

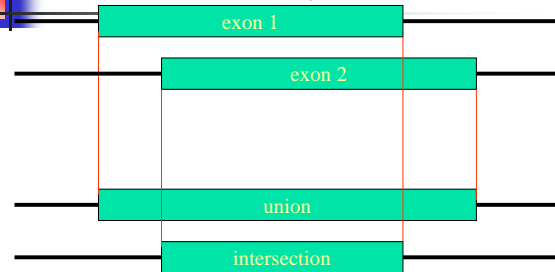
## Integrated approaches for gene finding

- Programs that integrate results of similarity searches with *ab initio* techniques (GenomeScan, FGENESH+, Procrustes)
- Programs that use synteny between organisms (ROSETTA, SLAM)
- Integration of programs predicting different elements of a gene (EuGene)
- Combining predictions from several gene finding programs (combination of experts)

(c) Devika Subramanian, 2006

122

## AND and OR Methods



(c) Devika Subramanian, 2006

123

## Combining Genscan and HMMgene

- High prediction accuracy as well as reliability of their exon probability make them good candidates.



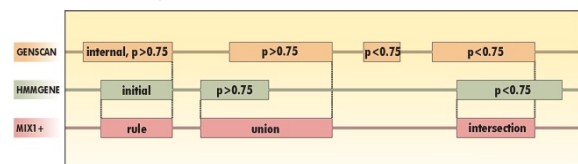
- Genscan predicted 77% of exons correctly, HMMgene 75%, both 87%

(c) Devika Subramanian, 2006

124

## EUI Method (exon union - intersection)

- Union of exons with  $p \geq 0.75$
- Intersection of exons with  $p < 0.75$
- Rule for initial exon

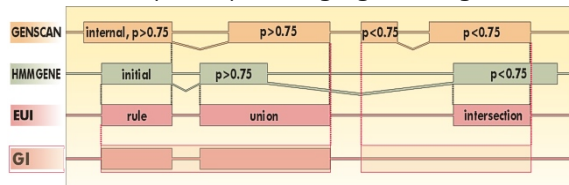


(c) Devika Subramanian, 2006

125

## Gene intersection (GI) method

1. Intersection of genes
2. Apply EUI method to exons completely belonging to GI genes



(c) Devika Subramanian, 2006

126

## EUI with reading frame consistency

1. Assign probabilities to GI genes. Determine position of acceptor and donor site in a reading frame.
2. GI gene with higher probability imposes the reading frame. Choose only EUI exons contained in GI genes that are in a chosen reading frame.

(c) Devika Subramanian, 2006

127

## Results - Burset/Guigó dataset

METHODS	#no prediction	Nucleotide accuracy			Exon accuracy				
		Sn	Sp	AC	ESn	ESp	$(ESn + ESp)^2$	ME	P/E
Genescan	8	0.94	0.93	0.92	0.78	0.81	0.80	0.09 (203)	0.05 (188)
HMMgene	38	0.93	0.94	0.92	0.81	0.83	0.82	0.14 (308)	0.04 (139)
EUI	20	0.94	<b>0.96</b>	<b>0.93</b>	<b>0.83</b>	<b>0.88</b>	<b>0.85</b>	0.12 (250)	<b>0.03</b> (95)
GI	43	0.91	<b>0.97</b>	<b>0.93</b>	<b>0.82</b>	<b>0.90</b>	<b>0.86</b>	0.18 (386)	<b>0.02</b> (67)
EUI_frame	27	0.93	<b>0.96</b>	<b>0.93</b>	<b>0.83</b>	<b>0.88</b>	<b>0.85</b>	0.13 (286)	<b>0.03</b> (87)

(c) Devika Subramanian, 2006

128

## Summary: Eukaryotic gene finding

- Overall accuracy usually below 50%
  - Human gene finding is hardest
  - Very long introns, and lots of them
- Leading methods: HMMs and variants
- New ideas needed
- New opportunity: use sequence of related species

(c) Devika Subramanian, 2006

129

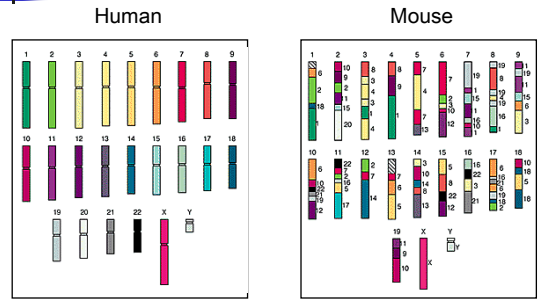
## Comparison of 1196 orthologous genes

- Sequence identity between genes in human/mouse
  - exons: 84.6%
  - protein: 85.4%
  - introns: 35%
  - 5' UTRs: 67%
  - 3' UTRs: 69%
- 27 proteins were 100% identical.

(c) Devika Subramanian, 2006

130

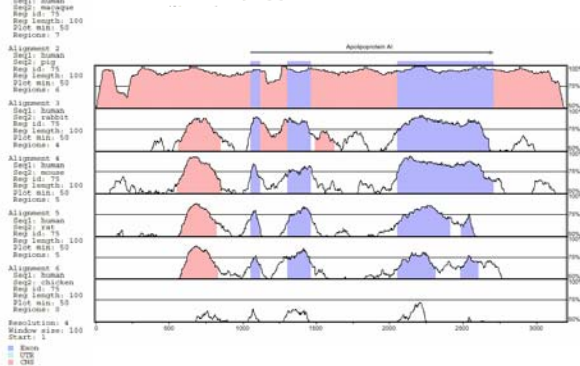
## Human-mouse homology



(c) Devika Subramanian, 2006

131

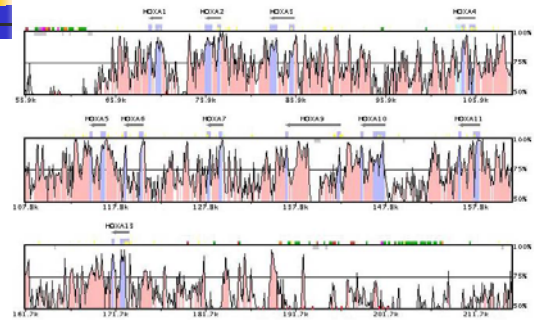
human vs macaque, pig, rabbit, mouse, rat, chicken



(c) Devika Subramanian, 2006

132

## HoxA human-mouse



(c) Devika Subramanian, 2006

133

## Alignment

```

50      . : . : . : . : . : . : . : . :
247 GGTGAGGTCGAGGACCCCTGCA  CCGAGCTGTATGGAGGGCA  AGAGC
   | : | | | | : | | | | --:| | | | : : | | | |---| | | |
368 GAGTCGGGGGAGGGGGCTGCTGTGGCTCTGGACAGCTTGCATTGAGAGG

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

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

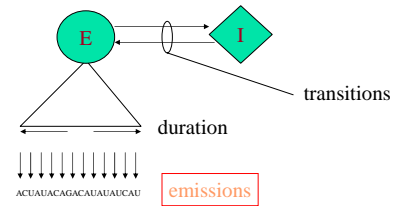
200     . : . : . : . : . : . : . : . :
367 TTCAACCTCAGCAATGCCATCATGGGCAGCGGCATCCTGGGACTCGCCTA
   ||| | : | | | | | | | | | | | | | | | | | | | | | | | | |
517 TTCAATCTCAGCAACGCCATCATGGGCAGTCGGAATTCTGGGGCTCGCCTA
  
```

(c) Devika Subramanian, 2006

134

## Twinscan

- Twinscan is an augmented version of the Genscan HMM.



(c) Devika Subramanian, 2006

135

## Twinscan Algorithm

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

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

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

(c) Devika Subramanian, 2006

136

## Twinscan Algorithm

- Run Viterbi using emissions  $e_j(k)$  where  $k \in \{ A-, A:, A|, \dots, T| \}$

Emission distributions  $e_j(k)$  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

(c) Devika Subramanian, 2006

137

## Example

Human: **ACGGCGACUGUGCACGU**  
 Mouse: **ACUGUGAC GUGCACUU**  
 Alignment: **||:|:| |-|||:|**

Input to Twinscan HMM:

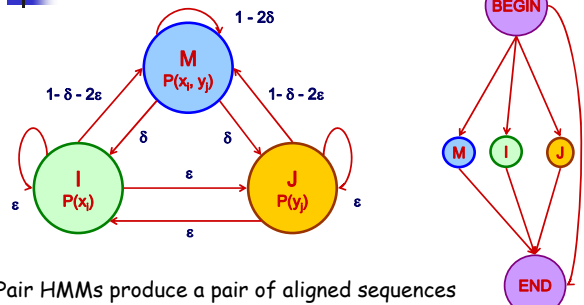
A| C| G: G| C: G| A| C| U- G| U| G| C| A| C| G: U|

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

(c) Devika Subramanian, 2006

138

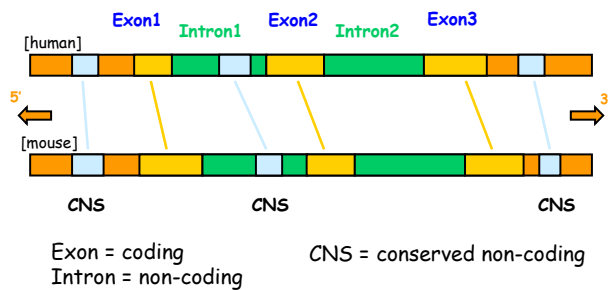
## A Pair HMM for alignments



(c) Devika Subramanian, 2006

139

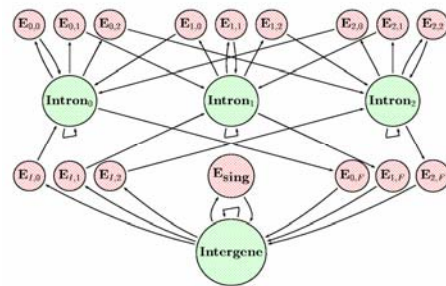
## Cross-species gene finding



(c) Devika Subramanian, 2006

140

## Generalized Pair HMMs



(c) Devika Subramanian, 2006

141

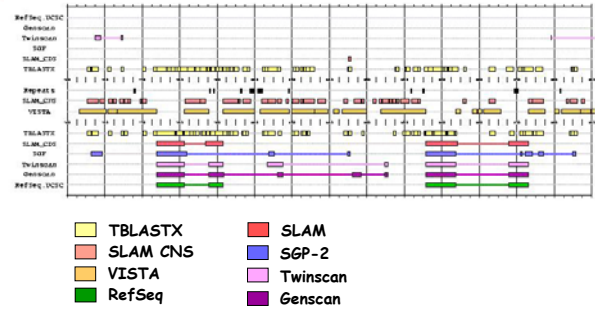
## Ingredients in exon scores

- Splice site detection (VLMM)
- Length distribution (generalized)
- Coding potential (codon freq. tables)
- Isochore group

(c) Devika Subramanian, 2006

142

## Example: HoxA2 and HoxA3



(c) Devika Subramanian, 2006

143

## What have we learned from comparative gene finding?

- conservation is a stronger splice site indicator than consensus
- intron lengths have diverged
- gene structure conservation is more powerful than sequence conservation for prediction
- consensus for GC splice sites

(c) Devika Subramanian, 2006

144

## Measuring Performance

Test set	Nucleotide level			Exon level				
	SN	SP	AC	SN	SP	(SN+SP)/2	ME	WE
<b>The ROSETTA set</b>								
ROSETTA	0.935	0.978	0.949	0.833	0.829	0.831	0.048	0.047
SGP-1	0.940	0.960	0.940	0.700	0.760	0.730	0.120	0.040
SLAM	0.951	0.981	0.960	0.783	0.755	0.769	0.038	0.057
TWILSCAN.p	0.960	0.941	0.940	0.855	0.824	0.840	0.045	0.081
TWILSCAN	0.984	0.889	0.923	0.839	0.767	0.803	0.034	0.118
GENSCAN	0.975	0.908	0.929	0.817	0.770	0.793	0.057	0.107
<b>HoxA</b>								
SLAM	0.852	0.896	0.864	0.727	0.533	0.630	0.000	0.333
TWILSCAN.p	0.976	0.829	0.896	0.773	0.531	0.652	0.000	0.312
TWILSCAN	0.949	0.511	0.704	0.591	0.173	0.382	0.000	0.707
SGP-2	0.640	0.637	0.619	0.409	0.173	0.291	0.091	0.596
GENSCAN	0.932	0.687	0.796	0.545	0.235	0.390	0.000	0.569
<b>Elastin</b>								
SLAM	0.876	0.981	0.926	0.802	0.859	0.831	0.121	0.059
TWILSCAN.p	0.942	0.950	0.945	0.879	0.889	0.884	0.066	0.036
TWILSCAN	0.933	0.877	0.903	0.835	0.826	0.831	0.110	0.120
SGP-2	0.755	0.998	0.873	0.533	0.900	0.720	0.352	0.017
GENSCAN	0.947	0.766	0.852	0.835	0.731	0.783	0.121	0.231

## Priority organisms

### Human-mouse gene finding not very high-impact

- lots of ancillary data gives better evidence
- most genes now known
- nonetheless, this problem is getting all the attention

### Countless other species really need gene finders:

- Brugia malayi (causes lymphatic filariasis)
- Toxoplasma gondii
- Schistosoma mansoni (Schistosomiasis)
- Entamoeba histolytica (50 million cases/year)
- Tetrahymena thermophila (model organism)
- Plants: potato, maize, sorghum
- Mammals: chimp, dog, cow, pig

From the TIGR web site.

(c) Devika Subramanian, 2006

146

## Genome scale gene finding

Strategy	Based on	Examples
Ab initio prediction	Models of gene structure/comp	Genscan, GRAIL GenLang, hmngene
Microarray	Hybridization	Exon-scanning array
Gene inference	Homology	GenomeScan
Genomic:genomic alignment	Homology	ExoFish GLASS/Rosetta
DNA:protein alignment	Homology	GeneWise
cDNA sequencing	Sequencing	RIKEN

C. Burge Nature Genet. 27, 5-7, 2001

(c) Devika Subramanian, 2006

147