

---

How does eukaryotic gene  
prediction work?

---

Age Tats  
JClub Feb 26th, 2008

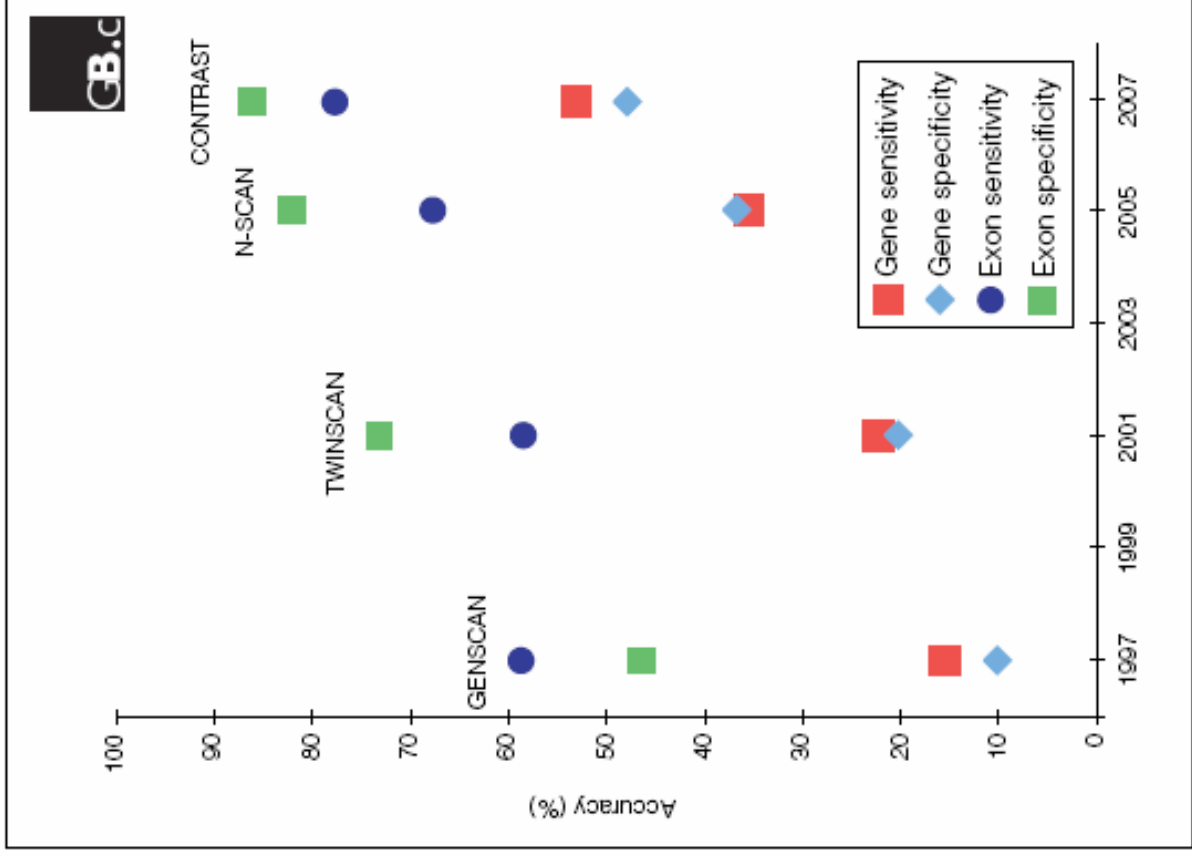
---

# Gene prediction

- complete exon-intron structures of the protein-encoding portions of transcripts (open reading frames)
  - 5' untranslated regions
  - only the boundaries of isolated exons
-

# Major approaches

	expression-based	<i>de novo</i>
input	genome sequence + cDNA sequences and/or their predicted translations ↓ cis alignment      ↓ trans alignment	the sequences of one or more genomes
pros	quite accurate	no need for additional data
cons	cannot predict genes expressed at low levels or under rare conditions (20-40%)	false-positives



# Generative (GHMM) versus discriminative models (CRF)

**Generative models**       $\Pr(\text{parse}|\text{seq}) = C \Pr(\text{parse})\Pr(\text{seq}|\text{parse})$

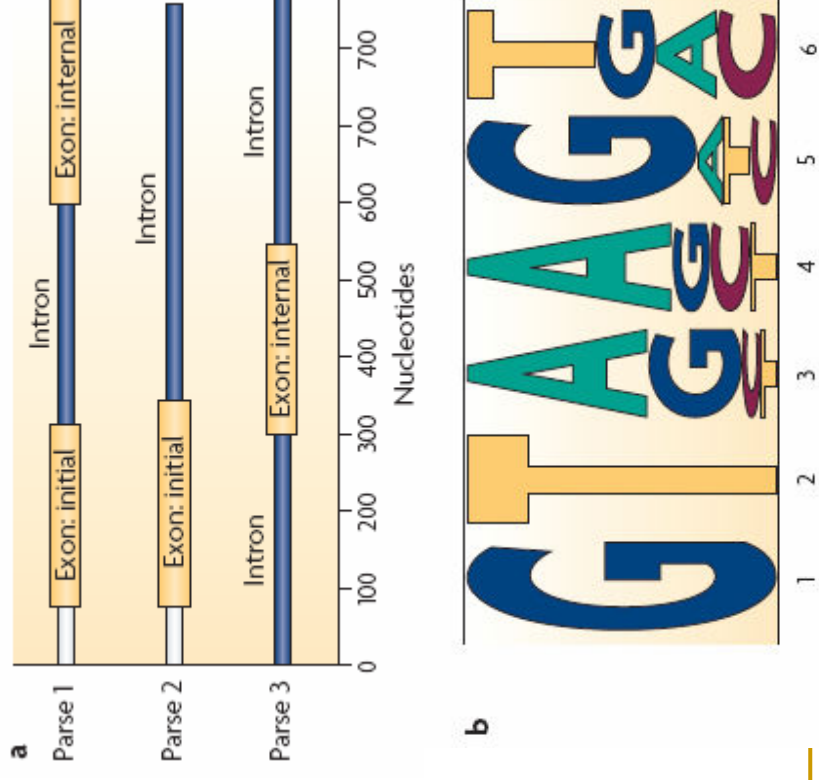
**$\Pr(\text{seq}|\text{parse 1}) =$**

$\Pr(\text{seq}(1,80)|\text{Intergenic}(1,80)) \times$

$\Pr(\text{seq}(80,320)|\text{Exon:initial}(80,320)) \times$

$\Pr(\text{seq}(320,600)|\text{Intron}(320,600)) \times$

$\Pr(\text{seq}(600,800)|\text{Exon:internal}(600,800))$



# Generative (GHMM) versus discriminative models (CRF)

**Generative models**       $\Pr(\text{parse}|\text{seq}) = C \Pr(\text{parse})\Pr(\text{seq}|\text{parse})$

**Discriminative models**       $\Pr(\text{parse}|\text{seq}) = \text{Ce}^{WF(\text{seq},\text{parse})}$

$WF(\text{seq},\text{parse})$  is a weighted sum of feature functions:

$$\sum_j W_j F_j(\text{seq}, \text{parse})$$

$F_j$  – predefined function, can depend on any part of the input sequence

$W_j$  – the weight of each function, learned from examples of correct parses during training

---

# CONTRAST: CONditionally TRAINED Search for

## Transcripts

Gross et al (2007) Genome Biology 8:R269

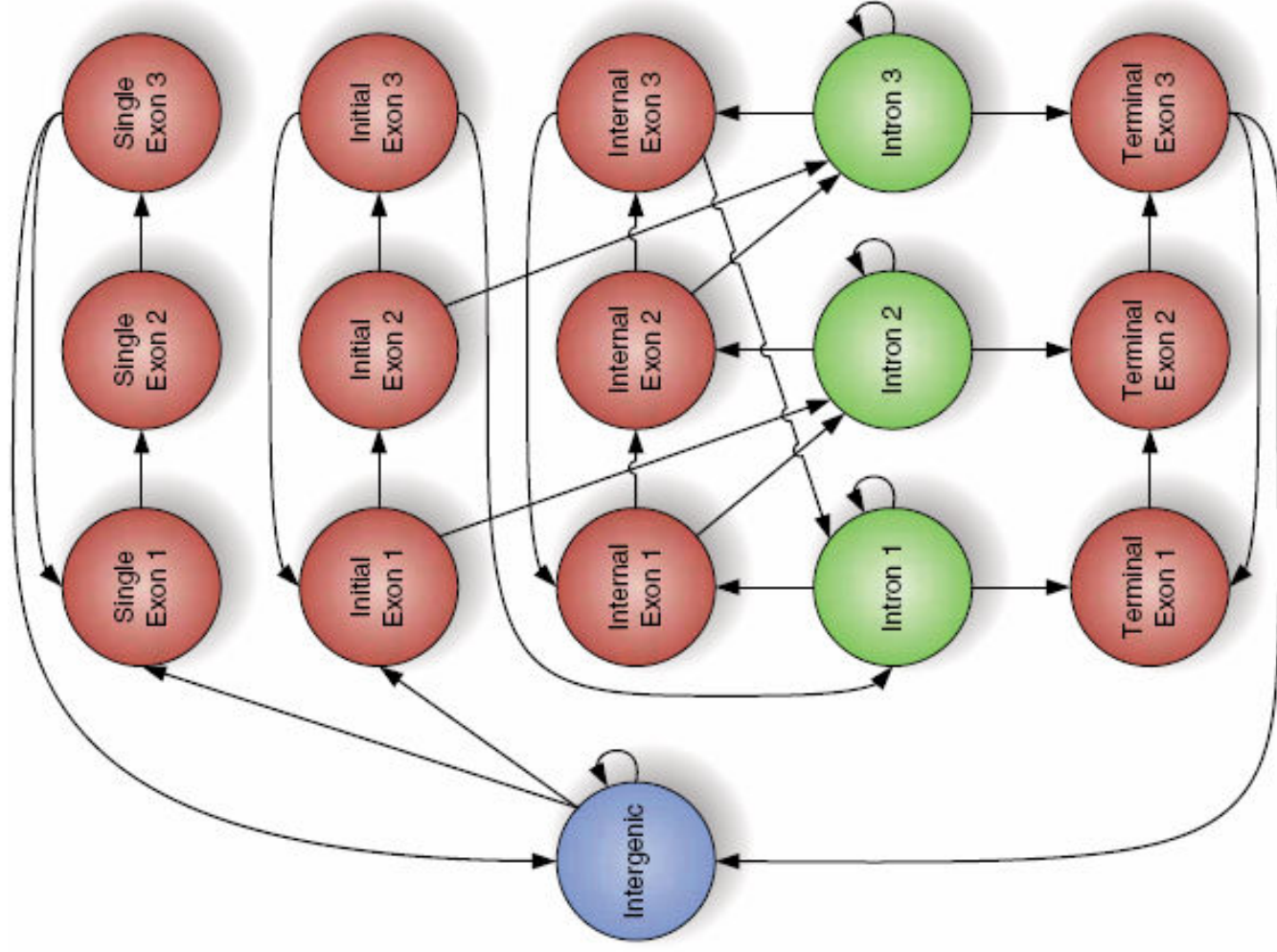
Directly incorporates information from multiple alignments rather than employing phylogenetic models.

### Uses

- SVMs for its coding boundary classifiers
  - CRF (conditional random field) for its global model of gene structure
-

human	ACAGGTGAGGAGGCG
macaque	.....
mouse	.....
rat	ACAGGTGAGAAAG..
rabbit	.....
dog	ACAGGTGAGGAGTCCG
cow	ACAGGTGAGCAGTCCG
armadillo	ACAGGTGAGGAG_CA
elephant	.....
tenrec	.....
opossum	CCAGGGAAG.....
chicken	CCAGGTGA.....
EST	SSSSIIIIIIIIII





# Boundaries of coding regions

Table 6

Coding region boundary classifiers. Coding region boundary classifiers. Window sizes and positions are shown for the five coding region boundary classifiers used by CONTRAST. Coordinates are defined such that the boundary occurs between the adjacent positions -1 and 1 (that is, either position -1 is coding and position 1 is coding or the reverse is true), with coordinates increasing in the 5' to 3' direction. Each classifier's require consensus sequence is shown in the second column.

	Consensus	5' end	3' end	Length
Start codon	A <sub>1</sub> T <sub>2</sub> G <sub>3</sub>	-8	6	14
Stop codon	T <sub>1</sub> A <sub>2</sub> A <sub>3</sub> , T <sub>1</sub> A <sub>2</sub> G <sub>3</sub> , T <sub>1</sub> G <sub>2</sub> A <sub>3</sub>	1	6	6
Donor splice GT	G <sub>1</sub> T <sub>2</sub>	-3	8	11
Donor splice GC	G <sub>1</sub> C <sub>2</sub>	-3	8	11
Acceptor splice	A <sub>2</sub> G <sub>1</sub>	-27	3	30

# Global model of gene structure

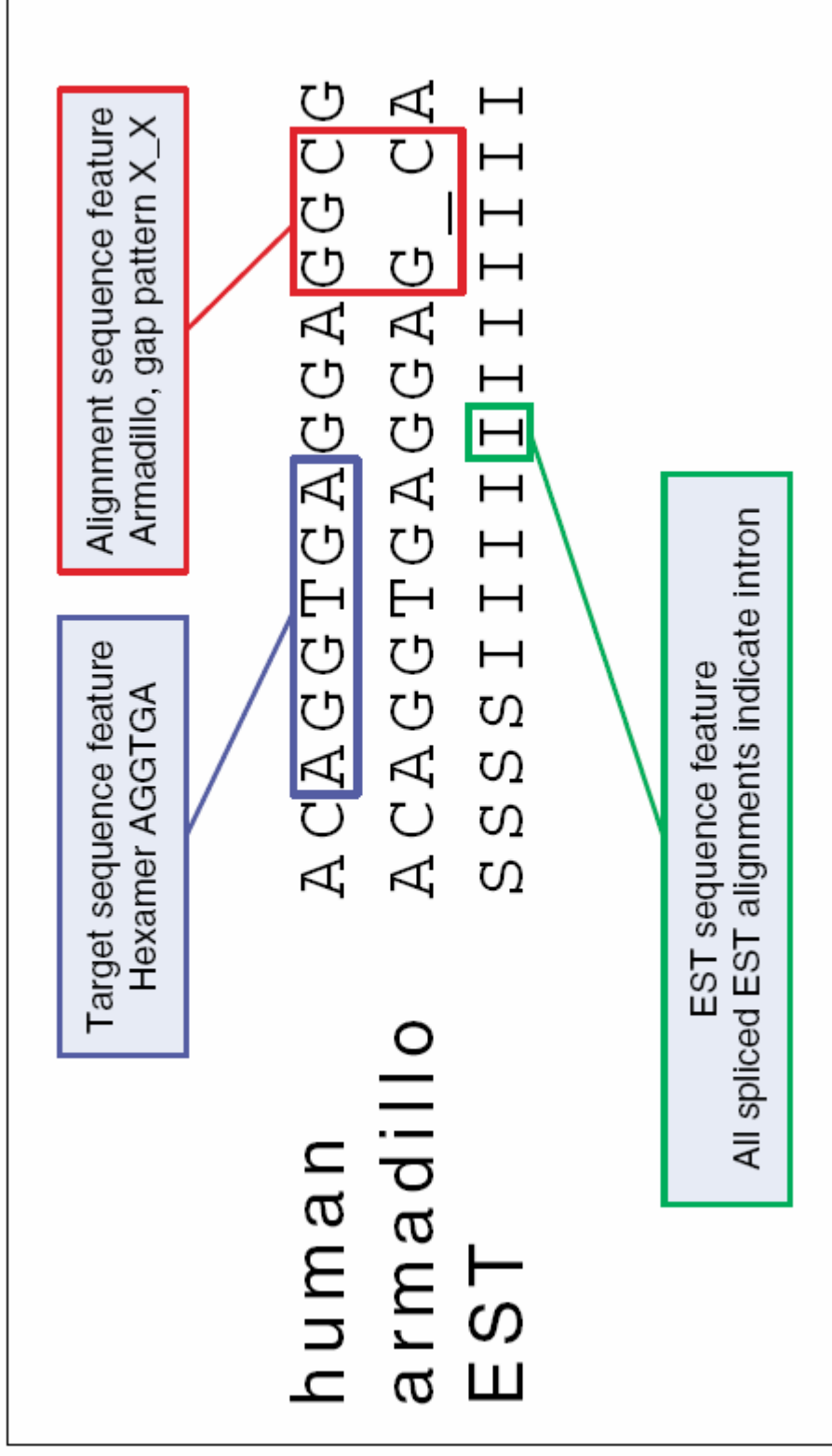
$$P(\mathbf{y} | \mathbf{x}) = \frac{e^{\mathbf{w}^T \mathbf{F}(\mathbf{x}, \mathbf{y})}}{\sum_{\mathbf{y}} e^{\mathbf{w}^T \mathbf{F}(\mathbf{x}, \mathbf{y})}}$$

$$\mathbf{F}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^L \mathbf{f}(y_{i-1}, y_i, i, \mathbf{x}).$$

$\mathbf{f}$  – feature mapping, a vector-valued function which determines the information used to calculate the score of a position.

$$\mathbf{f}(y_{i-1}, y_i, i, \mathbf{x}) = \begin{bmatrix} 1\{y_{i-1} = \text{Intron and } y_i = \text{Exon}\} \\ 1\{y_i = \text{Exon and } x_i = 'G'\} \\ \dots \end{bmatrix}$$

# Feature mapping in CONTRAST



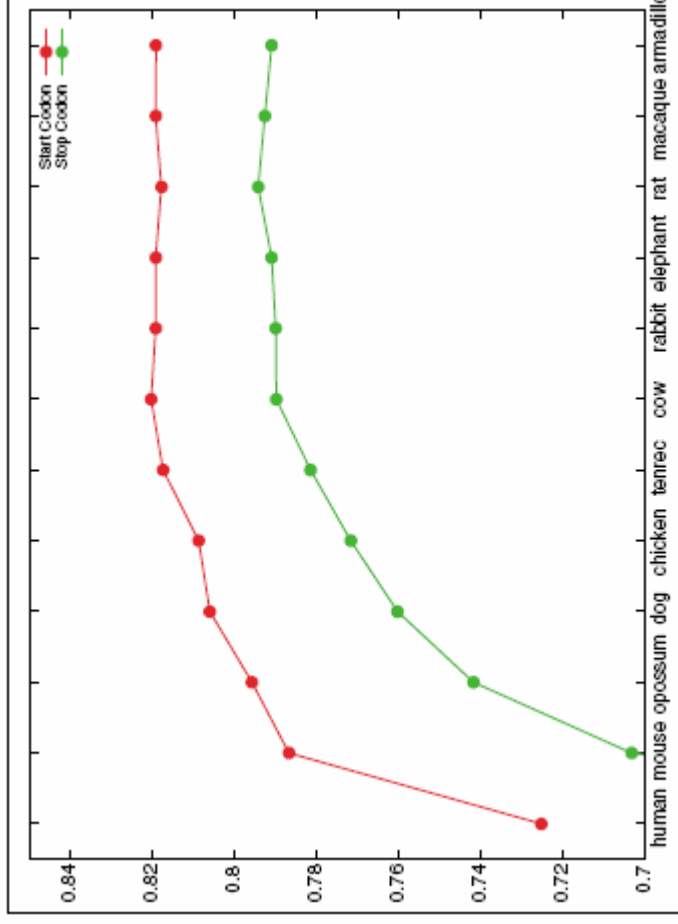
# Human gene prediction

Table 1

**De novo gene prediction performance for human. Sensitivity (Sn) and specificity (Sp) were evaluated at the gene, exon and nucleotide levels and reported as percentages. Also shown are the average number of genes and exons predicted for each cross-validation fold. The column headings indicate the predictor and informants used.**

	N-SCAN (mouse)	CONTRAST (mouse)	CONTRAST (11 informants)
Gene Sn	35.6	50.8	58.6
Gene Sp	25.1	29.3	35.5
Exon Sn	84.2	90.8	92.8
Exon Sp	64.6	70.5	72.5
Nucleotide Sn	90.8	96.0	96.9
Nucleotide Sp	67.9	70.0	72.0
Genes predicted	22,596	27,614	26,260
Exons predicted	196,643	211,431	210,180

## Start and stop codon classifier accuracy



## Splice site classifier accuracy

