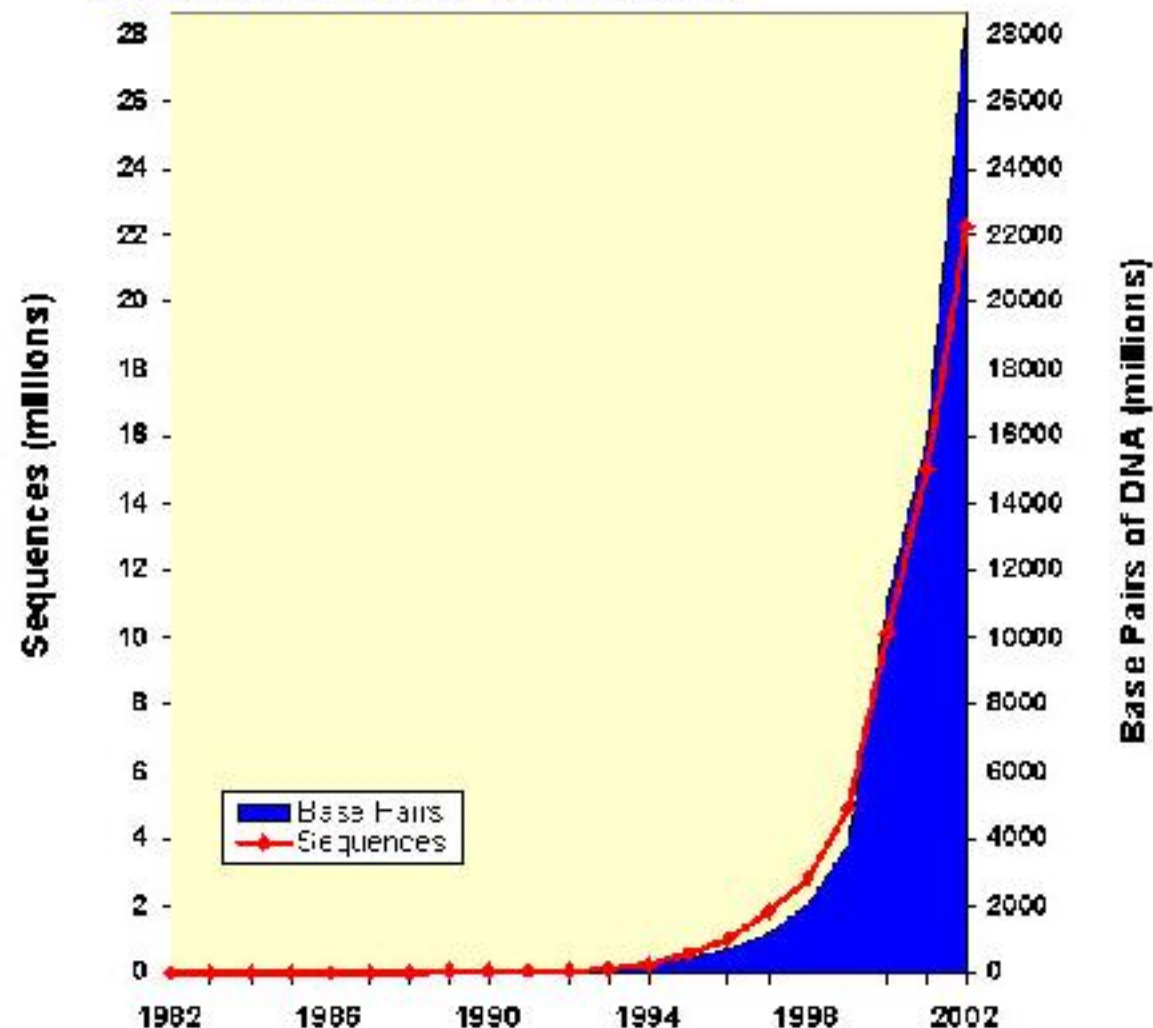
*Олег Владимирович Витневский*

Кафедра информативной биологии ФЕННГУ

# Международные проекты геномных исследований

Стремительно растут темпы исследований по секвенированию геномной ДНК (Benson et al., 2000; Wheeler et al., 2000). На сентябрь 2003 г. доступны 139 полных геномов прокариот, включая 16 видов археобактерий и 123 вида бактерий.

## Рост объема GenBank



**GenBank за 2002 год:**

**28,507,990,166 п.о.**

**22,318,883  
последовательностей**

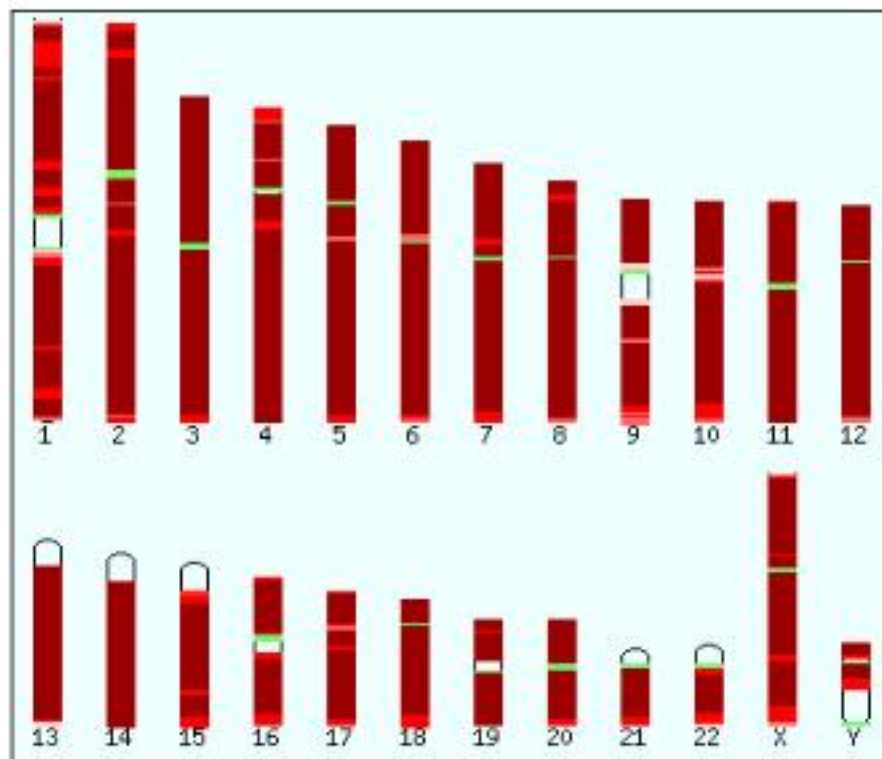
Полностью секвенированы эукариотические геномы мышевидного салата *Arabidopsis thaliana*, червя *Caenorhabditis elegans*, плодовой мушки *Drosophila melanogaster*, дрожжей *Saccharomyces cerevisiae* и *Schizosaccharomyces pombe*, некоторых внутриклеточных паразитических организмов (*Plasmodium falciparum*, *Encephalitozoon cuniculi*).



## Ensembl Stats: Chromosome status

The following images indicate the status of the Golden Paths

### NCBI34



### Key

- N50 length > 5mb
- N50 length > 1mb
- N50 length > 500kb
- N50 length > 100kb
- N50 length > 10kb
- More than 50% of bin
- Less than 50% of bin
- Clear- None of bin in golden

### Stats

**Freeze date:**

July 2003

**Estimated size:**

3 069.43 Mb

**Total mapped:**

2 843.41 Mb (92.64%)

**No. of supercontigs:**

350

**Super contig N50 length:**

29 104 799 bps

**In super contigs > 10Mb**

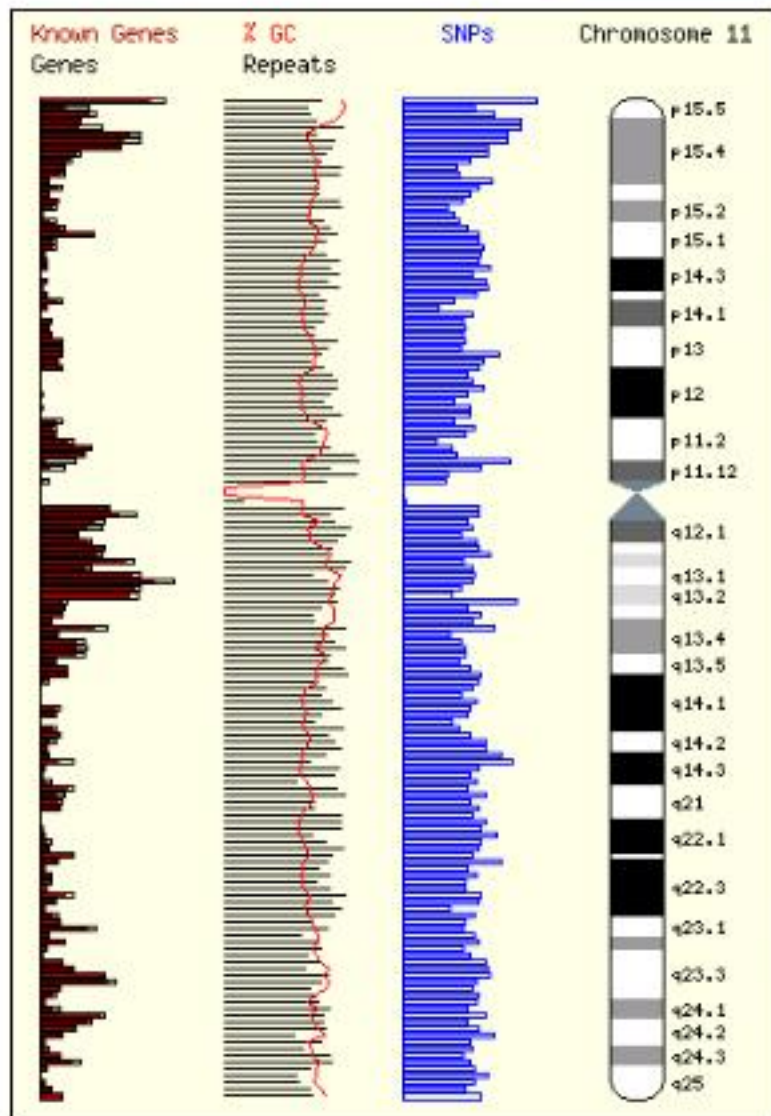
2 307.65 Mb (76 s'ctgs)

**In super contigs > 1Mb**

2 789.20 Mb (199 s'ctgs)

**In super contigs > 100Kb**

2 842.38 Mb (332 s'ctgs)



## Chromosome 11

Length: 134,482,954 bps  
Known Ensembl genes: 1,296  
Novel Ensembl genes: 227  
Snps: 245,108

## Change Chromosome

Chromosome:

## Jump to Contigview

Click anywhere on the chromosome ideogram or one of the feature distribution level view of features at that point. Alternatively, you can jump to contigview between any two features on this chromosome:

Between:   
and:

[Display contig-level view between any two features.](#)

## Synteny

View Human Chr 11 vs



Address [http://www.ensembl.org/Homo\\_sapiens/contigview?chr=11&highlight=&vc\\_start=3404185&vc\\_end=3504184&bottom=%7Cmarker%3Aoff](http://www.ensembl.org/Homo_sapiens/contigview?chr=11&highlight=&vc_start=3404185&vc_end=3504184&bottom=%7Cmarker%3Aoff)

# Ensembl Human ContigView



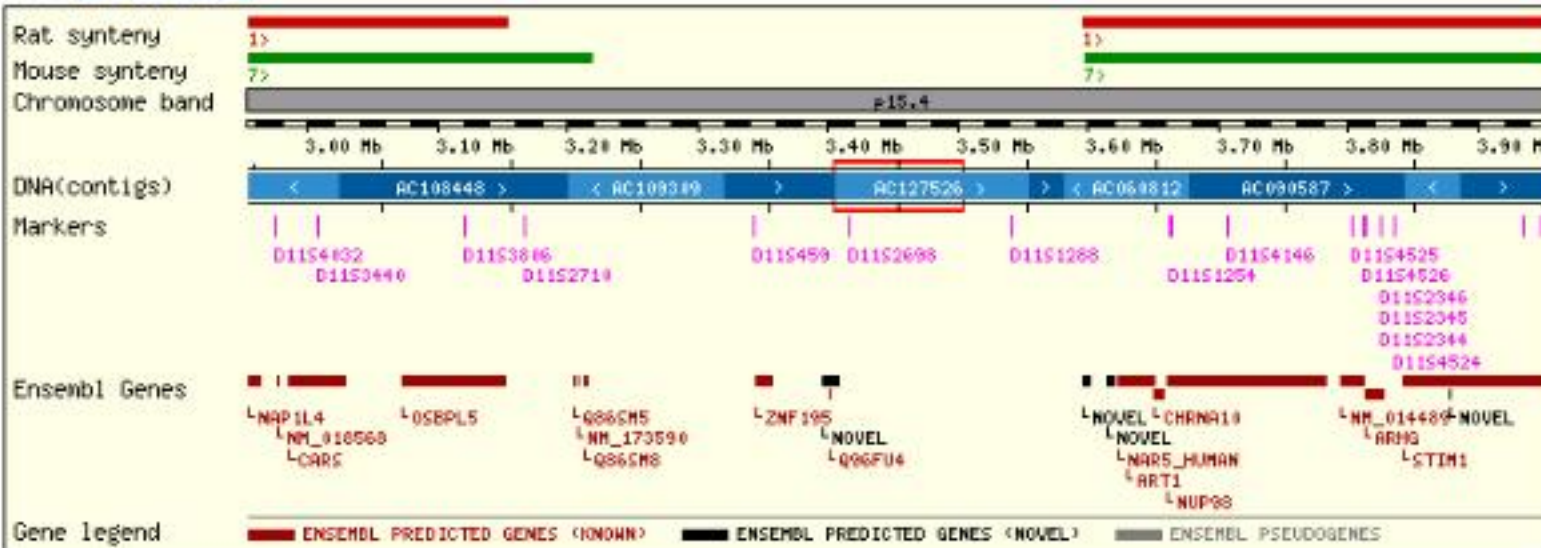
Home Human What's New TextSearch BlastSearch MapSearch Export Data Download Disease Browser Docs

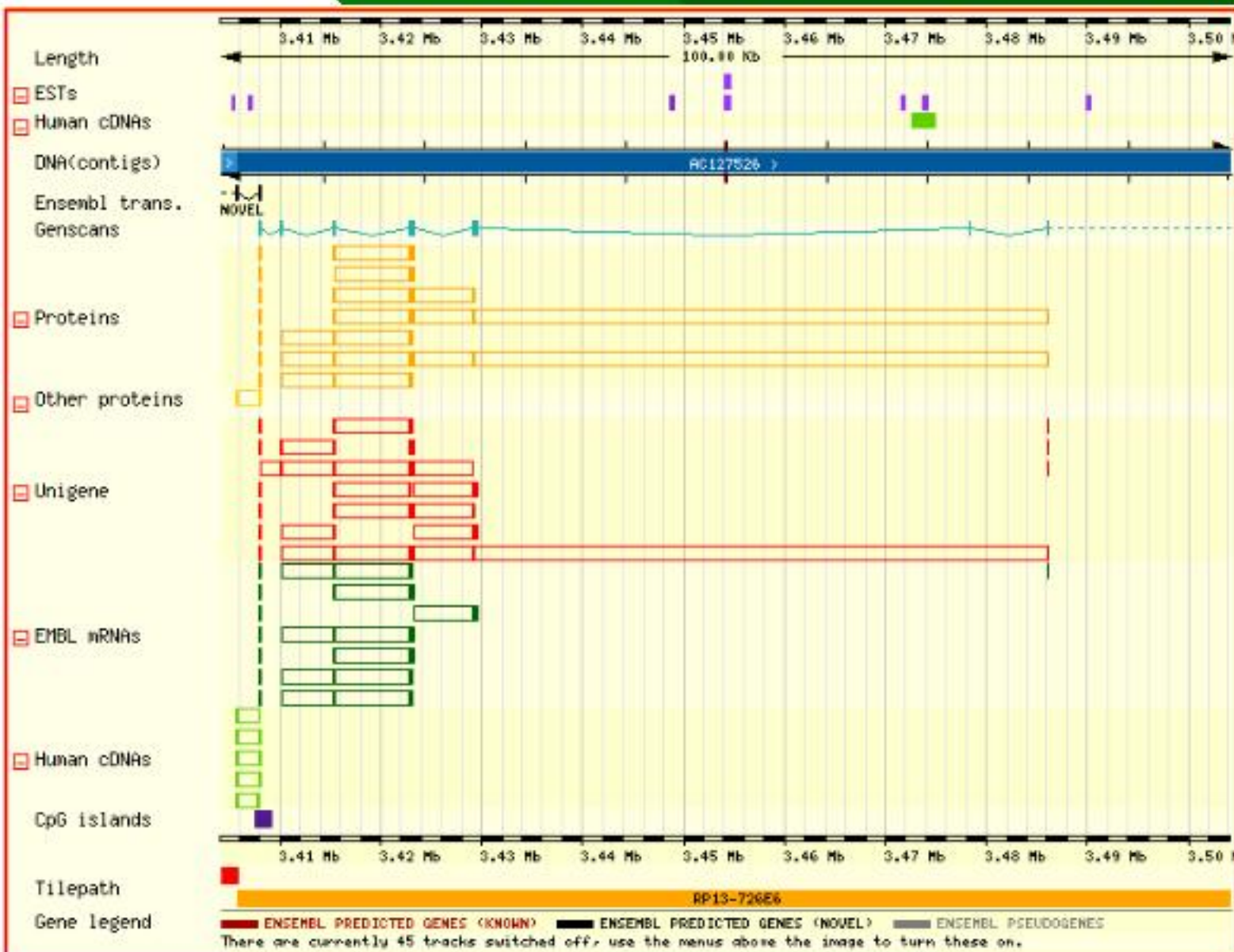
Find Sequence   [e.g. [AC067852](#), [AP003171](#)]

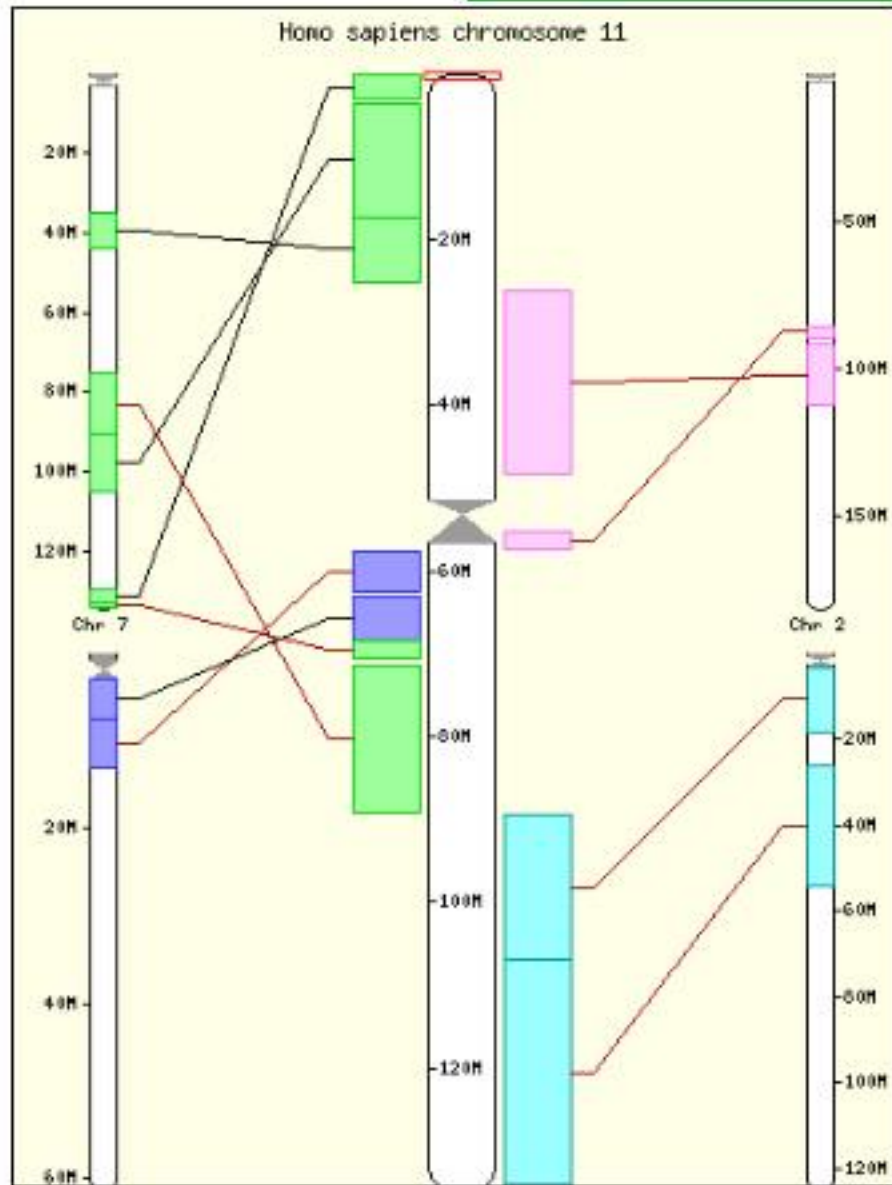
## Chromosome 11



## Overview







## Human Chromosome 11

Jump to chromosome

[Jump to mapview](#) for chromosome statistics.

### Homology Matches

*Homo\_sapiens* Genes      *Mus\_musculus* Homologues

**ENSG00000170061**

(149.64 Kb)

**ENSG00000177934**

(152.41 Kb)

**NM\_145651**

(183.08 Kb)

**ODF3**

(186.79 Kb)

**ENSG00000177951**

(193.48 Kb)

**BET1L**

(195.30 Kb)

**NM\_021932**

(198.85 Kb)

**SIRT3**

(206.14 Kb)

**PSMD13**

(227.05 Kb)

**PYA5\_HUMAN**

(268.57 Kb)

**NM\_025092**

(279.13 Kb)

**ENSG00000142056**

(288.80 Kb)

-> **ENSMUSG00000038801**

(chr 7 : 129.87 Mb)

-> **1700011004Rik**

(chr 7 : 129.87 Mb)

-> **Bet1l**

(chr 7 : 129.87 Mb)

-> **Al114950**

(chr 7 : 129.88 Mb)

-> **Sirt3**

(chr 7 : 129.88 Mb)

-> **Psm13**

(chr 7 : 129.90 Mb)

-> **Pypa5**

(chr 7 : 129.94 Mb)

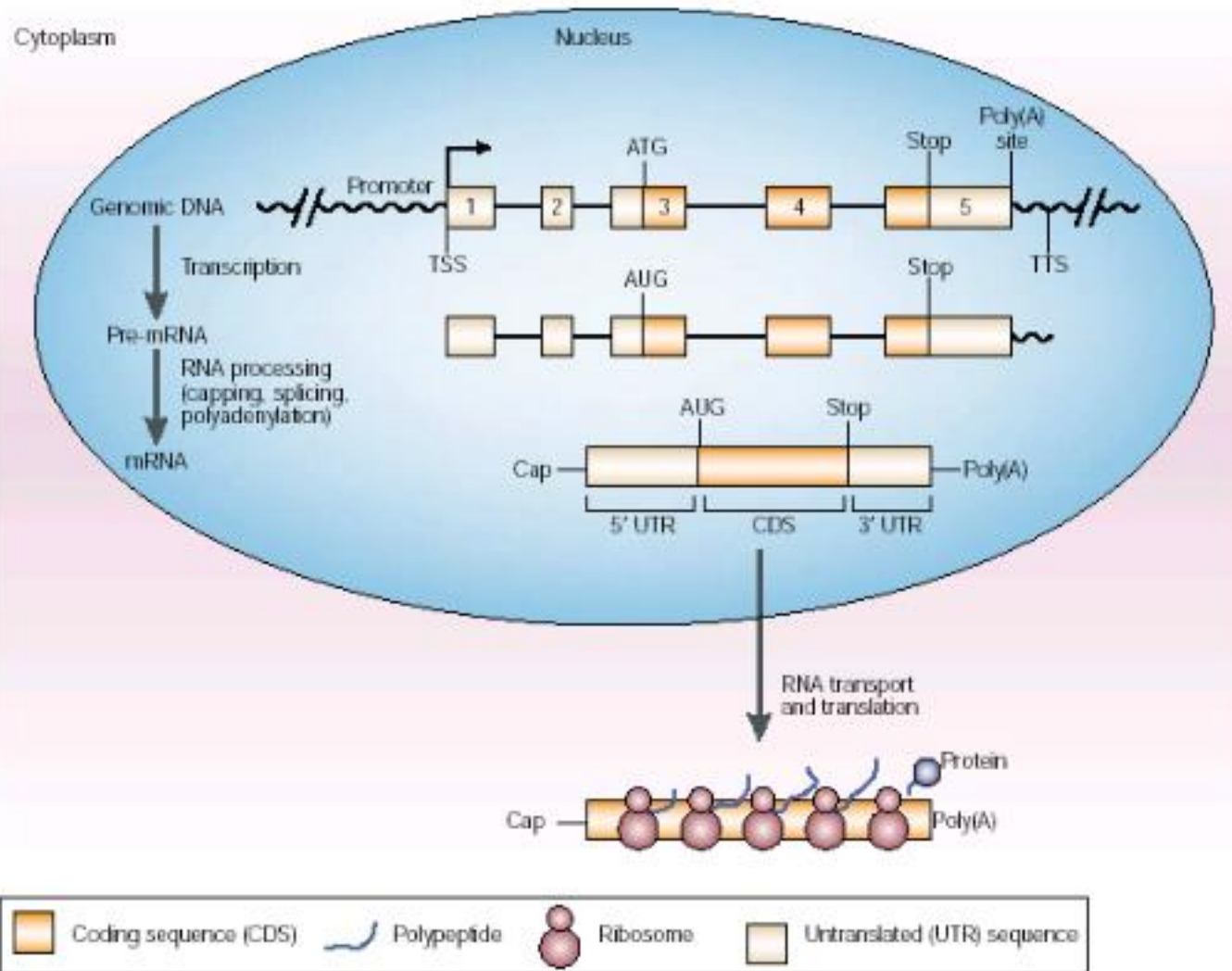
-> **NM\_145387**

(chr 7 : 129.96 Mb)

-> **Hrmp1-pending**

(chr 7 : 129.97 Mb)

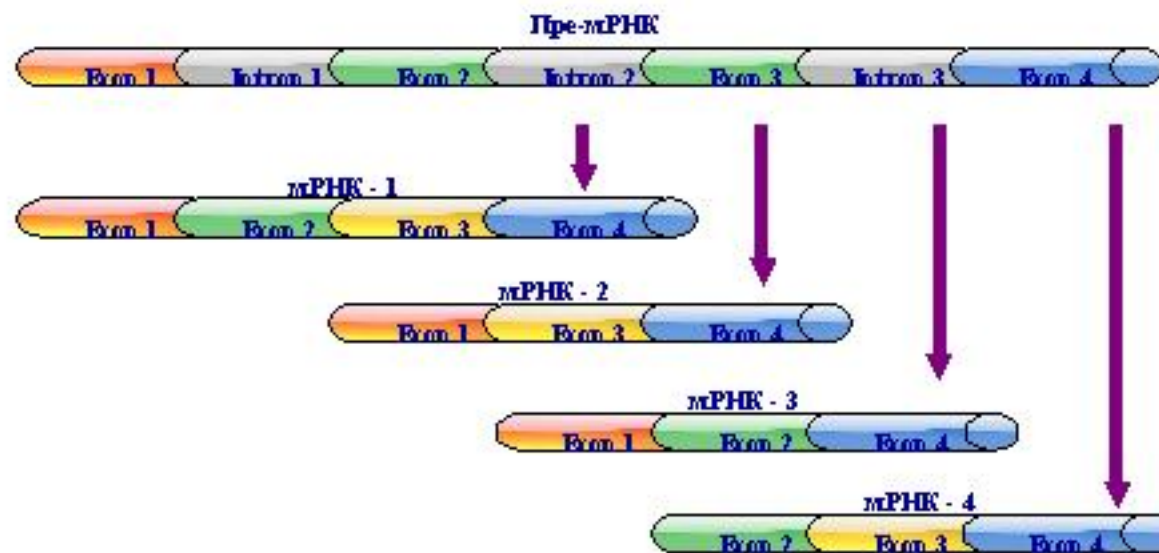
# Модель эукариотического гена



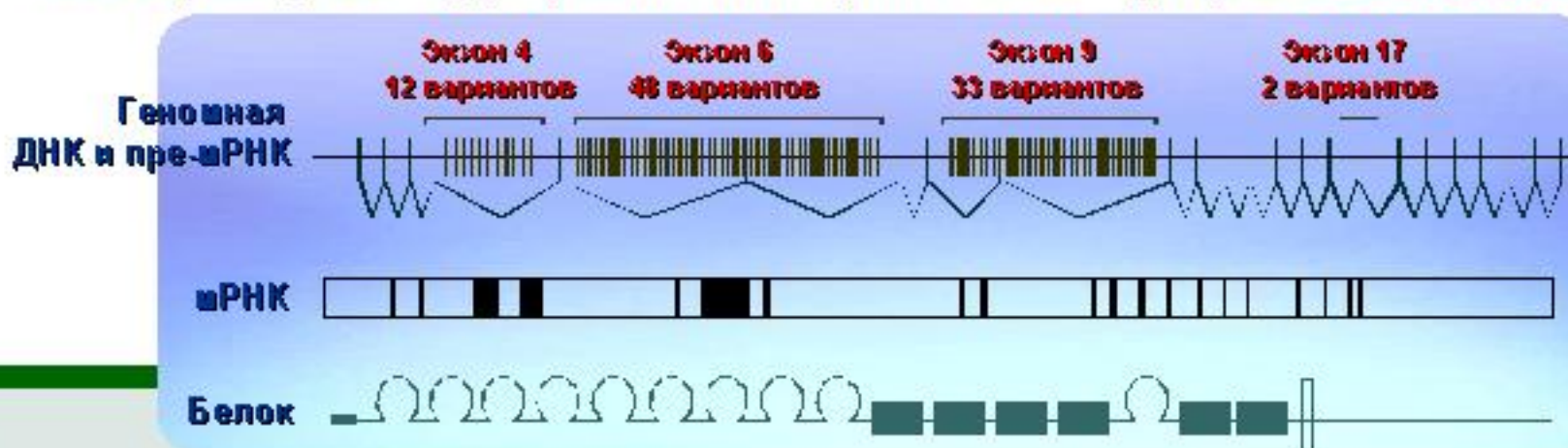




# Экзон-интронная структура и альтернативный сплайсинг обеспечивают огромную емкость кодирования генетической информации



Альтернативный сплайсинг обеспечивает продукцию огромного количества (более 30 000) вариантов белка DSCAM, участвующего в формировании тонкой нервной системы дрозофилы:  $N = 12 \times 48 \times 33 \times 2 \dots$





	Length of gene	Length of mRNA	Number of introns
--	----------------	----------------	-------------------

$\beta$ -globine	1,5	0,6	2
Insuline	1,7	0,4	2
Protéine kinase C	11	1,4	7
Albumine	25	2,1	14
Catalase	34	1,6	12
Récepteur des I.DI.	45	5,5	17
Facteur VIII	186	9	25
Thyroglobuline	300	8,7	36
Dystrophine*	plus de 2 000	17	plus de 50

\* Une forme modifiée de ce gène provoque la dystrophie musculaire de Duchenne.  
La taille des gènes indiquée ici comprend à la fois la partie transcrite du gène et les séquences d'ADN régulatrices voisines. (D'après des données fournies par Victor McKusick.)

**Lengths are in kbp**

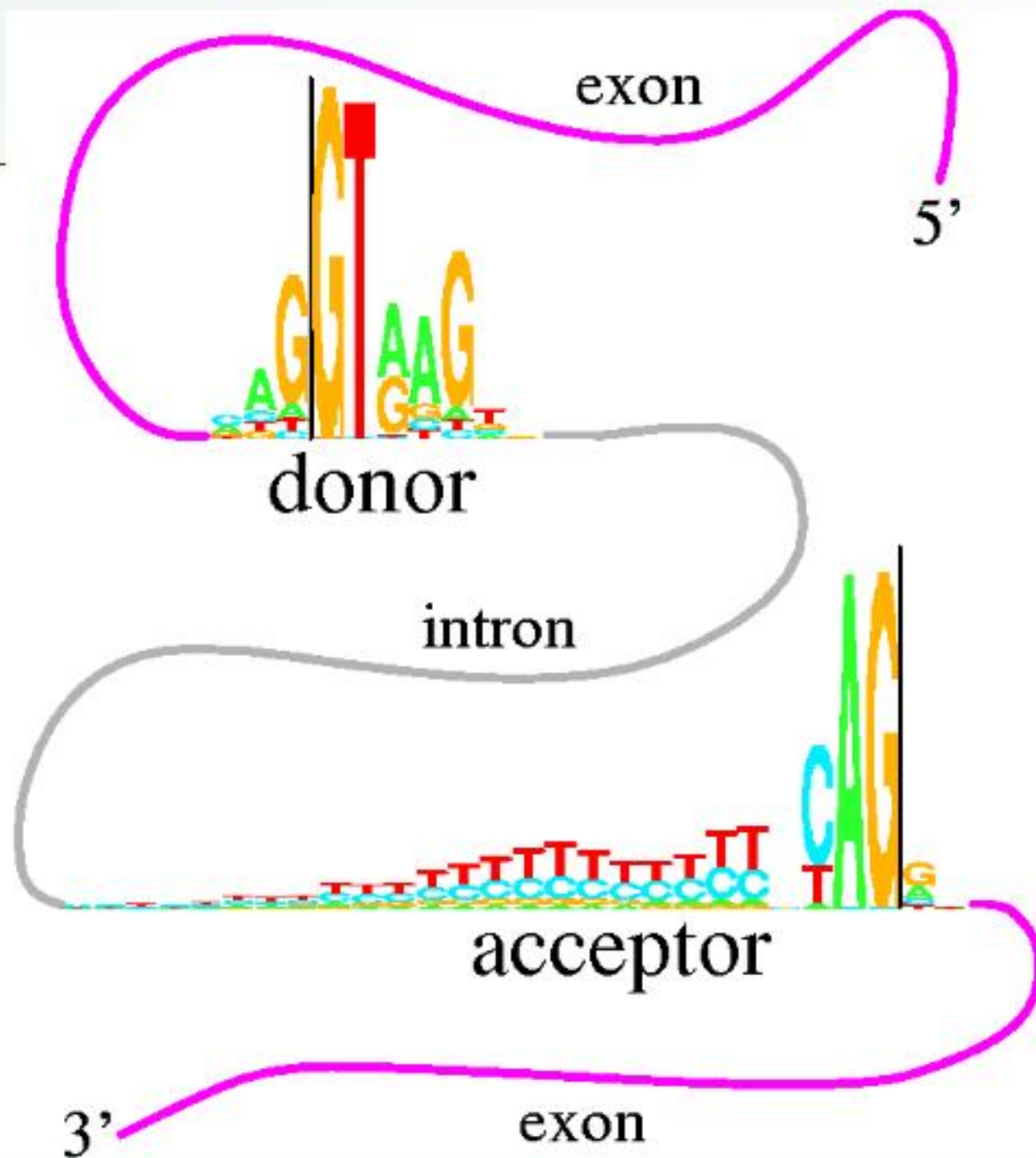


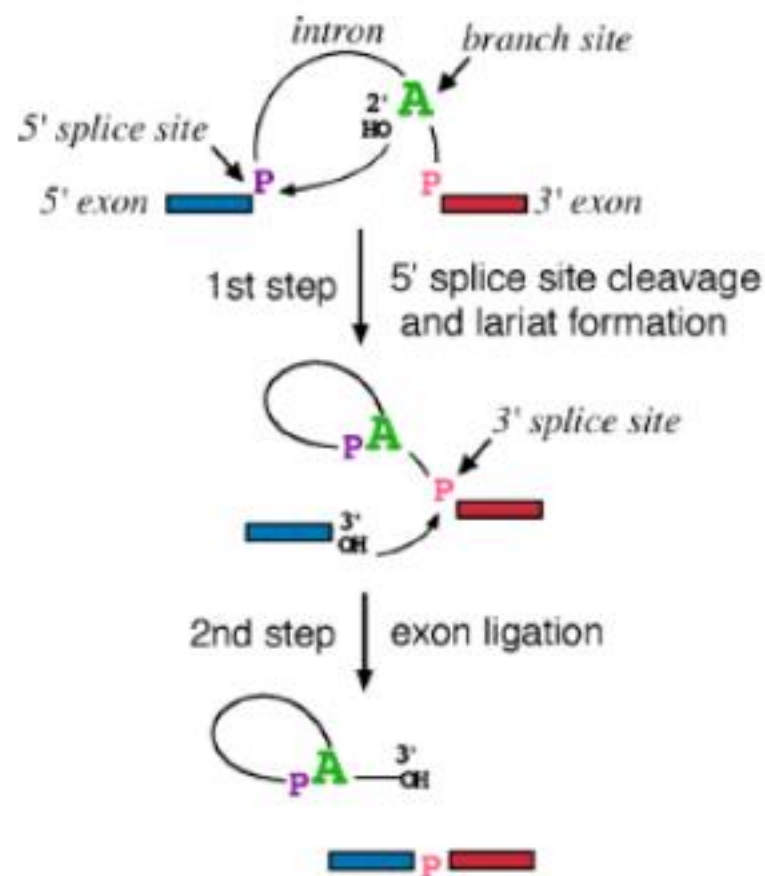
- Protein coding
  - most genes
- RNA genes
  - rRNA
  - tRNA
  - snRNA (small nuclear RNA)
  - snoRNA (small nucleolar RNA)



## Основные типы информации, использующихся для распознавания генов

- Информационные свойства сигналов, входящих в структуру генов (донорные и акцепторные сайты, сайты связывания транскрипционных факторов и т.д.)
- Свойства контекста (смещения частот использования кодонов в кодирующих районах, длина открытой рамки считывания и т.д.)
- Сходство с известными гомологичными последовательностями представленными в базах данных.





Chemical mechanism of intron excision by the spliceosome and group II introns.

- Consensus 'GT-AG'
  - Donor 5'
    - (A,C)AG/GT(A,G)AGT
  - Acceptor 3'
    - TTTTNCAG/GCCCCC
  - Branch
    - CT(G,A)A(C,T)

# Весовые матрицы, описывающие донорный и акцепторный сайты сплайсинга.

Весовая матрица донорного сайта сплайсинга

	-3	-2	-1	1	2	3	4	5	6
<b>A</b>	34.0	60.4	9.2	0.0	0.0	52.6	71.3	7.1	16.0
<b>C</b>	36.3	12.9	3.3	0.0	0.0	2.8	7.6	5.5	16.5
<b>G</b>	18.3	12.5	80.3	100	0.0	41.9	11.8	81.4	20.9
<b>U</b>	11.4	14.2	7.3	0.0	100	2.5	9.3	5.9	46.2

Весовая матрица акцепторного сайта сплайсинга

	-14	-13	-12	-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	1
<b>A</b>	11.1	12.7	3.2	4.8	12.7	8.7	16.7	16.7	12.7	9.5	26.2	6.3	100	0.0	21.4
<b>C</b>	36.5	30.9	19.1	23.0	34.9	39.7	34.9	40.5	40.5	36.5	33.3	68.2	0.0	0.0	7.9
<b>G</b>	9.5	10.3	15.1	12.7	8.7	9.5	16.7	4.8	2.4	6.3	13.5	0.0	0.0	100	62.7
<b>U</b>	38.9	41.3	58.7	55.6	42.1	40.5	30.9	37.3	44.4	47.6	27.0	25.4	0.0	0.0	7.9



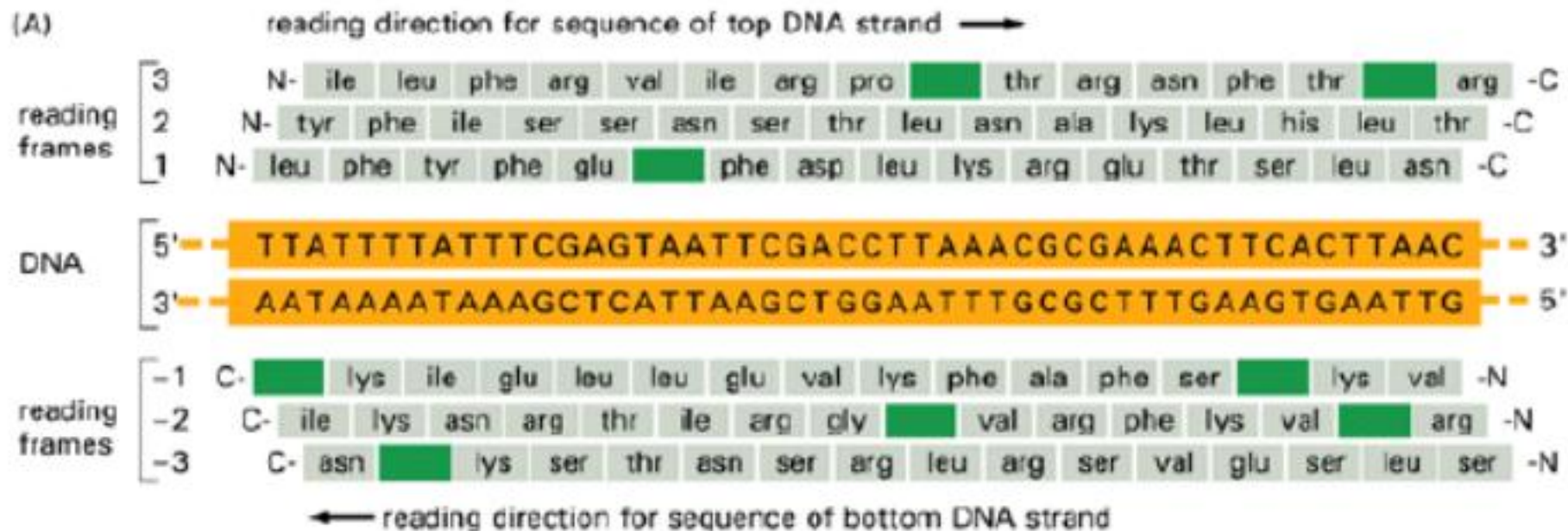
- Стартовый кодон

**ATG**

- Стоп - кодоны:

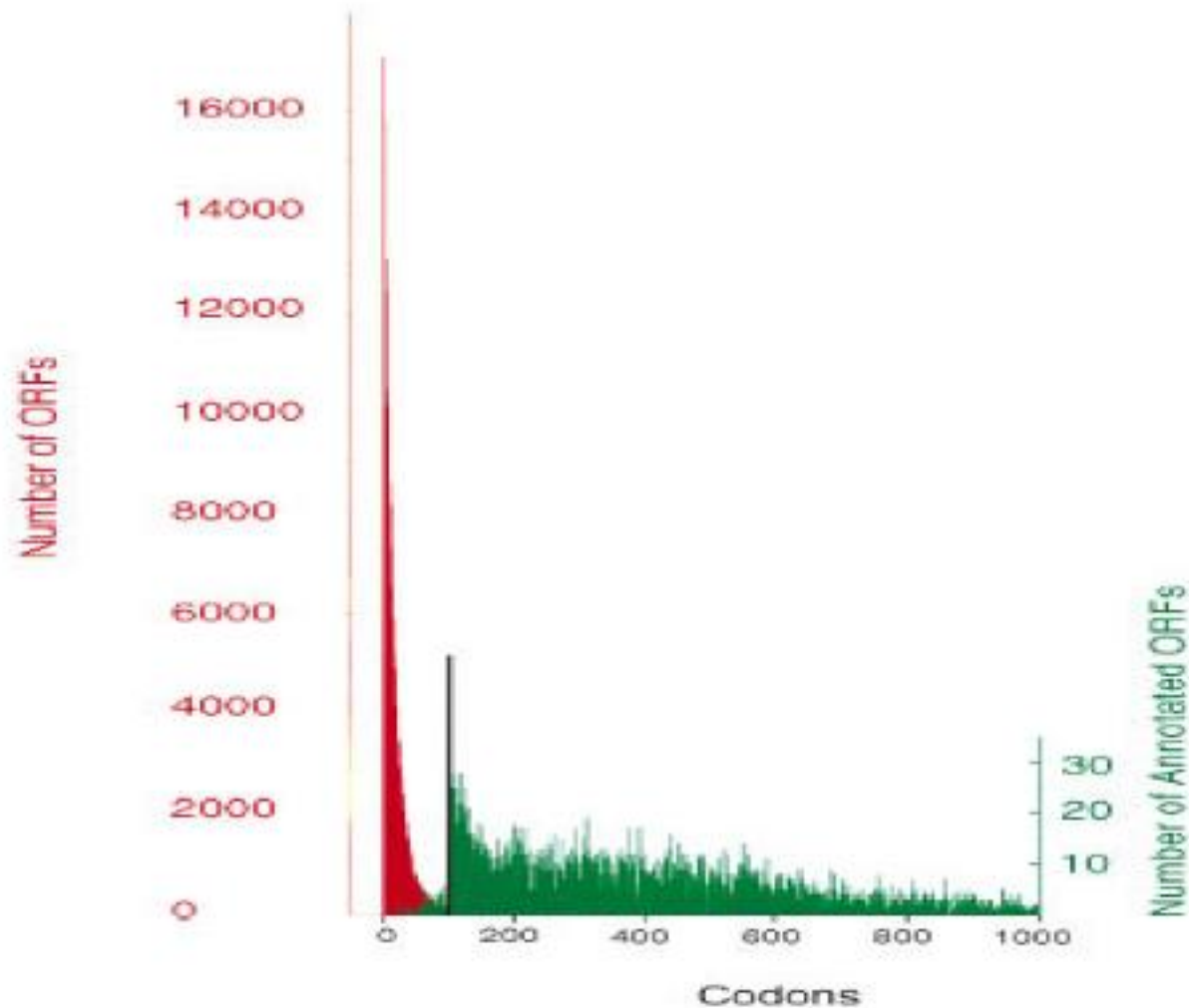
**– TAA, TAG, TGA**







# Распределение длин ORF в геноме *S. cerevisiae*



Basrai MA, Heter P, and Boeke J Genome Research 1997 7:768-771



# Частоты использования кодонов кодирующих серин у различных организмов

Example: Codon frequencies (%) for serine (Ser or S) codons in different organisms

Codon	<i>E. coli</i>	Fruitfly	Man	Maise	Yeast
AGT	3	1	10	3	5
AGC	20	23	34	30	4
TCG	4	17	9	22	1
TCA	2	2	5	4	6
TCT	34	9	13	4	52
TCC	37	48	28	37	33

The most frequently used codons have the highest chance of being present in the CDS -> this is used in gene prediction

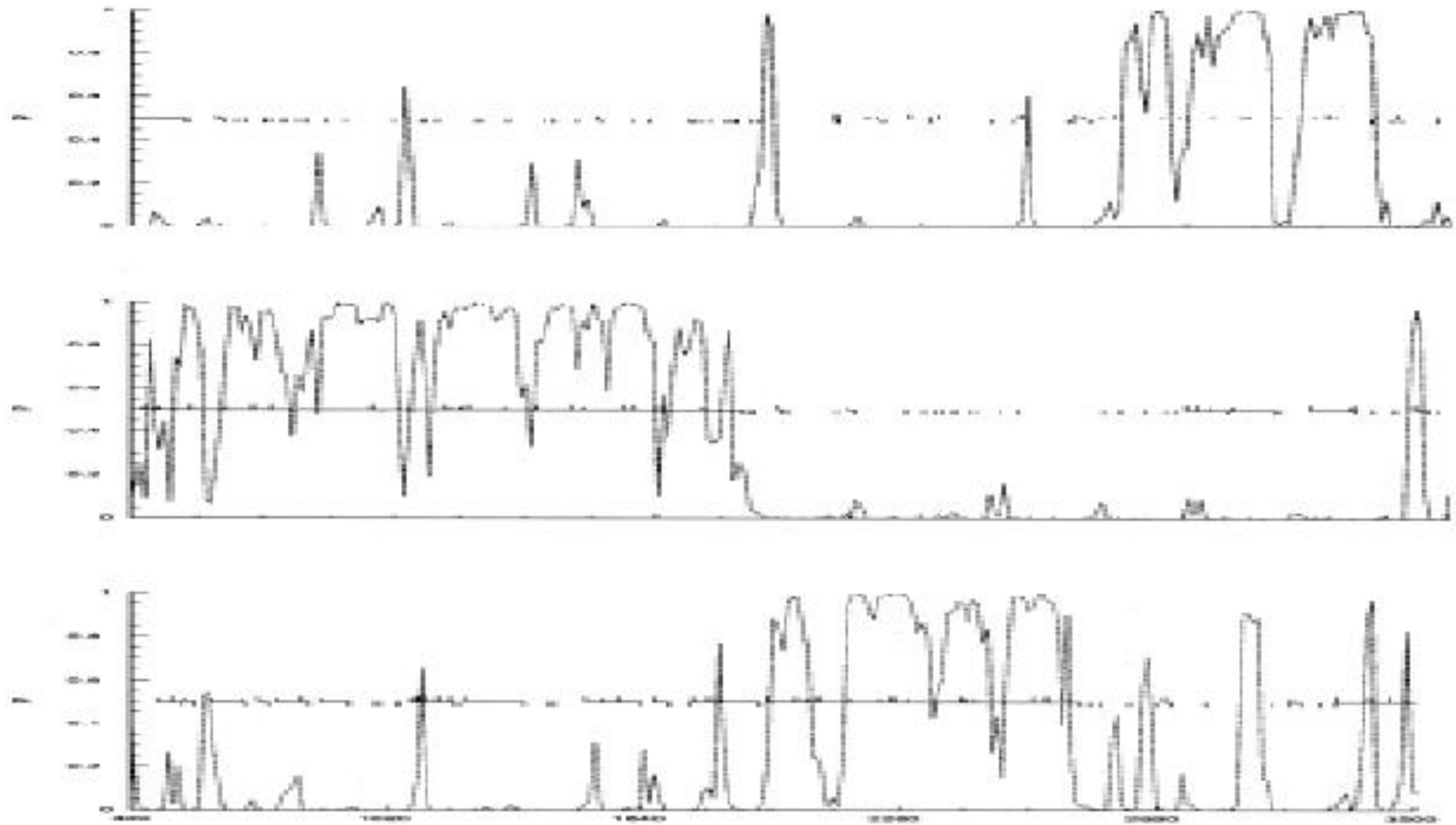


## Индекс адаптации кодонов (CAI)

$$CAI = \prod_{i=codons} \left[ \frac{f_{codon_i}}{f_{(codon_i)_{max}}} \right]$$

# Индексы адаптации кодонов, рассчитанные для человека и дрожжей

		hs	sc			hs	sc			hs	sc			hs	sc
UUU	Phe	16.6	26.0	UCU	Ser	14.5	23.6	UAU	Tyr	12.1	18.8	UGU	Cys	9.7	8.0
UUC	Phe	20.7	18.2	UCC	Ser	17.7	14.2	UAC	Tyr	16.3	14.7	UGC	Cys	12.4	4.7
UUA	Leu	7.0	26.3	UCA	Ser	11.4	18.8	UAA	stop	0.7	1.0	UGA	stop	1.3	0.6
UUG	Leu	12.0	27.1	UCG	Ser	4.5	8.6	UAG	stop	0.5	0.5	UGG	Trp	13.0	10.3
CUU	Leu	12.4	12.2	CCU	Pro	17.2	13.6	CAU	His	10.1	13.7	CGU	Arg	4.7	6.5
CUC	Leu	19.3	5.4	CCC	Pro	20.3	6.8	CAC	His	14.9	7.8	CGC	Arg	11.0	2.6
CUA	Leu	6.8	13.4	CCA	Pro	16.5	18.2	CAA	Gln	11.8	27.5	CGA	Arg	6.2	3.0
CUG	Leu	40.0	10.4	CCG	Pro	7.1	5.3	CAG	Gln	34.4	12.2	CGG	Arg	11.6	1.7
AUU	Ile	15.7	30.2	ACU	Thr	12.7	20.2	AAU	Asn	16.8	36.0	AGU	Ser	11.7	14.2
AUC	Ile	22.3	17.1	ACC	Thr	19.9	12.6	AAC	Asn	20.2	24.9	AGC	Ser	19.3	9.7
AUA	Ile	7.0	17.8	ACA	Thr	14.7	17.7	AAA	Lys	23.6	42.1	AGA	Arg	11.2	21.3
AUG	Met	22.2	20.9	ACG	Thr	6.4	8.0	AAG	Lys	33.2	30.8	AGG	Arg	11.1	9.3
GUU	Val	10.7	22.0	GCU	Ala	18.4	21.1	GAU	Asp	22.2	37.8	GGU	Gly	10.9	23.9
GUC	Val	14.8	11.6	GCC	Ala	28.6	12.6	GAC	Asp	26.5	20.4	GGC	Gly	23.1	9.7
GUA	Val	6.8	11.7	GCA	Ala	15.6	16.2	GAA	Glu	28.6	45.9	GGA	Gly	16.4	10.9
GUG	Val	29.3	10.7	GCG	Ala	7.7	6.1	GAG	Glu	40.6	19.1	GGG	Gly	16.5	6.0



$$S = s_1 s_2 s_3 s_4 \dots$$

$$P_0(s) = p(s_1) \cdot p(s_2) \cdot p(s_3) \cdot \dots = \prod_{i=1}^N p(s_i)$$

$$P_1(s) = p(s_1) \cdot p(s_2 | s_1) \cdot p(s_3 | s_2) \cdot \dots = p(s_1) \cdot \prod_{i=2}^N p(s_i | s_{i-1})$$

$$P_2(s) = p(s_1 s_2) \cdot p(s_3 | s_1 s_2) \cdot p(s_4 | s_2 s_3) \cdot \dots = p(s_1 s_2) \cdot \prod_{i=3}^N p(s_i | s_{i-2} s_{i-1})$$



- Bayes' Rule

$$\text{posterior} \rightarrow P(M | D) = \frac{\overset{\text{likelihood}}{P(D | M)} \overset{\text{prior}}{P(M)}}{\underset{\text{marginal}}{P(D)}}$$

- $M$ : the model,  $D$ : data or evidence

$$P(D) = \sum P(D | M)P(M) \cdot \text{discrete}$$

$$= \int P(D | M)P(M)dM \cdot \text{continuous}$$





## LOG-ODDS Ratio

$$\log \frac{P(\mathbf{c}|\mathit{data})}{P(\mathbf{nc}|\mathit{data})} = \log \frac{P(\mathit{data}|\mathbf{c})}{P(\mathit{data}|\mathbf{nc})} + \log \frac{P(\mathbf{c})}{P(\mathbf{nc})}$$



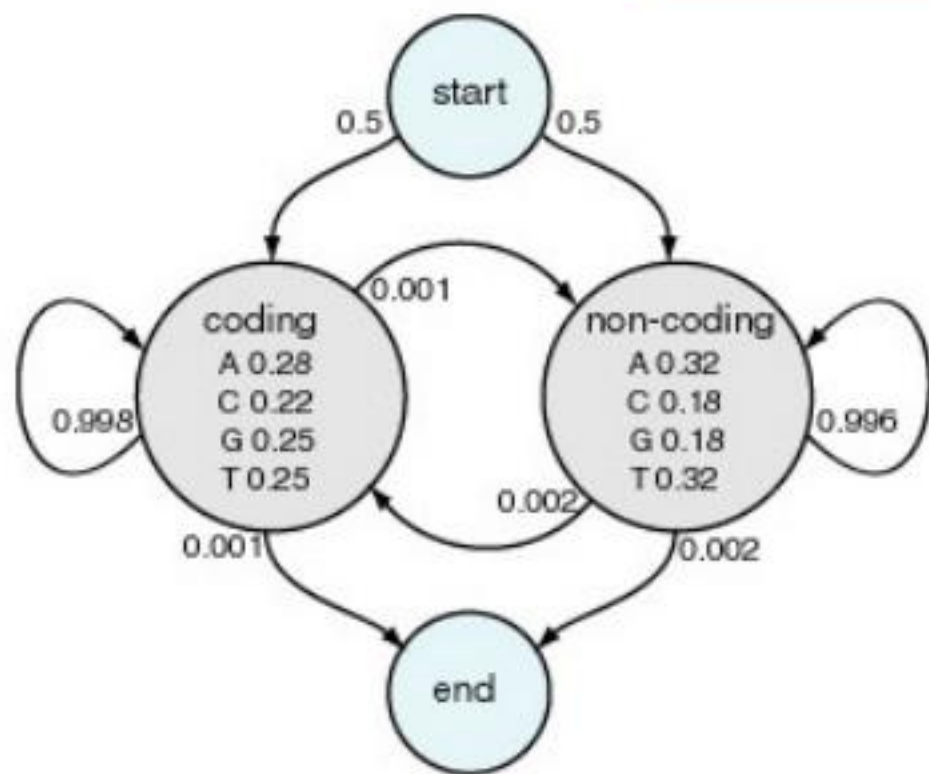
- GeneMark uses a **5th-order Markov chain**:

$x = \underline{GCTAC}GCGCTAGGAT\dots$

$$\Pr(x) = \Pr(GCTAC) \times \Pr(G | \underline{GCTAC}) \times \Pr(C | \underline{CTACG}) \dots$$



# Модель скрытой Марковской цепи (HMM)



$$\Phi = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0.5 & 0.998 & 0.002 & 0 \\ 0.5 & 0.001 & 0.996 & 0 \\ 0 & 0.001 & 0.002 & 0 \end{bmatrix}$$

$$H = \begin{bmatrix} 0.28 & 0.32 \\ 0.22 & 0.18 \\ 0.25 & 0.18 \\ 0.25 & 0.32 \end{bmatrix}$$

$H(m, y_i)$  = probability of emitting character  $y_i$  in state  $m$ ;

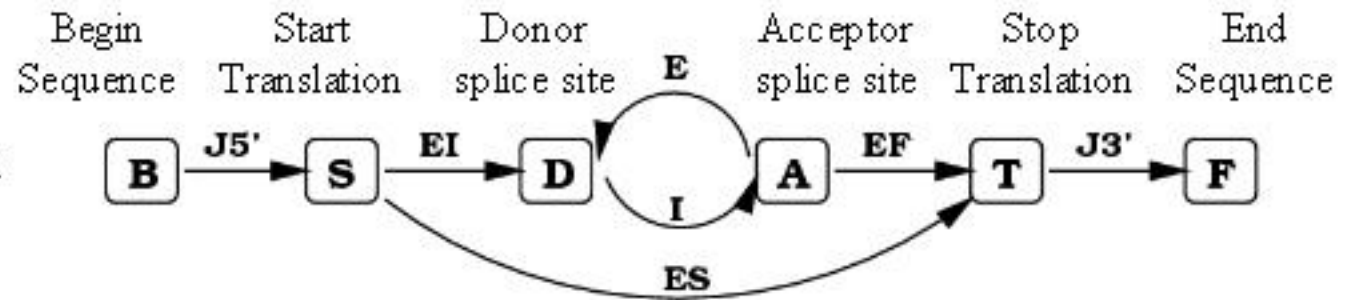
$\Phi_{mk}$  = probability of transition from state  $k$  to  $m$ .

## Некоторые программы распознавания генов использующие скрытые Марковские цепи

- GENSCAN (Burge 1997)
- FGENESH (Solovyev 1997)
- HMMgene (Krogh 1997)
- GENIE (Kulp 1996)
- GENMARK (Borodovsky & McIninch 1993)
- VEIL (Henderson, Salzberg, & Fasman 1997)



- J5' – 5' UTR
- EI – Initial Exon
- E – Exon, Internal Exon
- I – Intron
- EF – Final Exon
- ES – Single Exon
- J3' – 3'UTR

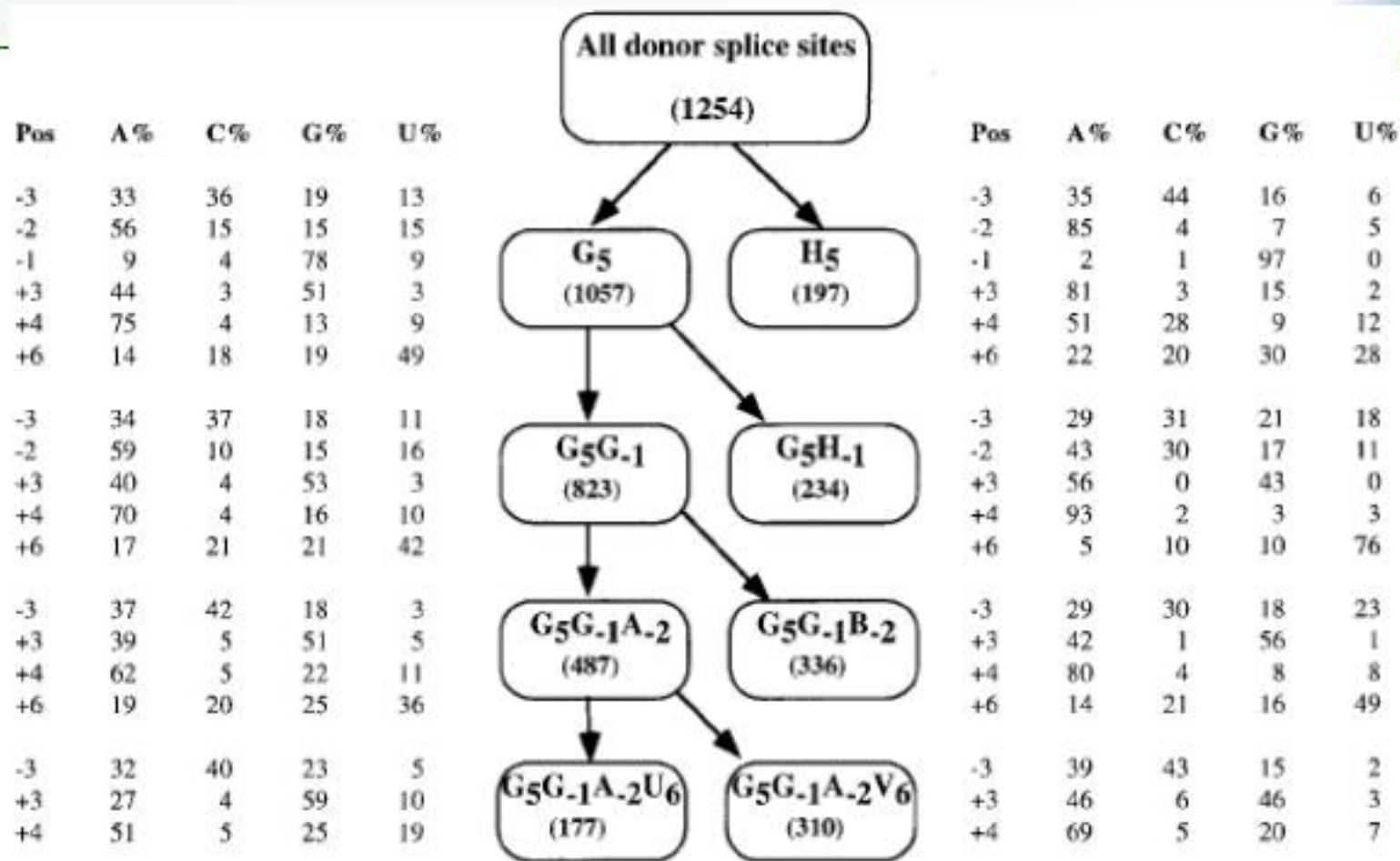




- Учитывает распределение длин интронов и экзонов
- Весовая позиционная матрица (WMM) для описания TATA box, PolyA signal, CAP end and Transcription Initiation End (TIE) of 5'UTR.
- Weight Array model (WAM) для описания акцепторного сайта связывания.
- Дерево решений (maximal dependence decomposition) для моделирования донорных сайтов сплайсинга



# Maximal Dependence Decomposition (MDD)



All sites:

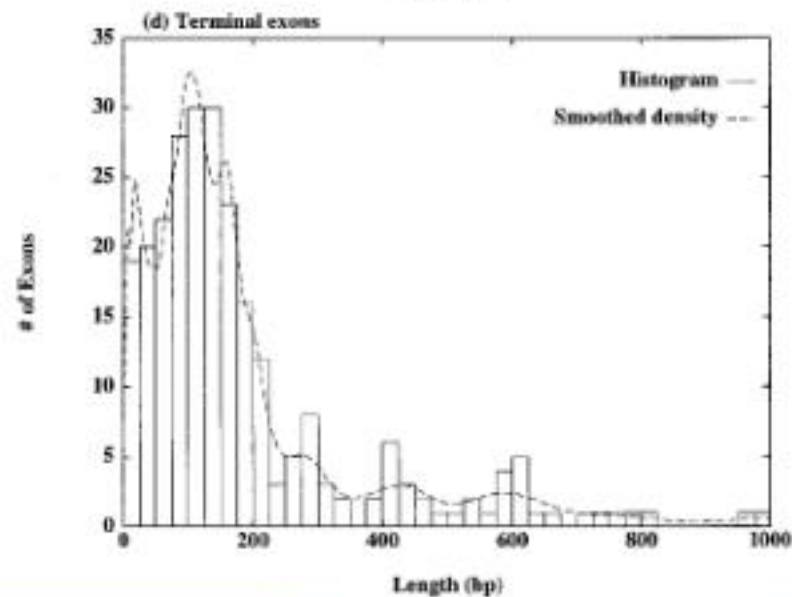
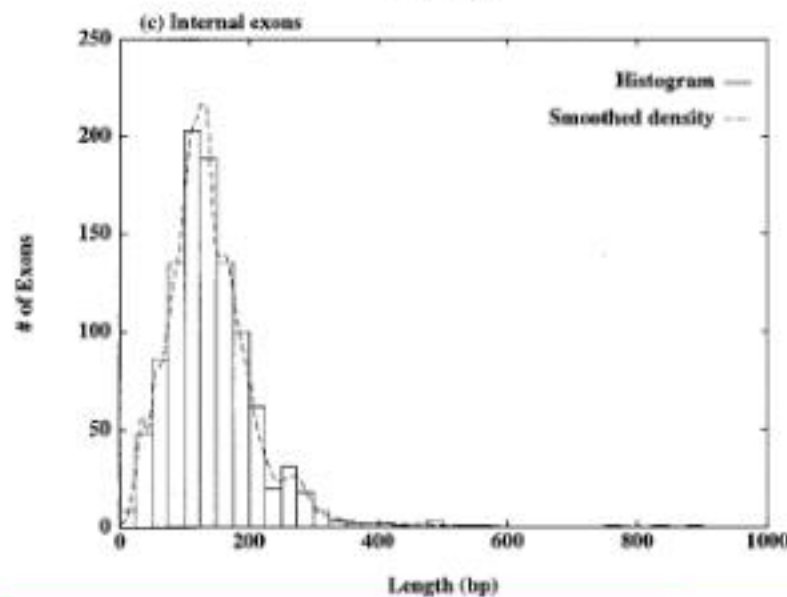
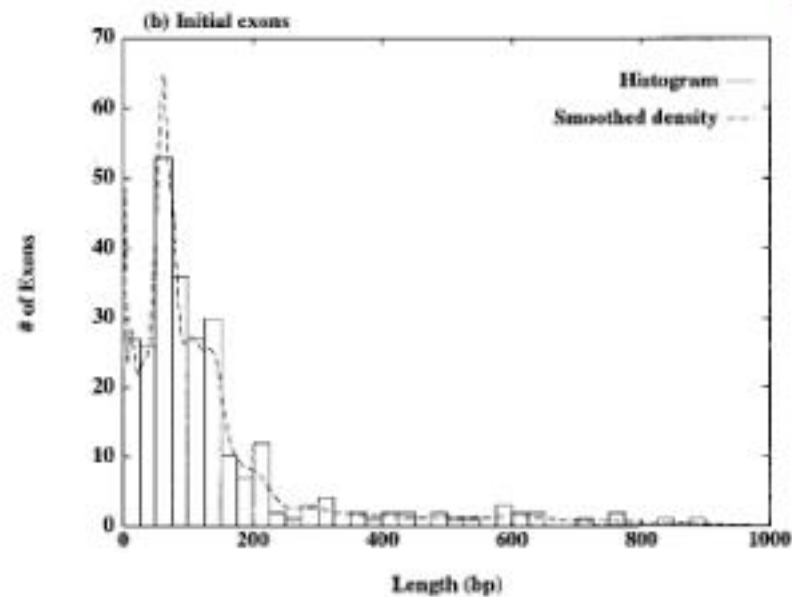
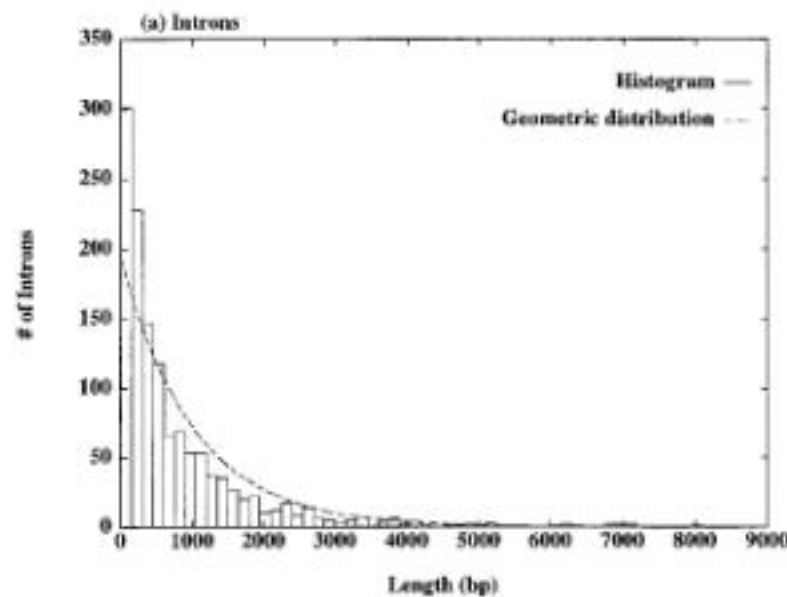
Position

Base	-3	-2	-1	+1	+2	+3	+4	+5	+6
A%	33	60	8	0	0	49	71	6	15
C%	37	13	4	0	0	3	7	5	19
G%	18	14	81	100	0	45	12	84	20
U%	12	13	7	0	100	3	9	5	46

UI snRNA: 3' G U C C A U U C A 5'



# Распределение длин экзонов и интронов, входящих в структуру человеческих генов

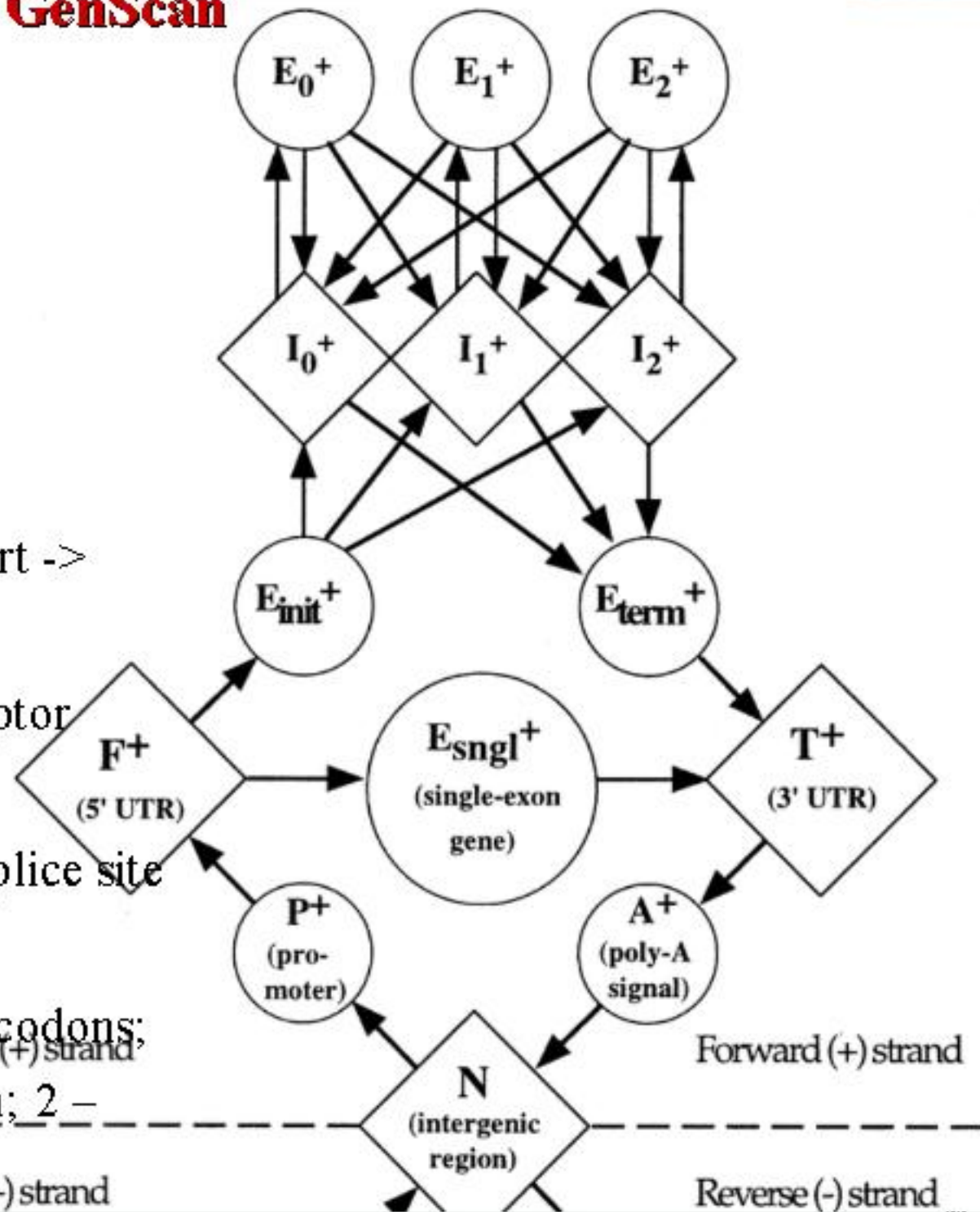






# GenScan

- N - intergenic region
- P - promoter
- F - 5' untranslated region
- $E_{sngl}$  - single exon (intronless)  
(translation start -> stop codon)
- $E_{init}$  - initial exon (translation start -> donor splice site)
- $E_k$  - phase k internal exon (acceptor splice site -> donor splice site)
- $E_{term}$  - terminal exon (acceptor splice site -> stop codon)
- $I_k$  - phase k intron: 0 - between codons;  
1 - after the first base of a codon; 2 -  
after the second base of a codon





GENSCAN - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites Media Print Mail News RSS

Address <http://bioweb.pasteur.fr/seqanal/interfaces/genSCAN.html> Go Links

## GENSCAN : Gene Identification (C. Burge)

Reset Run genSCAN  your e-mail

(● = required, ● = conditionally required)

● DNA Sequence File: please enter [either](#)

1. the name of a file:  Browse...

2. or the actual data here:

(sequence [format](#))

● [Parameter file](#) ?  Arabidopsis  Human/3e  Mouse

[Output parameters](#)

---

### Output parameters

Verbose output (-v)

Print predicted coding sequences (-cde)

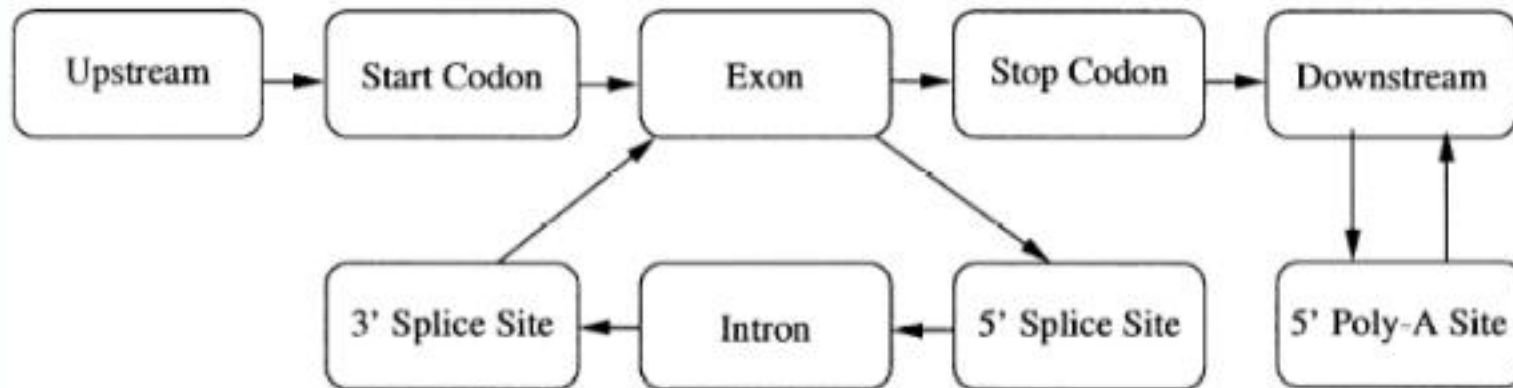
Identify [suboptimal](#) exons (-subopt)

[Cutoff](#) for sub-optimal exons

Done Internet



## VEIL (Viterbi Exon-Intron Locator)

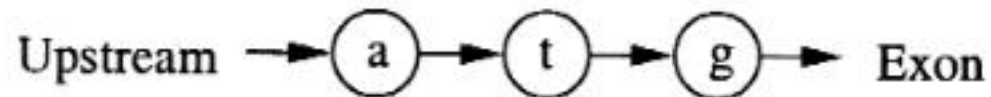


- **Each box represents a separate HMM**
- **Edges represent multiple edges between HMMs**
- **Total model has 241 states and 1003 edges**



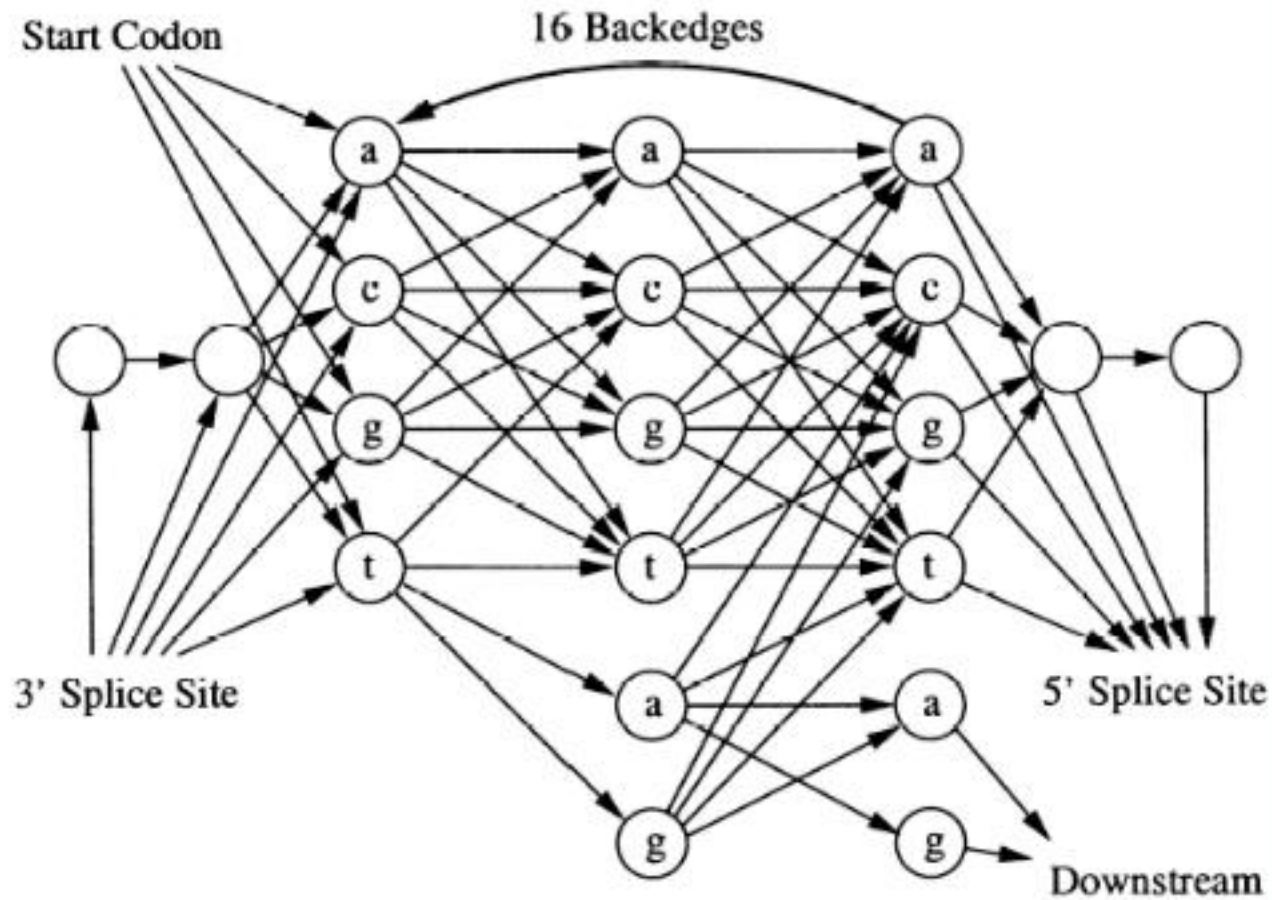
## Upstream/Downstream and Start Codons

- **Upstream model:**
  - 15-stage chain
  - Loops at end to absorb extra bases
- **Start Codon Model:**



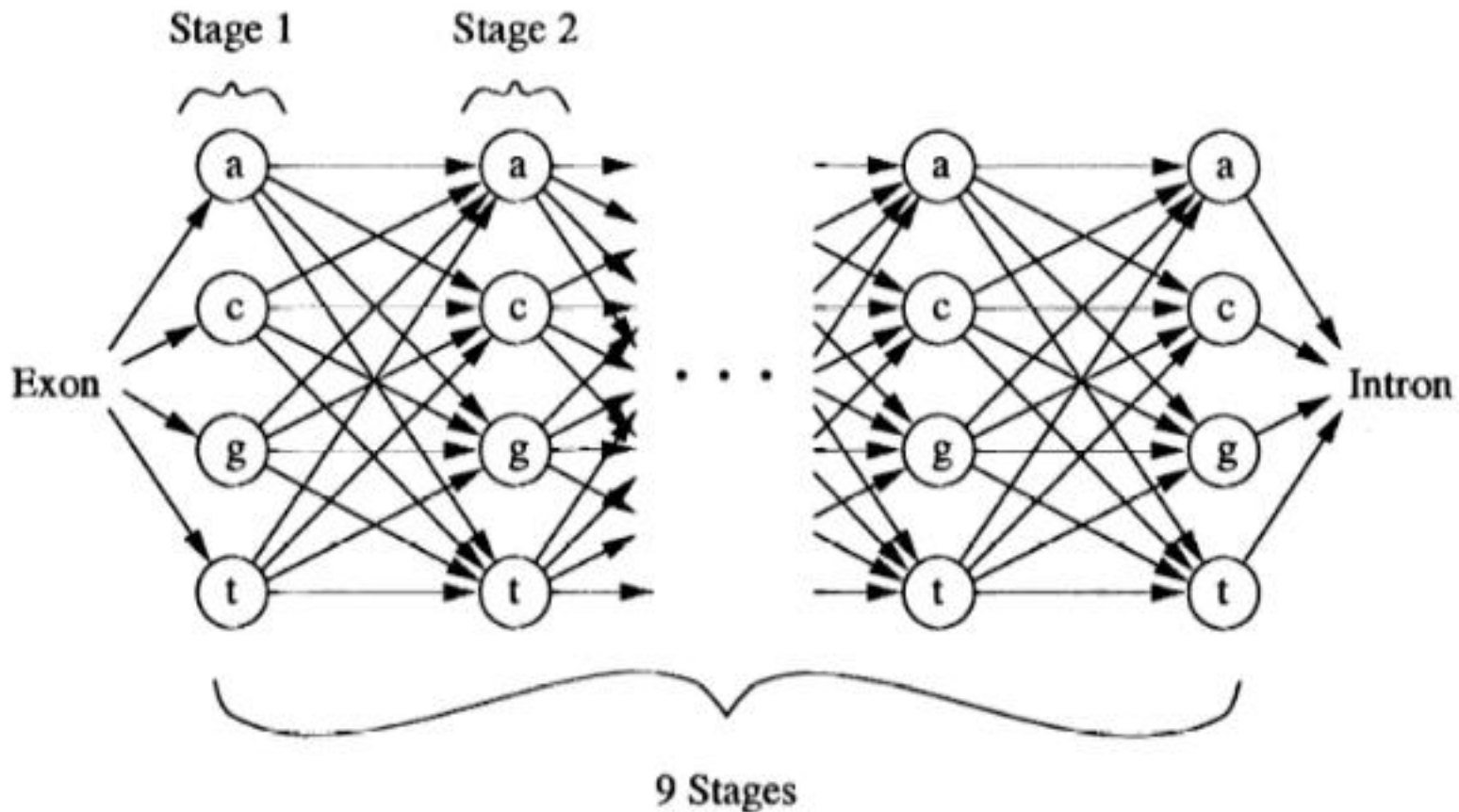
- **Downstream model:**
  - 10-stage chain
  - Loops to absorb extra bases

# Exon and Stop Codon Models in VEIL





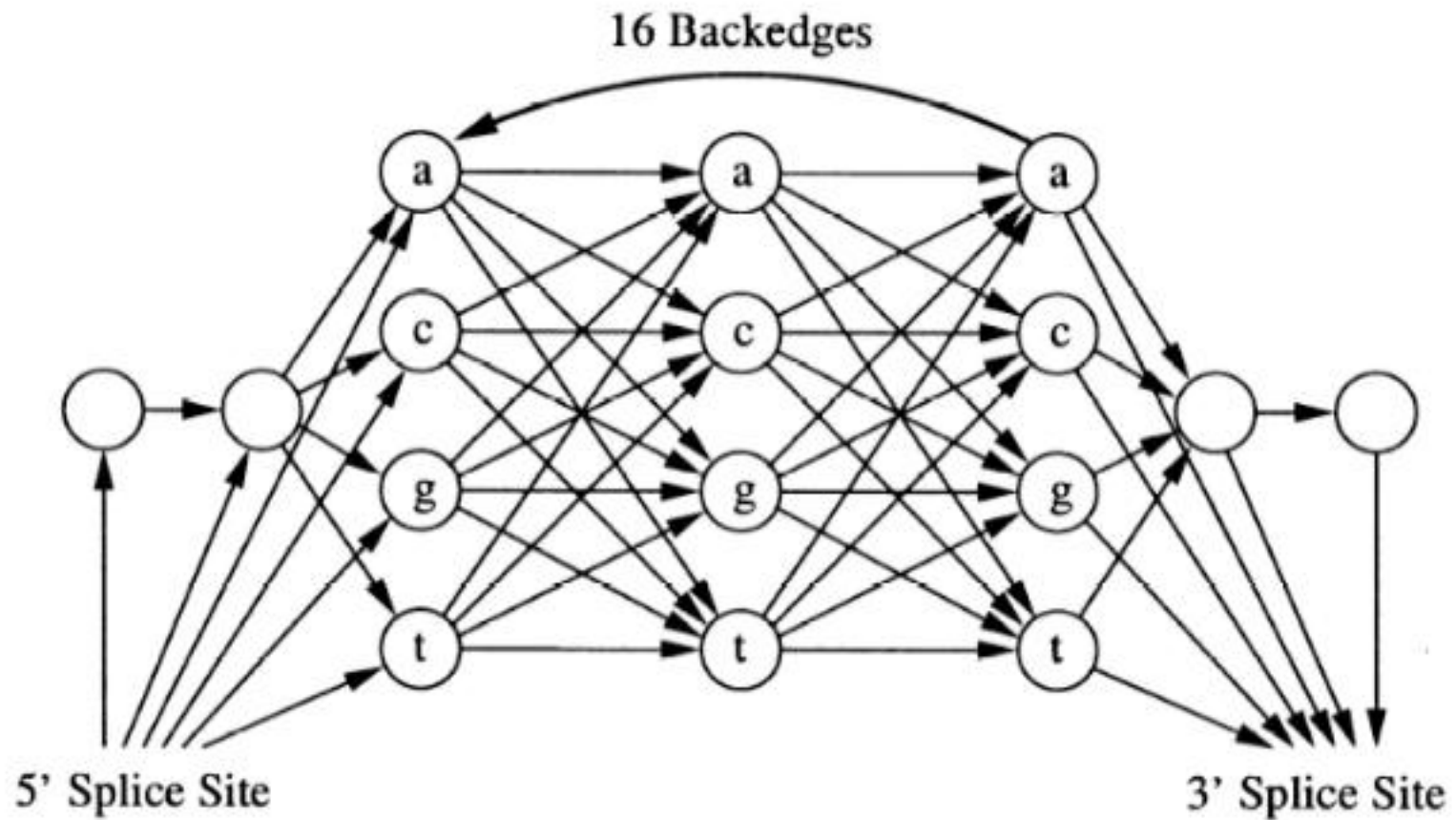
## Donor Site (5' splice site) Model



- Acceptor site (3' splice site) has 15 stages
- Length based on consensus sequences (Mount et. al, 1995)



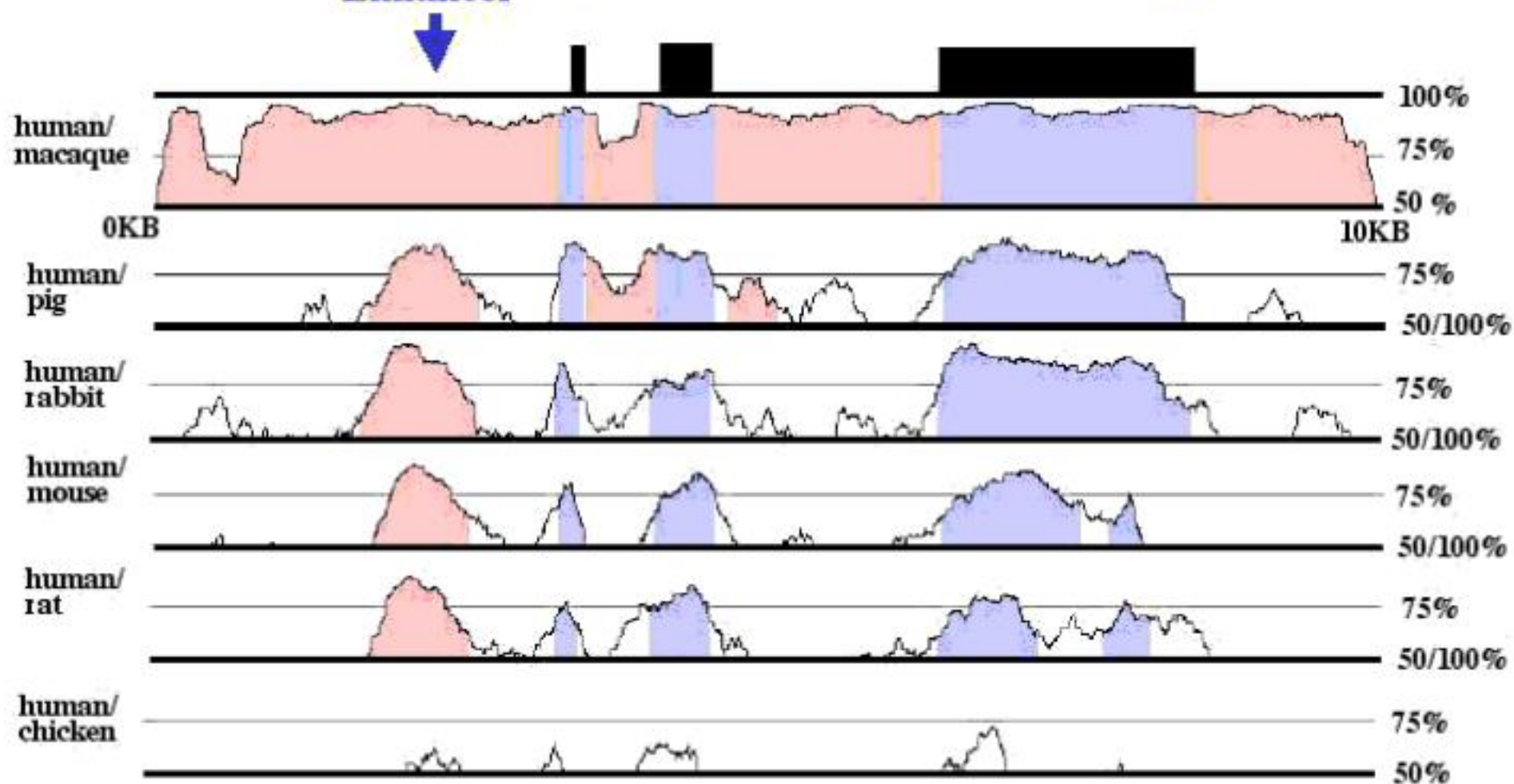
# Intron Model in VEIL



## Multi-Species Comparative Analysis

Liver  
Enhancer

Apolipoprotein AI gene







# Выравнивание ортологических генов



```
50      .      :      .      :      .      :      .      :      .      :
247 GGTGAGGTCGAGGACCCTGCA  CGGAGCTGTATGGAGGGCA  AGAGC
      |:  ||  ||||:  ||||  --:||  |||  |:|  |||---|||
368 GAGTCGGGGGAGGGGGCTGCTGTTGGCTCTGGACAGCTTGCATTGAGAGG

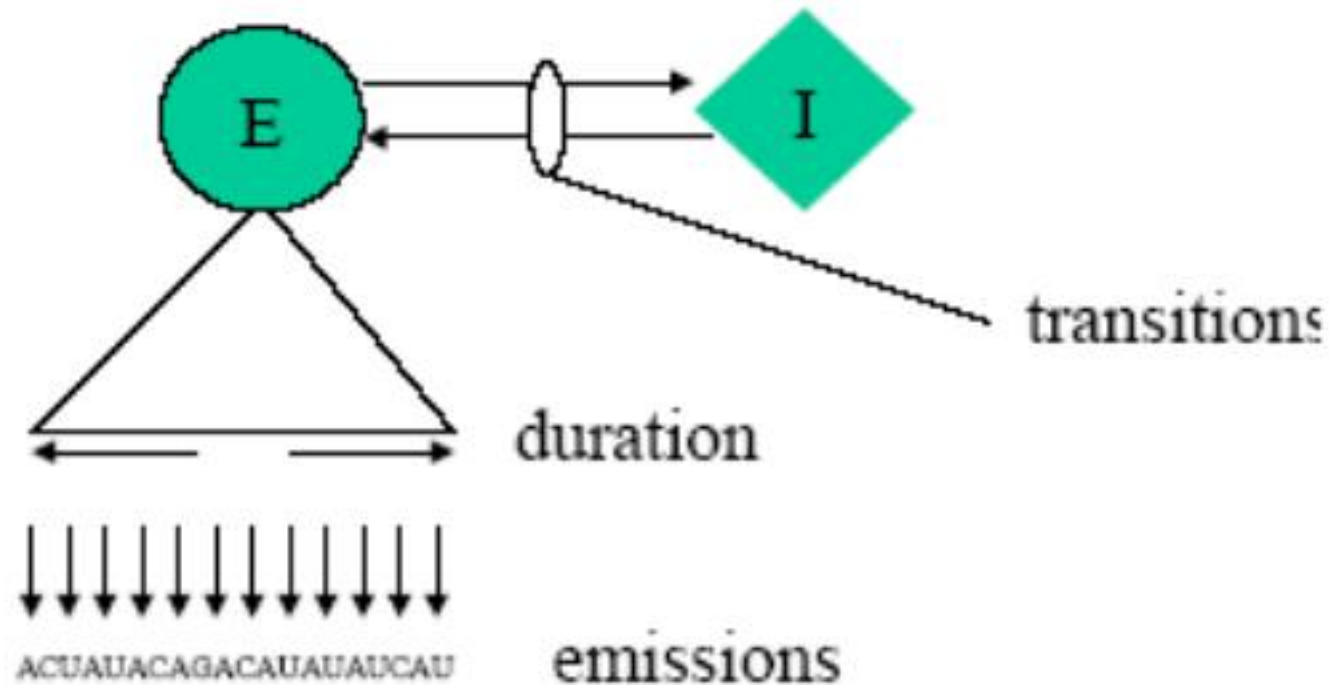
100     .      :      .      :      .      :      .      :      .      :
292 TTC                CTACAGAAAAGTCCCAGCAAGGAGCCACACTTCACTG
      |||-----||  |  |:|  |:  |||::|:|:-||  ||:|  |
418 TTCTGGCTACGCTCTCCCTTAGGGACTGAGCAGAGGGCT  CAGGTCGCGG

150     .      :      .      :      .      :      .      :      .      :
332                ATGTCGAGGGGAAGACATCATTCGGGATGTCAGTG
      -----|||:|||||:|||||:|||||:|||||:|||||:|||||
467 TGGGAGATGAGGCCAATGTCGAGGGGAAGACATCATTTGGGATGTCAGTG

200     .      :      .      :      .      :      .      :      .      :
367 TTCAACCTCAGCAATGCCATCATGGGCAGCGGCATCCTGGGACTCGCCTA
      ||||:|||||:|||||:|||||:||  ||:||||:|||||
517 TTCAATCTCAGCAACGCCATCATGGGCAGTGGAAATCTGGGGCTCGCCTA
```



# Twinscan





1. Align the two sequences (eg. from human and mouse)
2. Mark each human base as gap ( - ), mismatch ( : ), match ( | )

New "alphabet":  $4 \times 3 = 12$  letters

$\Sigma = \{ A-, A:, A|, C-, C:, C|, G-, G:, G|, U-, U:, U| \}$



3. Run Viterbi using emissions  $e_k(b)$   
where  $b \in \{A-, A:, A|, \dots, T| \}$

Note:

Emission distributions  $e_k(b)$  estimated from real genes  
from human/mouse

$e_I(x|) < e_E(x|)$ : matches favored in exons

$e_I(x-) > e_E(x-)$ : gaps (and mismatches) favored in  
introns



Human:    ACGGCGACUGUGCACGU  
Mouse:     ACUGUGAC  GUGCACUU  
Align:     ||:|:||||-|||||||:|

Input to Twinscan HMM:

A| C| G: G| C: G| A| C| U- G| U| G| C| A| C| G: U|

Recall,      $e_G(A|) > e_I(A|)$   
            $e_G(A-) < e_I(A-)$

Likely exon



# ORFScan



```

a)
EST-1:   ATCGAT-GGATGAA-TAGAGC-CTAAC
EST-2:   ATCGAT-GGATGAA-TANAGCACTAACTG-ATACGG
EST-3:   ATCGATTGNAT
EST-4:   ATCGAT-GGATGAA-TAGAGCACTAACTG-ATACGGATGC-GGA
EST-5:                               GC-CTANCTG-ATACGGATGCCGGA
EST-6:   ATCGAT-GGATGAAAGTAGAGCTCTAACTG-ATACGGATGC-GGA
EST-7:   T-GGATGAA-TAGAGC-CTAACTGCATACGGATGC-GGA
EST-8:                               GG-TGC-GGA

```

Consensus: ATCGAT-GGATGAA-TAGAGC-CTAACTG-ATACGGATGC-GGA

```

|
b)
=ATCGAT-GGATGAA-TAGAGC=CTAACTG-ATACGGATGC-GGA=

```

```

|
c)
=ATCGATGGATGAATAGAGC=CTAACTGATACGGATGCCGGA=

```

```

|
d)
ATCGATGGATGAATAGAGC      CTAACTGATACGGATGCCGGA

```

```

|
e)
S1,1 ATC GAT GGA TGA ATA GAG C   S1,1 CT AAC TGA TAC GGA TGC GGA
S1,2  AT CGA TGG ATG AAT AGA GC   S1,2  C TAA CTG ATA CGG ATG CGG A
S1,3  A TCG ATG GAT GAA TAG AGC   S1,3  CTA ACT GAT ACG GAT GCG GA

```

```

|
f)
. . . . .
ATGAATAGAGC-CTAACTGATACGGATGCCGGA

```

```

g)
M N R ? L T D T D A

```

## Предсказание кодирующих частей генов с помощью программ ORFScan, ESTScan, GENSCAN, GRAIL, GENFIND на примере DT\_101886.

Ген human ferrochelatase. Ошибочно определенные позиции маркированы цветом

```
true      : mrslganmaaalraagvllrdplassswrvccpwrwksgaaaaavttetaqhaqgagkpqvqpqkrkpktilmlnmggpetlgdvhdfllrlfldqdlmtlpiqklapf
ORFScan   : MRSLGANMAAALRAAGULLRDPLASSSWRVCCPWRWKSAAAAAVTTETAQHAQGAKPQVQPQKRKPKTGILMLNMGGPETLGDVHD FLLRLFLDQDLMTLP IQNKLAPF
ESTScan   : NAFIGANMAAALRAAGULLRXSAGIQQLEGLS AMEVEVUGAAAAAVTTETAQHAQGAKPQVQPQKRKPKTGILMLNMGGPETLGDVHD FLLRLFLDQDLMTLP IQNKLAPF
GENSCAN   : _____ MAAALRAAGULLRDPLASSSWRVCCPWRWKSAAAAAVTTETAQHAQGAKPQVQPQKRKPKTGILMLNMGGPETLGDVHD FLLRLFLDQDLMTLP IQNKLAPF
GRAIL     : MRSLGANMAAALRAAGULLRDPLASSSWRVCCPWRWKSAAAAAVTTETAQHAQGAKPQVQPQKRKPKTGILMLNMGGPETLGDVHD FLLRLFLDQDLMTLP IQNKLAPF
GenFind   : _____ MAAALRAAGULLRDPLASSSWRVCCPWRWKSAAAAAVTTETAQHAQGAKPQVQPQKRKPKTGILMLNMGGPETLGDVHD FLLRLFLDQDLMTLP IQNKLAPF
```

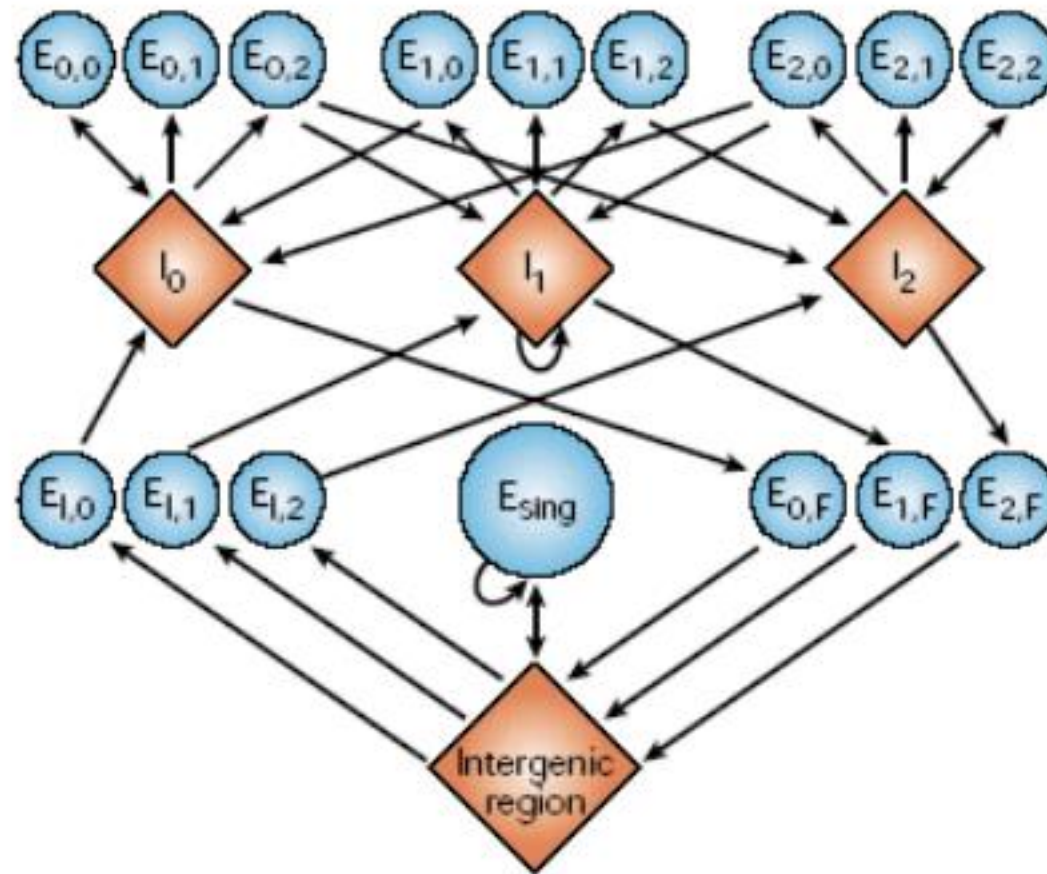
```
true      : iakrtpkiqeqyrriggspikivtskqgegmvklldel spntaphkyyigfryvhlpteeaeemer dgleraiaftqypqyscsttgsslnaiyryynqvgrkptmk
ORFScan   : IAKRTPKIQEQYRRLEADPPSRYY??SKQGE GNVKLLDEL SPNT APHKYYIGIRYVHPLTEE A IEEMERD GLERA I AFT QYP QYS CSTTGS SLNAIYRYYNQVGRKPTMK
ESTScan   : IAKRTPKIQEQYRRXWR IPHQIMDFQ QEGE GNVKLLDEL SPNT APHKYYIGIRYVHPLTEE A IEEMERD GLERA I AFT QYP QYS CSTTGS SLNAIYRYYNQVGRKPTMK
GENSCAN   : IAKRTPKIQEQYRRXWR IPHQIMDFQ AGRGH _____ APHKYYIGIRYVHPLTEE A IEEMERD GLERA I AFT QYP QYS CSTTGS SLNAIYRYYNQVGRKPTMK
GRAIL     : IAKRTPKIQEQYRRXWR IPHQIMDFQ AGRGH _____
GenFind   : IAKRTPKIQEQYRRXWR IPHQIMDFQ AGRGH _____ EE A IEEMERD GLERA I AFT QYP QYS CSTTGS SLNAIYRYYNQVGRKPTMK
```

```
true      : wstidrwpthhlliqcfadhilkel dhfplekrsewv ilfshslpmsvvnrgdyp qevs atvqkmerleycnp yrlvwoq skvgmpmlgp qtde s ikgl cergrkn i
ORFScan   : WSTIDRWPTHLL IQCFADHILKELDHFPLEKRSEWV ILFSSHSLPMSVVNRGDPYPQEVSATVQKMERLEYCNP YRLVWQ SKVGPMPMLGP QTDES IKGL CERGRKNI
ESTScan   : WSTIDRWPTHLL IQCFADHILKELDHFPLEKRSEWV ILFSSHSLPMSVVNRGDPYPQEVSATVQKMERLEYCNP YRLVWQ SKVGPMPMLGP QTDES IKGL CERGRKNI
GENSCAN   : WSTIDRWPTHLL IQCFADHILKELDHFPLEKRSEWV ILFSSHSLPMSVVNRGDPYPQEVSATVQKMERLEYCNP YRLVWQ SKVGPMPMLGP QTDES IKGL CERGRKNI
GRAIL     : _____ GE AAGQEVSATVQKMERLEYCNP YRLVWQ SKVGPMPMLGP _____
GenFind   : WSTIDRWPTHLL IQCFADHILKELDHFPLEKRSEWV ILFSSHSLPMSVVNRGDPYPQEVSATVQKMERLEYCNP YRLVWQ SKVGPMPMLGP QTDES IKGL CERGRKNI
```

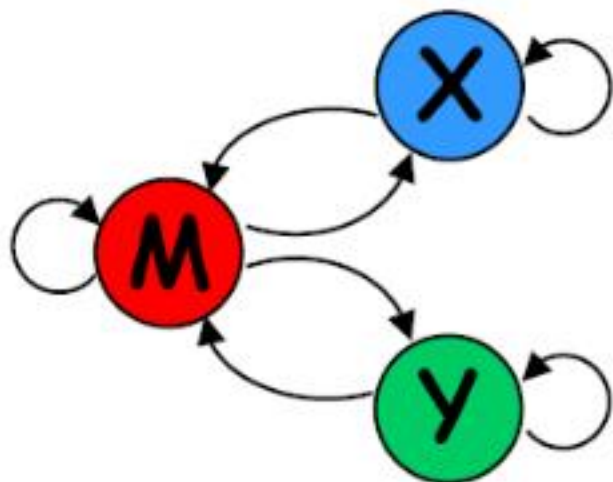
```
true      : llvpiafstdhietlyeldieysqvlakecgvnirraeslngnplfskal adlvshiqsnel cskqtlscplcwpvcvretksfftsqgl
ORFScan   : LLVPIAFTSDHIETLYELD IEYSQVLAKECGVNIRRAESLNG- INCS ?KALADLV?FTHPVKRAVFPQ AADPDCPLCWPVCVRETKEFFT ?P AAVTP AGGPRGUSKCPTRYLRCGEGVI
ESTScan   : LLVPIAFTSDHIETLYELD IEYSQVLAKDGVN DKELSLFWGKFKSLKAPGRLGAFTHPSQTS CVPKQLTPELWVPCRETKEFXHP ASS CEPPPUDPVALANAQPPDTSDUER
GENSCAN   : LLVPIAFTSDHIETLYELD IEYSQVLAKECGS _____
GRAIL     : _____
```



# Pair HMM







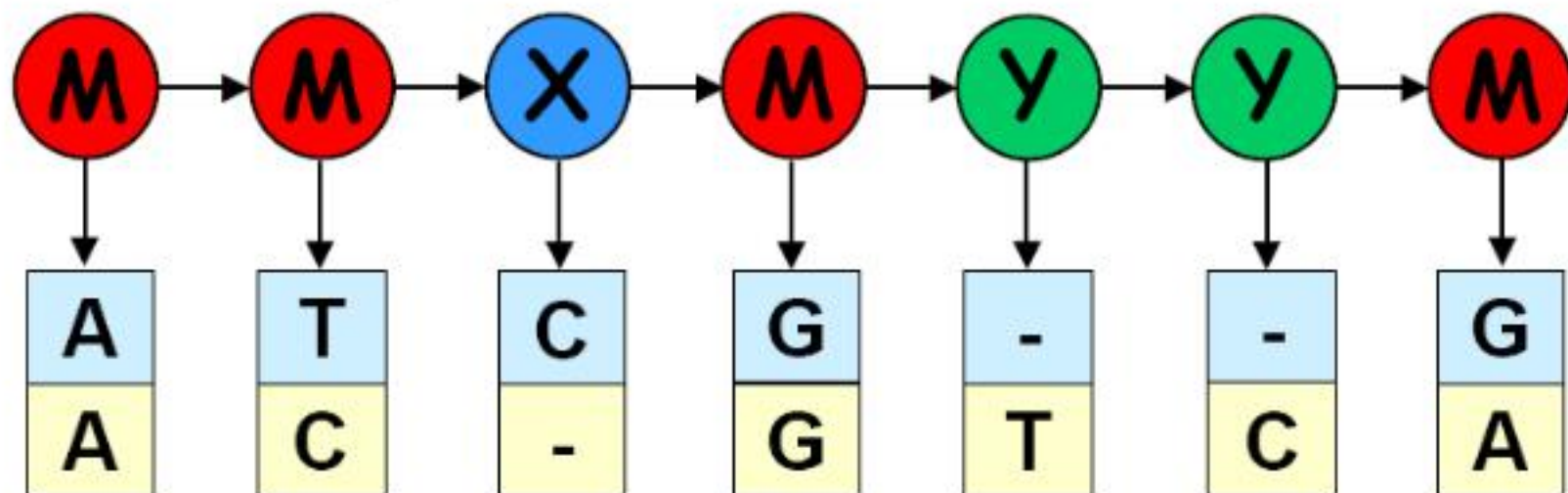
**M** = (mis)match

**X** = insert seq1

**Y** = insert seq2



Hidden sequence:



Hidden alignment:

ATCG--G  
AC-GTCA

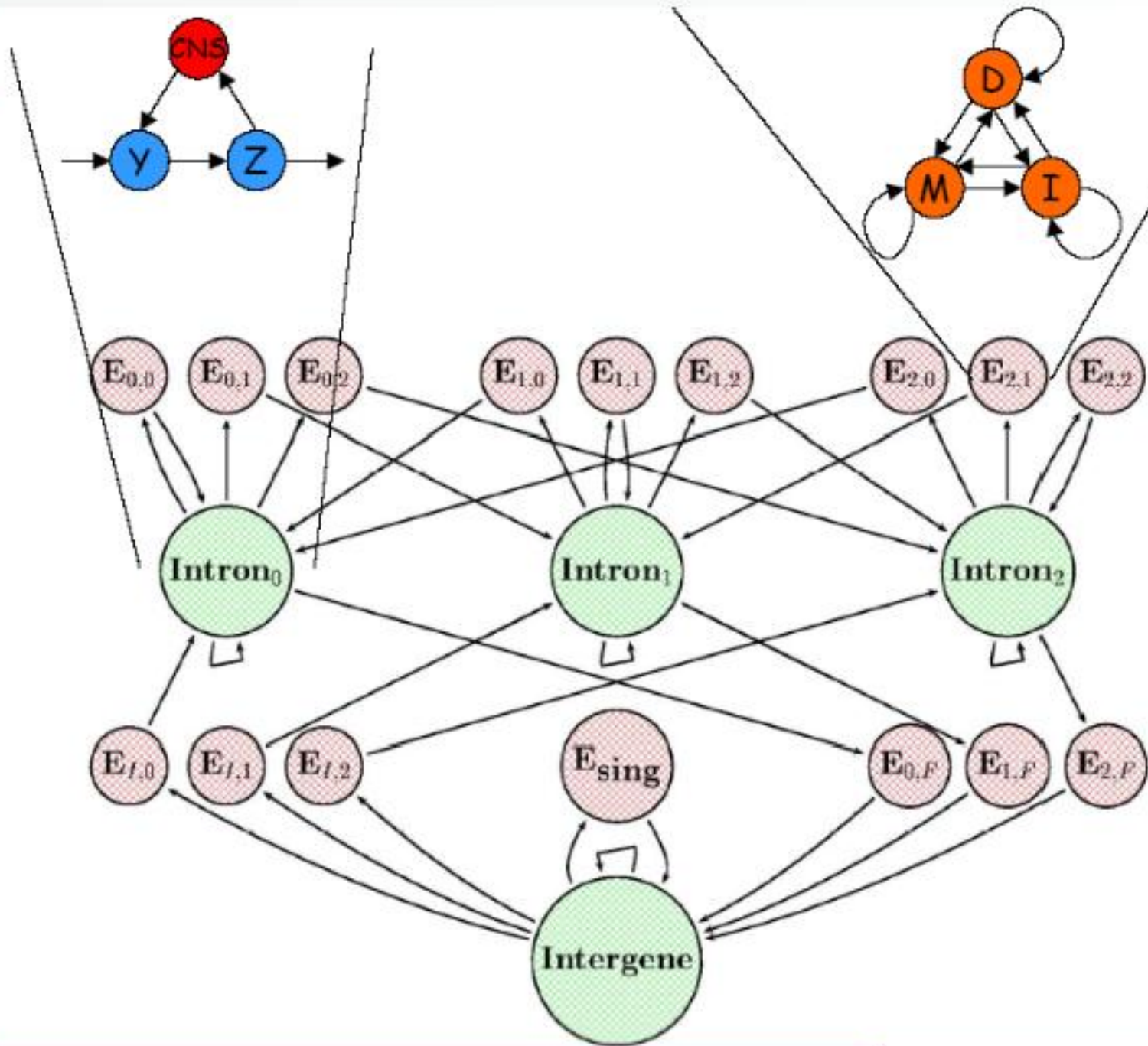
Observed sequence:

ATCGG  
ACGTCA



# SLAM

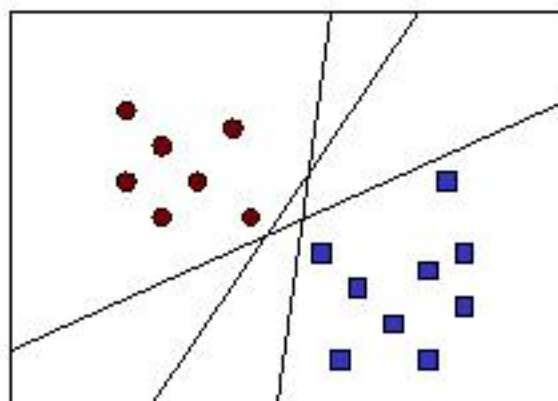
- SLAM components
  - Splice sites (Variable length Markov models).
  - Introns and Intergenic regions (2nd order Markov models, independent geometric lengths, CNS states).
  - Coding sequences (3-periodic Markov models, generalized length distributions, protein-based pairHMM.)
- Input
  - Pair of syntenic genomic sequences.
  - Approximate alignment.
- Output
  - CDS predictions in *both* sequences.



# Дискриминантный анализ.

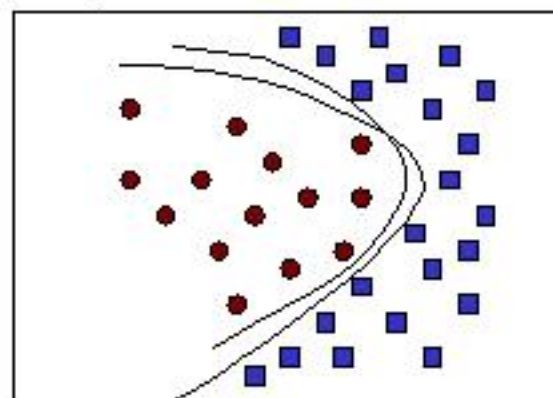
## Программа MZEF.

linear discriminator



$$Y = aX + b$$

quadratic discriminator



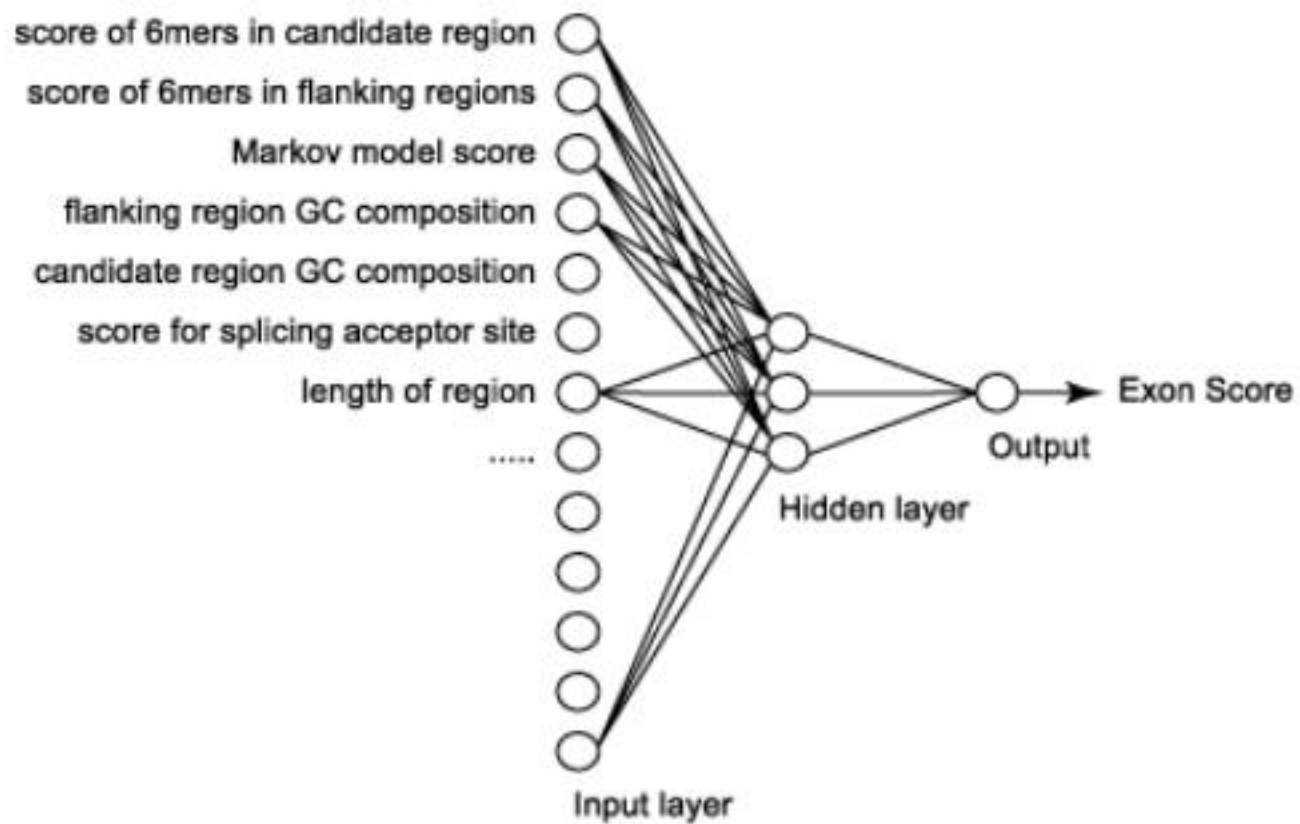
$$Y = aX^2 + bX + c$$



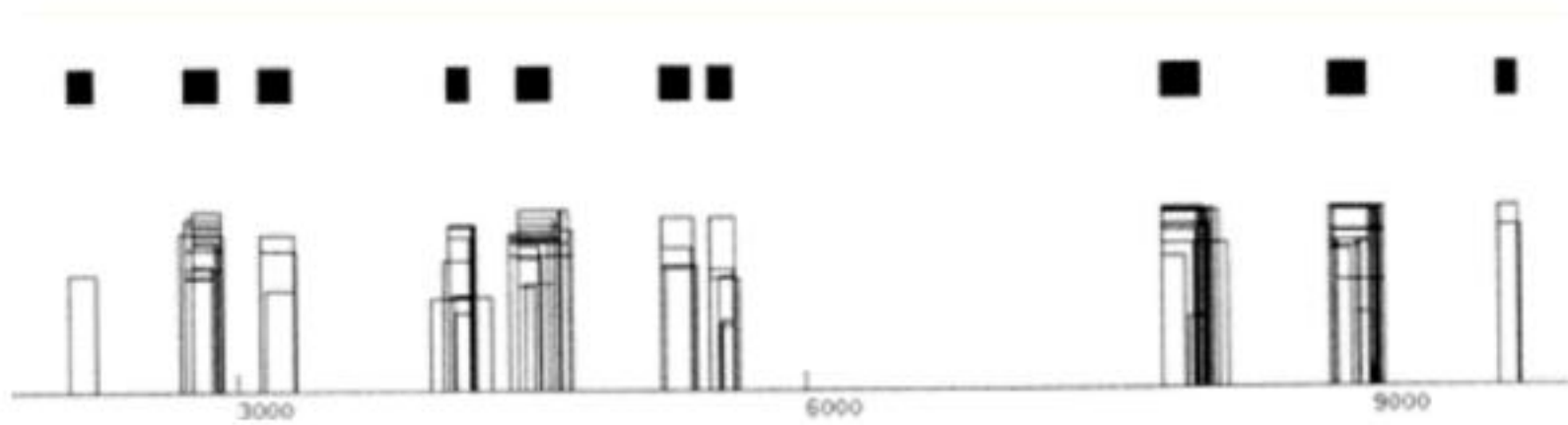
- ∇● Частоты олигонуклеотидов длины 6.
- Использование марковских цепей пятого порядка для трех рамок считывания, подобно GenMark.
- Учет распределения длин экзонов.
- Учет G+C состава 2kb прилежащих районов.
- Нейронные сети для распознавания сайтов сплайсинга.
- и т.д.



Feature vectors	Scores
$\langle f_{11}, f_{12}, f_{13}, \dots, f_{1n} \rangle$	score <sub>1</sub>
$\langle f_{21}, f_{22}, f_{23}, \dots, f_{2n} \rangle$	score <sub>2</sub>
...	
$\langle f_{m1}, f_{m2}, f_{m3}, \dots, f_{mn} \rangle$	score <sub>m</sub>

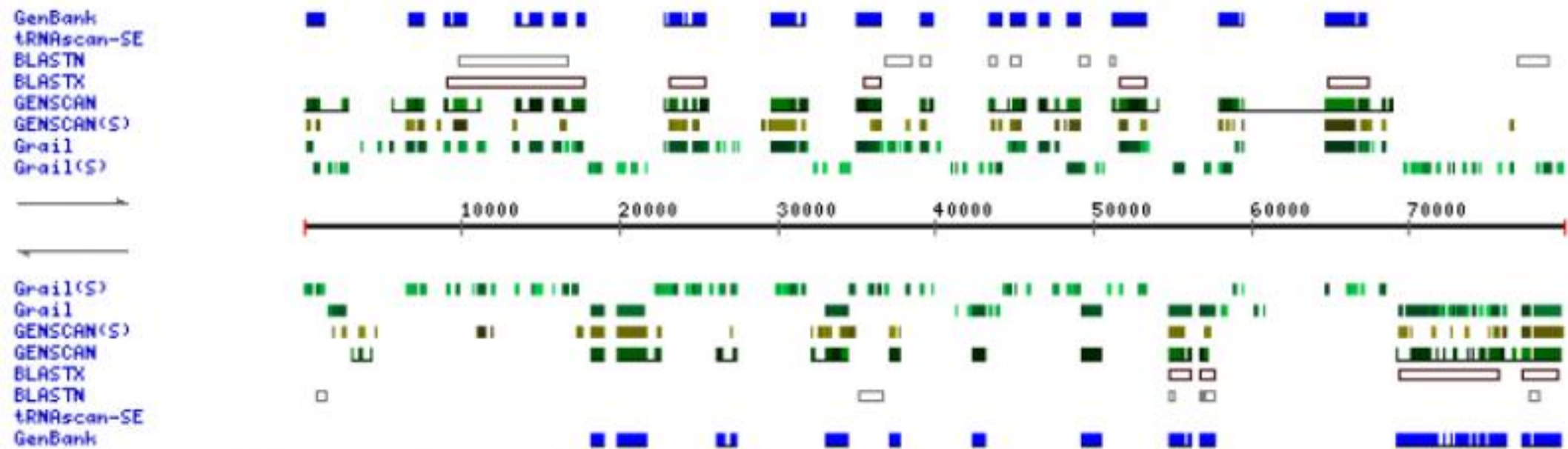






# Комбинированные методы распознавания генов

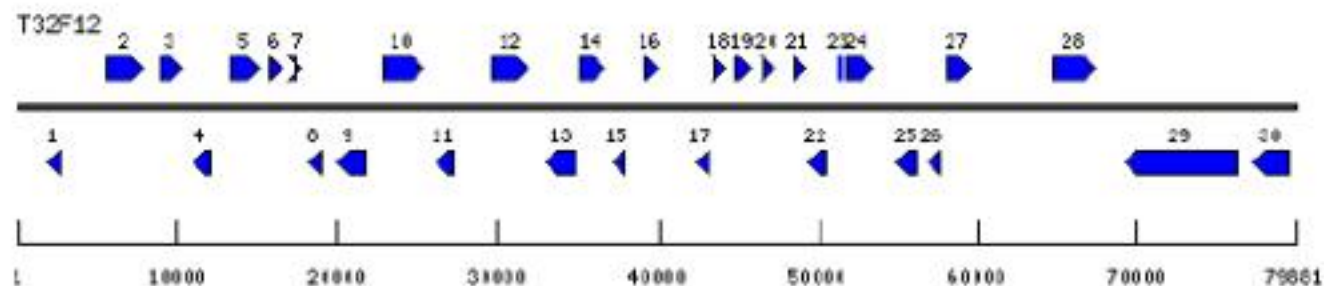
T32F12 (1..79881)



Different predictions from different programs.

TIGR "Combiner"

Final prediction is manually determined

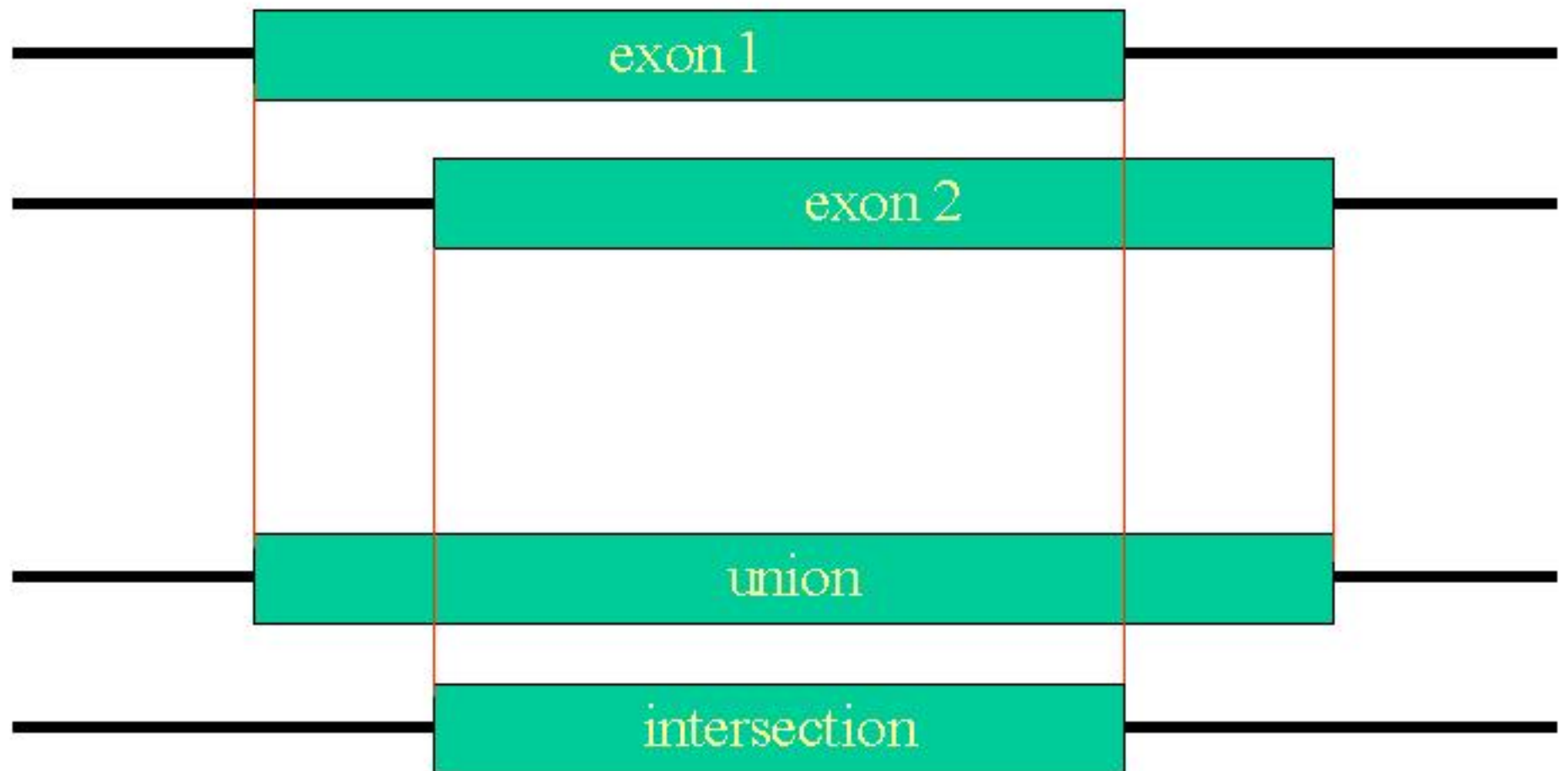


Genes on  
both strands

(this is one BAC clone drawn with the genes = 80,000 base pairs)

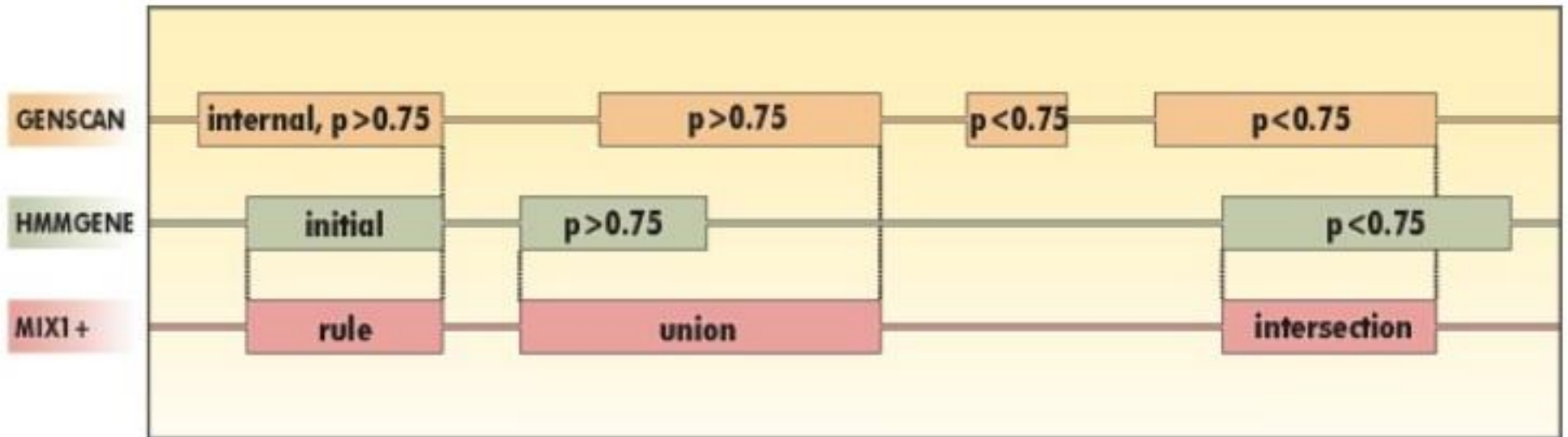


# AND и OR методы

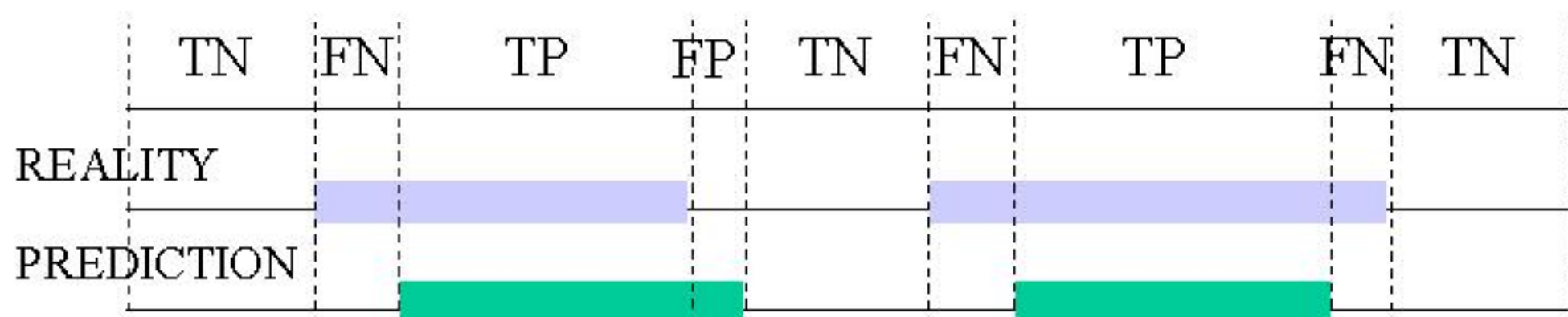




# EUI Метод



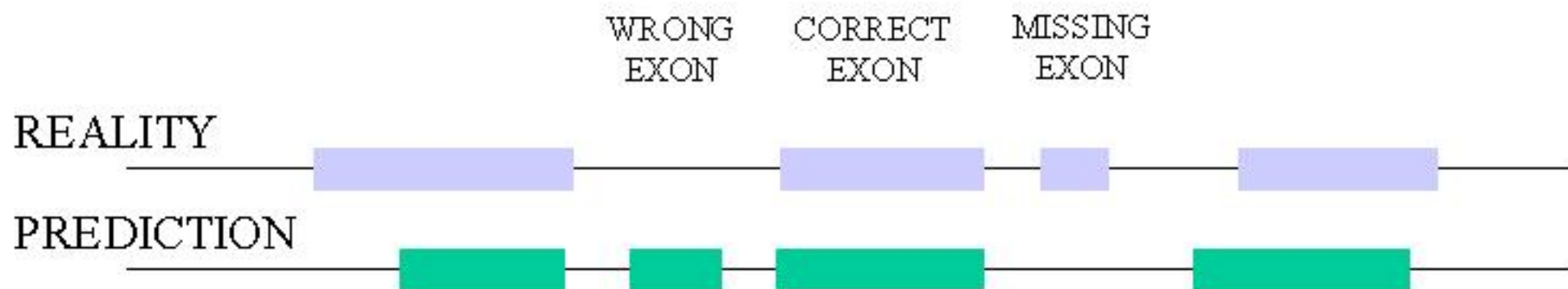
# Оценка качества работы программ предсказания на уровне нуклеотидов.



Sensitivity  $S_n = \frac{TP}{TP + FN}$   $\frac{\text{number of correct exons}}{\text{number of actual exons}}$

Specificity  $S_p = \frac{TP}{TP + FP}$   $\frac{\text{number of correct exons}}{\text{number of predicted exons}}$

# Оценка качества работы программ предсказания на уровне целых экзонов.



$$ESn = \frac{TE}{AE}$$

$$ESp = \frac{TE}{PE}$$

$$AC = \frac{1}{2} \left( \frac{TP}{TP + FN} + \frac{TP}{TP + FP} + \frac{TN}{TN + FP} + \frac{TN}{TN + FN} \right) - 1$$

$$CC = \frac{(TP * TN) - (FN * FP)}{((TP + FN) * (TN + FP) * (TP + FP) * (TN + FN))^{\frac{1}{2}}}$$

**Table 1. Nucleotide and Exon Level Accuracy**

Programs	No. of sequences	Nucleotide accuracy				Exon accuracy							
		Sn	Sp	AC	CC	ESn	ESp	(ESn+ESp)/2	ME	WE	PCa	PCp	OL
PCENES	195 (5)	0.86	0.88	0.84 ± 0.19	0.83	0.67	0.67	0.67 ± 0.32	0.12	0.09	0.20	0.17	0.02
GeneMark.hmm	195 (0)	0.87	0.89	0.84 ± 0.18	0.83	0.53	0.54	0.54 ± 0.36	0.13	0.11	0.29	0.27	0.09
Genie	195 (15)	0.91	0.90	0.89 ± 0.16	0.88	0.71	0.70	0.71 ± 0.30	0.19	0.11	0.15	0.15	0.02
GeneScan	195 (3)	0.95	0.90	0.91 ± 0.12	0.91	0.70	0.70	0.70 ± 0.32	0.08	0.09	0.21	0.19	0.02
HMMgene	195 (5)	0.93	0.93	0.91 ± 0.13	0.91	0.76	0.77	0.76 ± 0.30	0.12	0.07	0.14	0.14	0.02
Morgan	127 (0)	0.75	0.74	0.70 ± 0.21	0.69	0.46	0.41	0.43 ± 0.26	0.20	0.28	0.28	0.25	0.07
MZEF	119 (8)	0.70	0.73	0.68 ± 0.21	0.66	0.58	0.59	0.59 ± 0.28	0.32	0.23	0.08	0.16	0.01

For each sequence in the HMR195 dataset, the exons predicted on the forward (+) strand were compared to the annotated exons. The standard measures of predictive accuracy on nucleotide and exon level were calculated for each sequence and averaged over all sequences for which they were defined. This was done separately for each of the programs tested.

(No. of sequences) number of sequences effectively analyzed by each program; in parentheses is the number of sequences where the absence of gene was predicted; (Sn) nucleotide level sensitivity; (Sp) nucleotide level specificity; (AC) approximate correlation; (CC) correlation coefficient; (ESn) exon level sensitivity; (ESp) exon level specificity; (ME) missed exons; (WE) wrong exons; (PCa) proportion of real exons that were partially predicted (only one exon boundary correct); (PCp) proportion of predicted exons that were only partially correct; (OL) proportion of predicted exons that overlap an actual exon. AC and (ESn+ESp)/2 are given with standard deviation.



**Спасибо за внимание!**