

Current Bioinformatics Projects

Sören Sonnenburg[†], Mikio Braun,[†] Petra Philipps[‡], Cheng Soon Ong,[‡]
Alex Zien,[‡] Gunnar Rätsch[‡], Klaus-Robert Müller[†]



Fraunhofer
Institut
Rechnerarchitektur
und Softwaretechnik



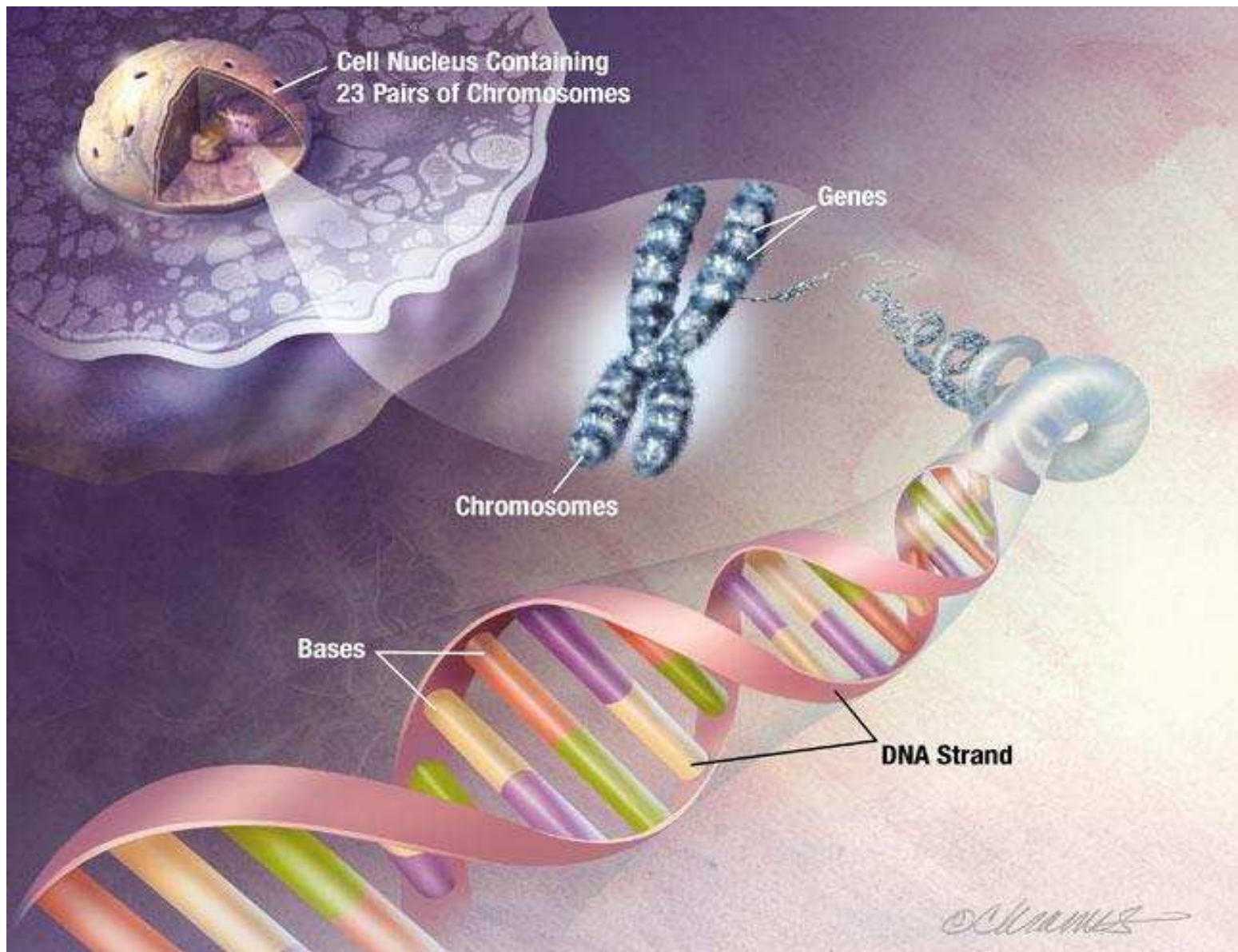
MAX-PLANCK-GESELLSCHAFT

- [†] Fraunhofer FIRST.IDA, Kekuléstr. 7, 12489 Berlin, Germany
- [‡] Friedrich Miescher Laboratory of the Max Planck Society,
Spemannstr. 37-39, 72076 Tübingen, Germany

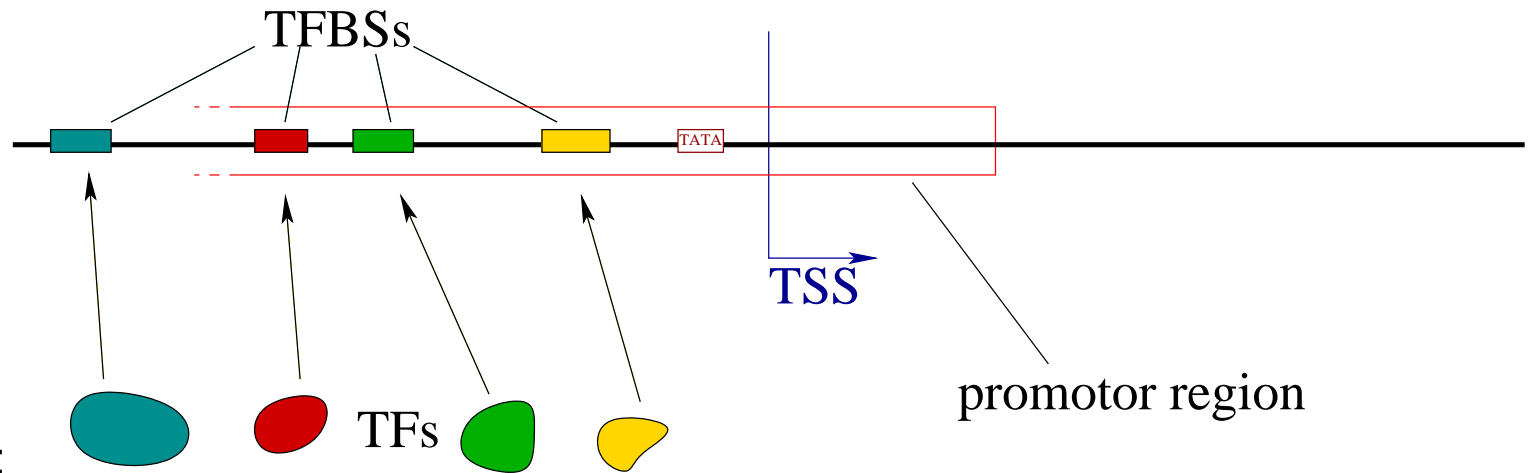
PROJECT OVERVIEW:

- Promoter Prediction
- Splice Form Prediction
- Alternative Splicing
- TFBS Module Discovery
- siRNA Prediction
- PolyA Prediction
- Genome wide SNP discovery on *A. Thaliana*
- Secondary Structure Prediction

PROMOTER PREDICTION



DEFINITION



Properties:

- TSS has no exact location, more like a range of $[-20, +20]$ base pairs
- TSS - no consensus sequence
- position, order and number of TFBS in Promotor region variable

⇒ **Promotor Prediction is non-trivial**

FEATURES

- TFBS in Promotor region
- condition: DNA should not be too twisted
- CpG islands (often over TSS/first exon, seem to be 2 general types CpG and non-CpG island promoters)
- TSS with TATA box (≈ -30 bp upstream)
- exon content in UTR 5' region

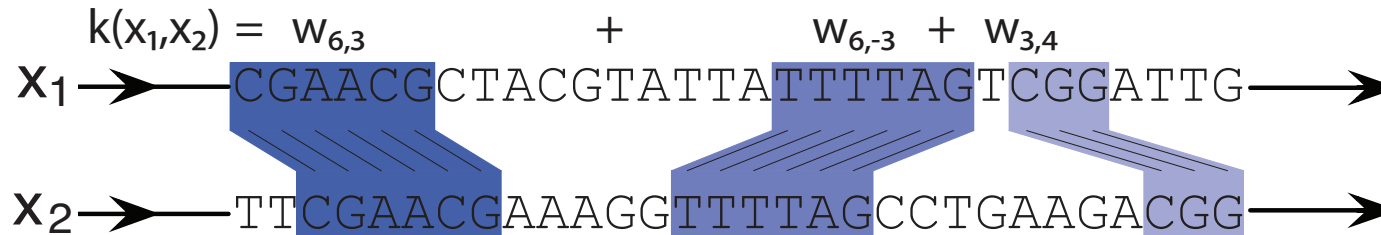
Idea: Combine weak features to build strong promotor predictor

METHOD

Use Support Vector Machine Classifier classifier

$$f(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^N y_i \alpha_i \mathbf{k}(\mathbf{x}, \mathbf{x}_i) + b \right)$$

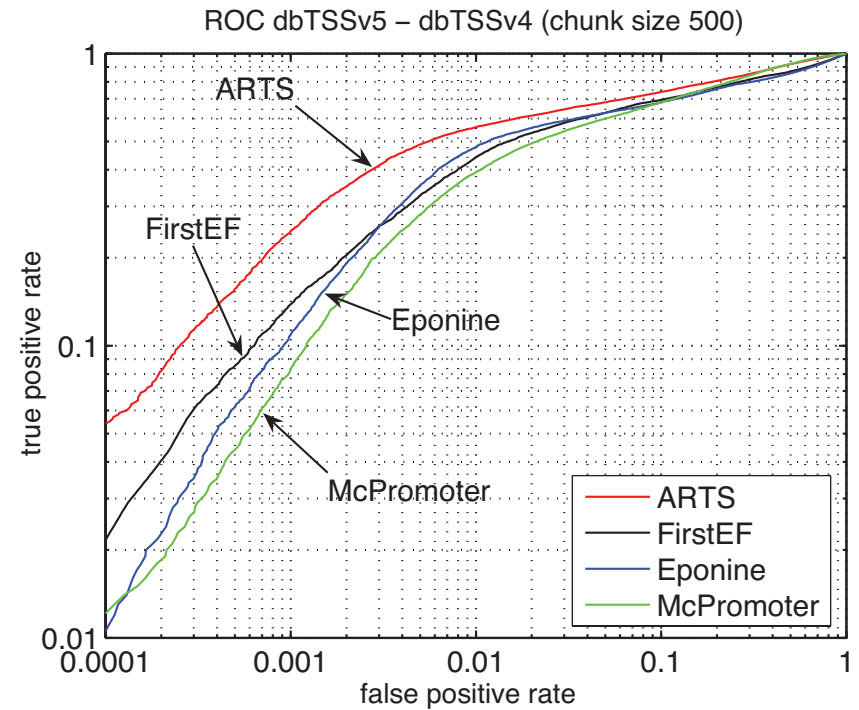
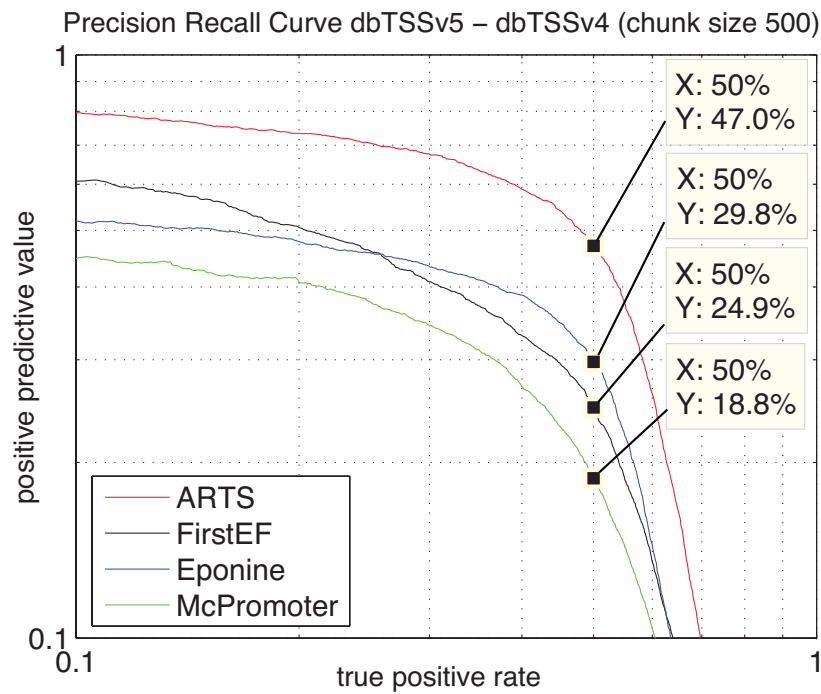
- key ingredient is kernel $k(\mathbf{x}, \mathbf{x}')$ – gives means to compare 2 sequences
- to detect TSS (including parts of core promotor with TATA box) – use Weighted Degree Shift kernel



- CpG Islands, distant Enhancers and TFBS upstream of TSS, coding sequence downstream (Spectrum kernel)
- stacking energy of DNA (linear kernel)

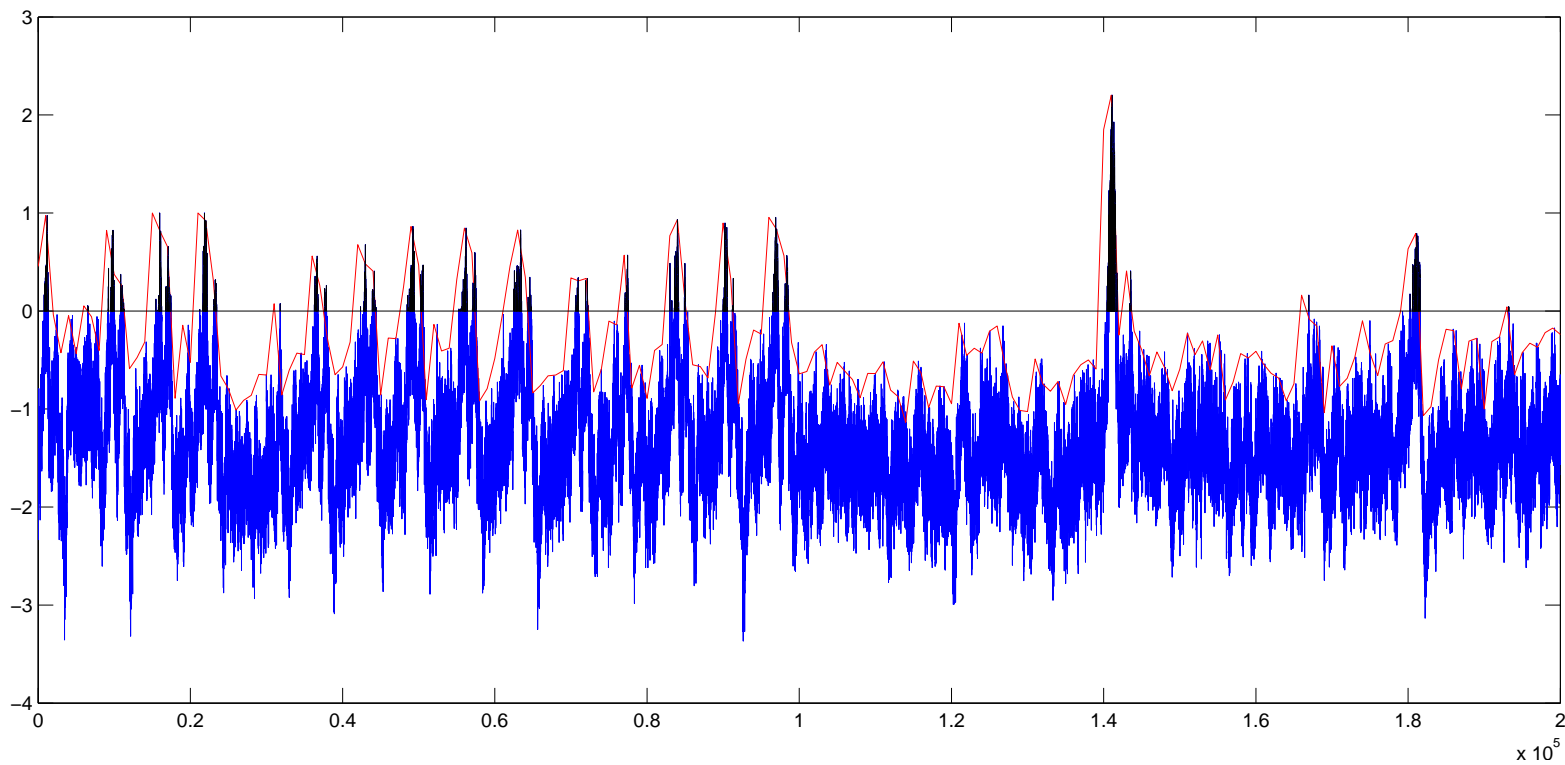
RESULTS

Genomewide (human) evaluation



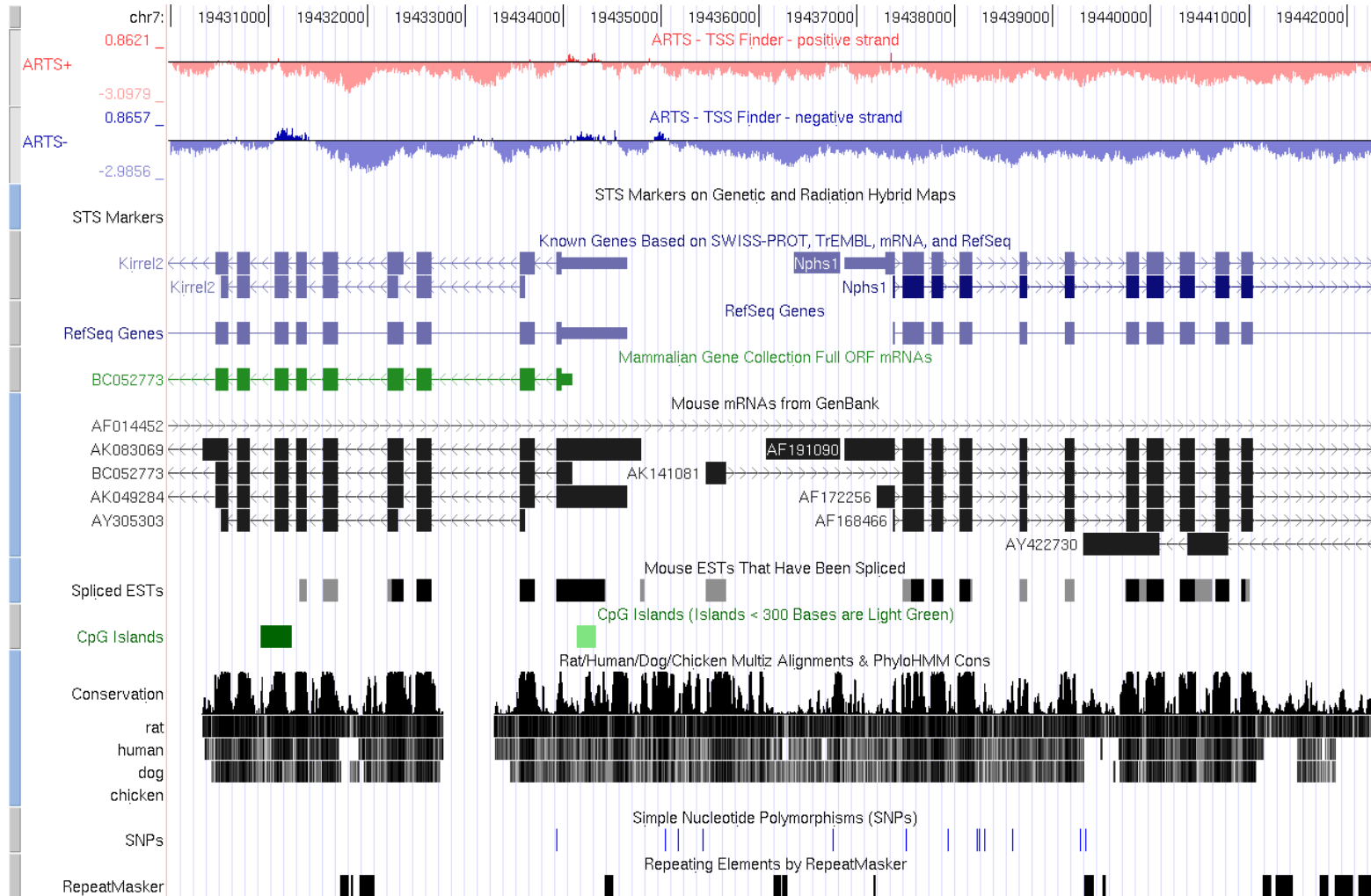
EXAMPLE

Protocadherin Alpha

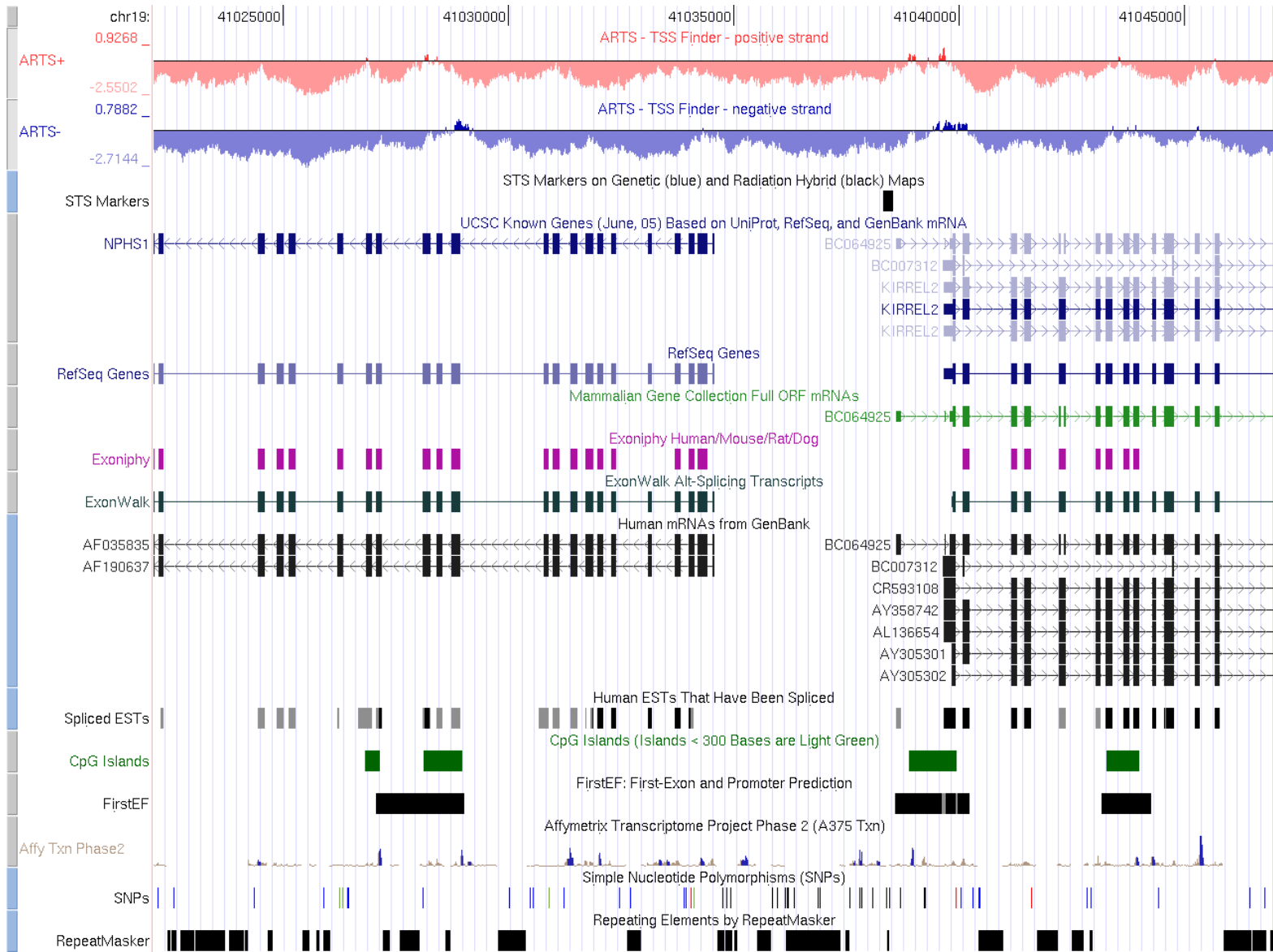


- ⇒ might even detect alternative promoters
- ⇒ accepted at ISMB 2006

NEPHRIN MM5

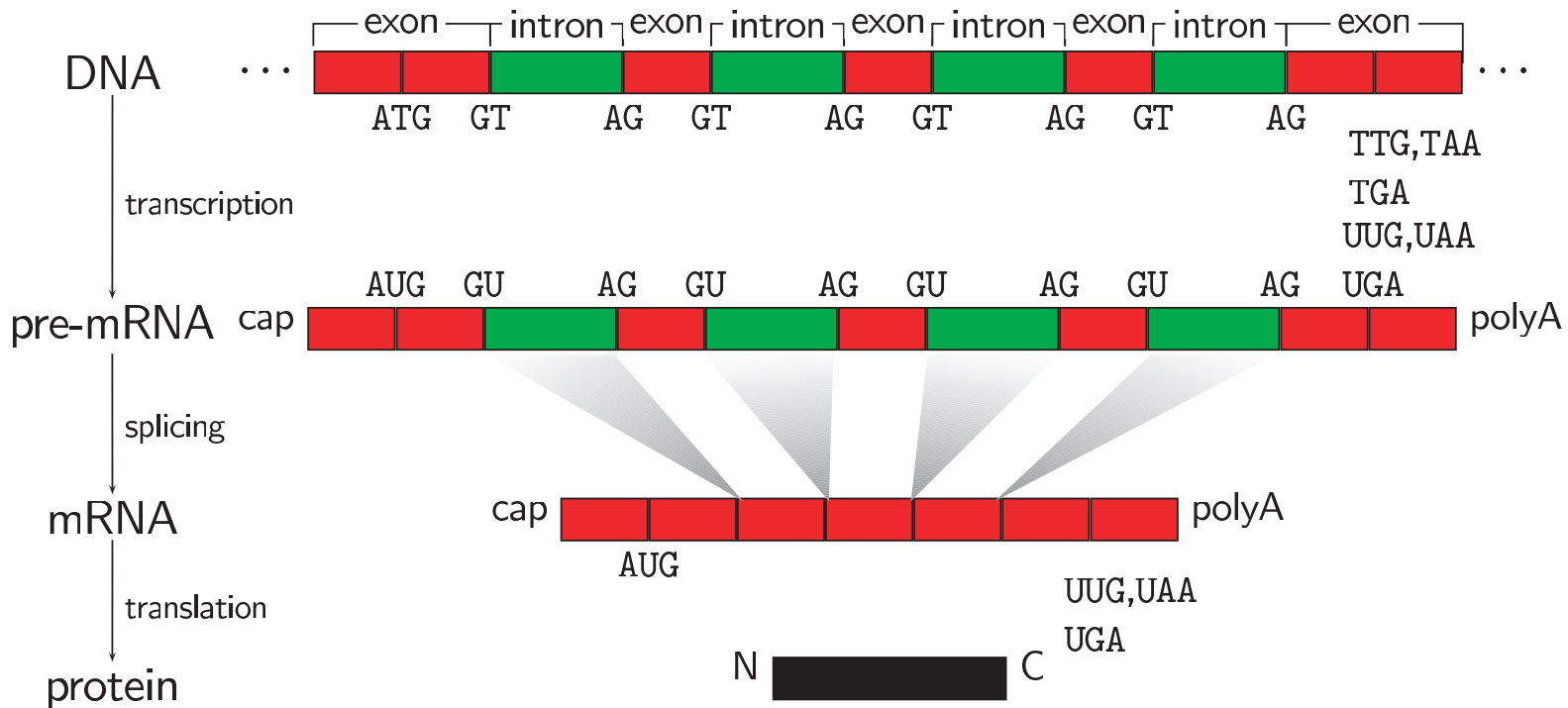


NEPHRIN HG17

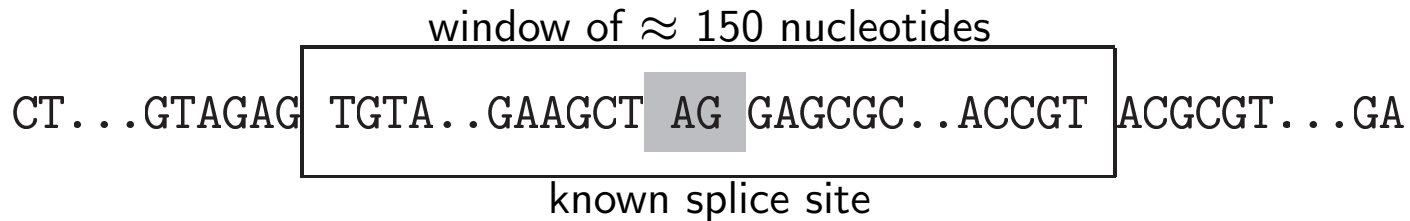




SPLICE FORM PREDICTION



2-CLASS SPLICE SITE DETECTION

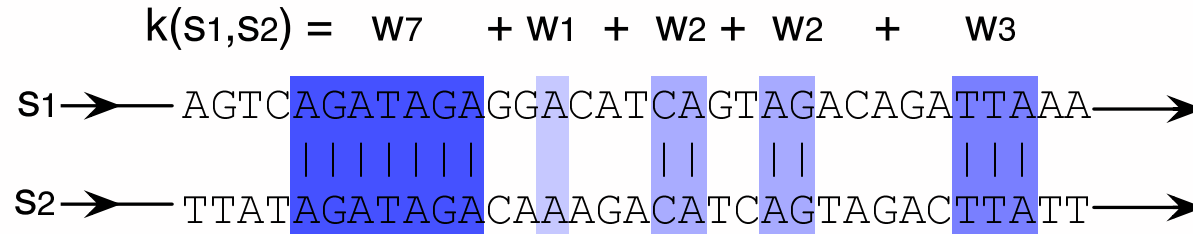


- **true sites:** fixed window around a true splice site
- **decoys sites:** generated by shifting the window

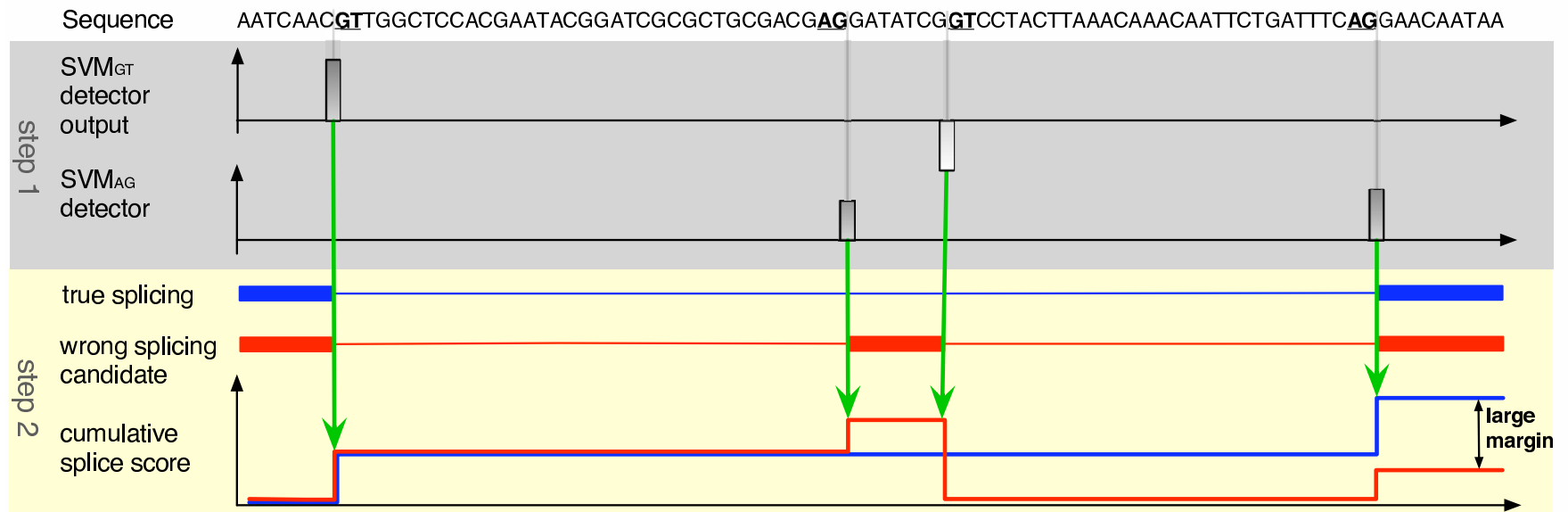
```

AAACAAATAAGTAACTAATCTTTTAGGAAGAACGTTTCAACCATTTTGAG
AAGATTAACAAAAACAAATTTTAGCATTACAGATATAATAATCTAATT
CACTCCCCAAATCAACGATATTTTAGTTCACTAACACATCCGTCTGTGCC
TTAATTTCACTTCCACATACTTCCAGATCATCAATCTCCAAAACCAACAC
    
```

Design of new Support Vector kernel that allows appropriate processing:

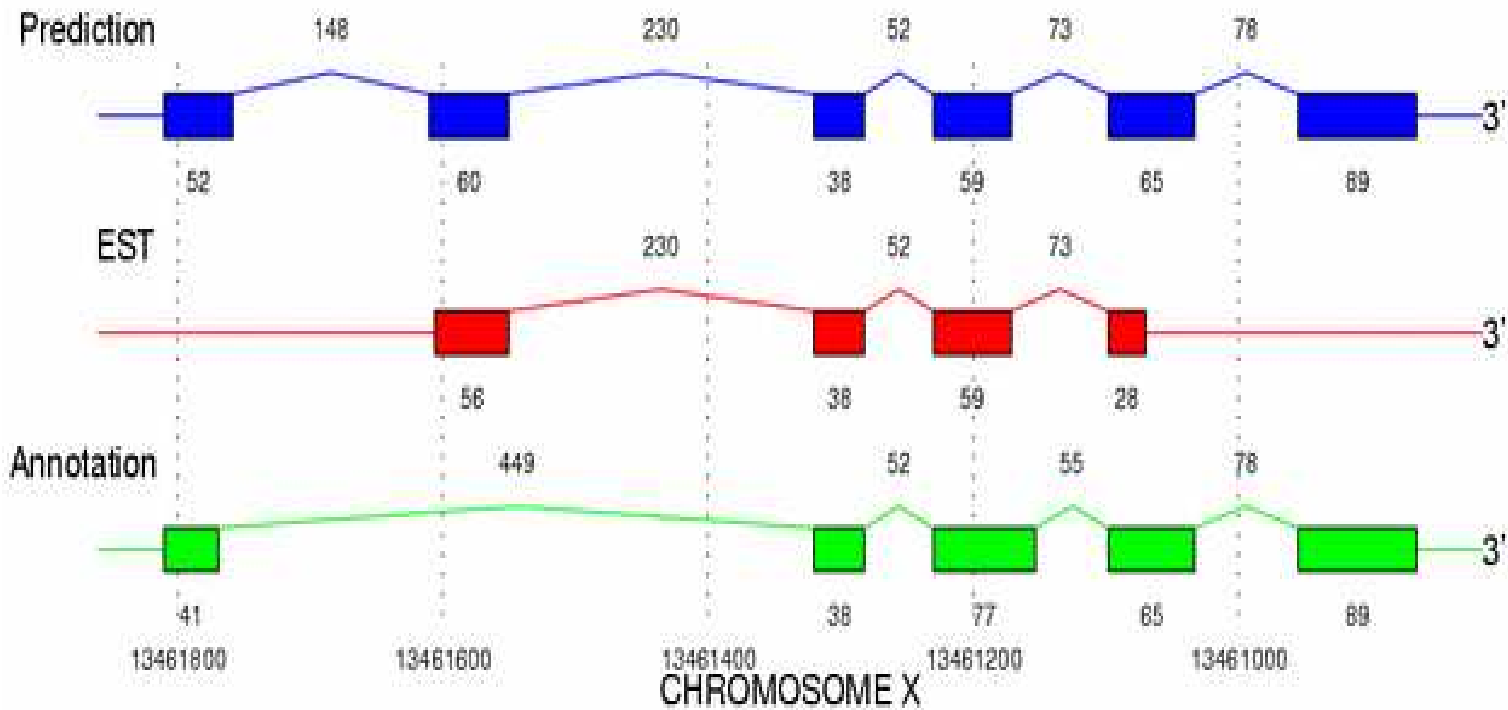


LARGE MARGIN SPLICE FORM PRED.



$$\begin{aligned}
 S := & S_{L_{E,f}}(p_1^{\text{GT}} - p_s) + S_E(s_{[p_s, p_1^{\text{GT}}]}) + S_{L_{E,l}}(p_e - p_n^{\text{AG}}) + S_E(s_{[p_n^{\text{AG}}, p_e]}) \\
 & + \sum_{i=1}^n \left[S_{L_I}(p_i^{\text{AG}} - p_i^{\text{GT}}) + S_I(s_{[p_i^{\text{GT}}, p_i^{\text{AG}}]}) + S_{\text{AG}}(p_i^{\text{AG}}) + S_{\text{GT}}(p_i^{\text{GT}}) \right] \\
 & + \sum_{i=1}^{n-1} \left[S_E(s_{[p_i^{\text{AG}}, p_{i+1}^{\text{GT}}]}) + S_{L_E}(p_i^{\text{AG}} - p_{i+1}^{\text{GT}}) \right]
 \end{aligned}$$

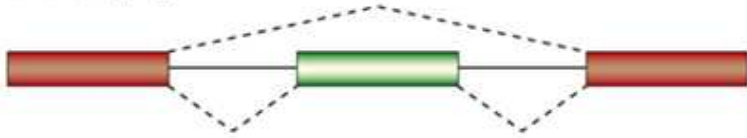
RESULTS



⇒ submitted to PLoS computational biology

ALTERNATIVE SPLICING

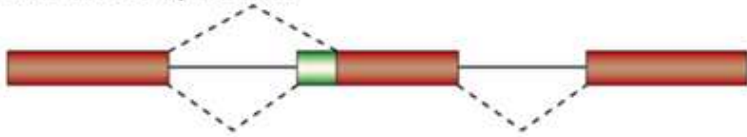
Exon skipping



Alternative 5' splice sites



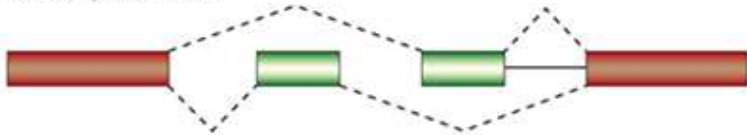
Alternative 3' splice sites



Intron retention



Mutually exclusive

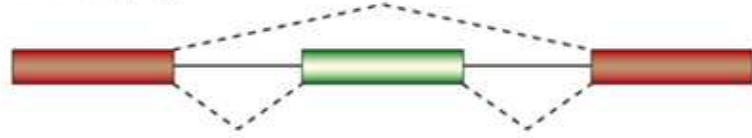


Idea: Use *Machine Learning* to

- analyze sequences near splice sites
- understand differences between alternative and constitutive splicing
- exploit and identify regulative splicing elements
- predict yet unknown alternative splicing events

FORMS OF ALTERNATIVE SPLICING

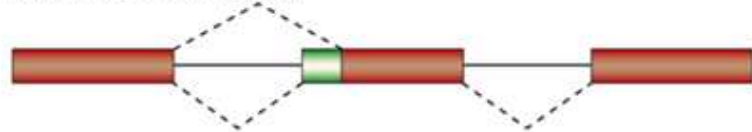
Exon skipping



Alternative 5' splice sites



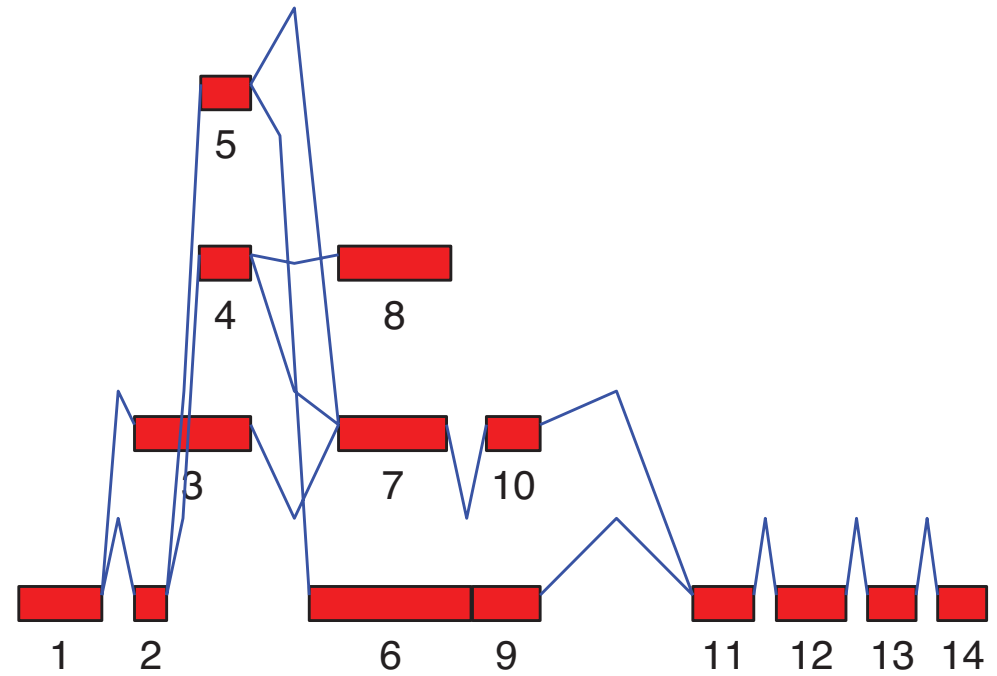
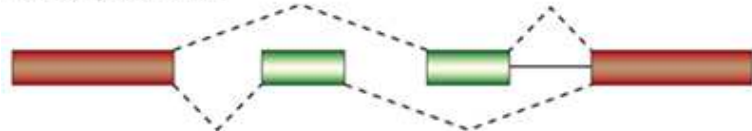
Alternative 3' splice sites



Intron retention



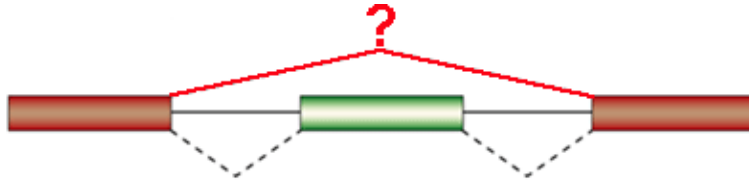
Mutually exclusive



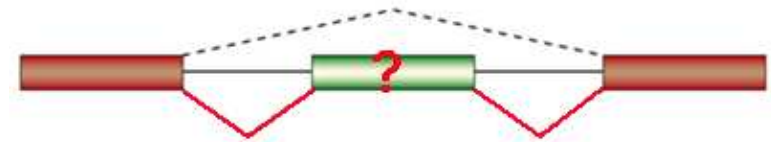
(*C. elegans* gene T08B2.5)

PREDICTION OF ALT. SPLICED EXONS

- Two different Tasks:
 - Exon is known
 - Can it be skipped?

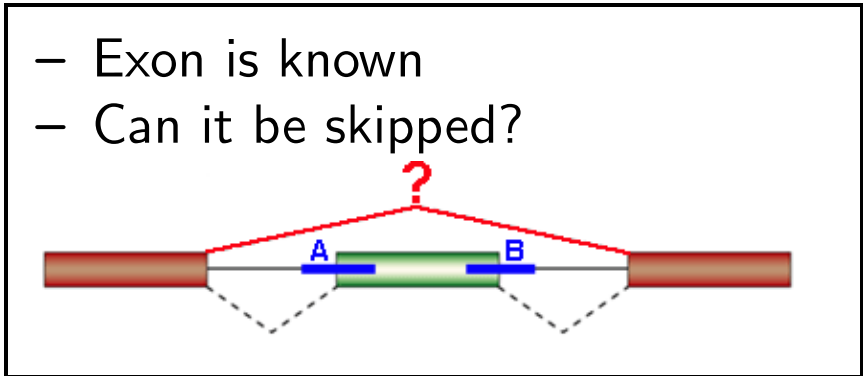


- Intron is known
- Does it contain an exon?

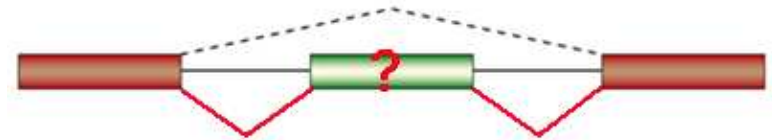


PREDICTION OF ALT. SPLICED EXONS

- Two different Tasks:



- Intron is known
- Does it contain an exon?

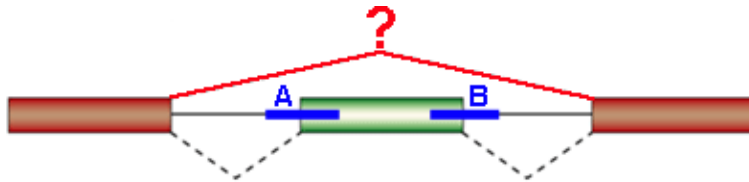


Solution to Task 1

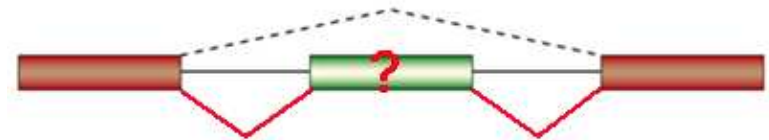
- Two-class Classification Problem
- ⇒ Use Support Vector Machines (SVMs) on
 - * sequences A & B (± 100 nt of splice sites)
 - * exon & intron lengths

PREDICTION OF ALT. SPLICED EXONS

- Two different Tasks:
 - Exon is known
 - Can it be skipped?



- Intron is known
- Does it contain an exon?

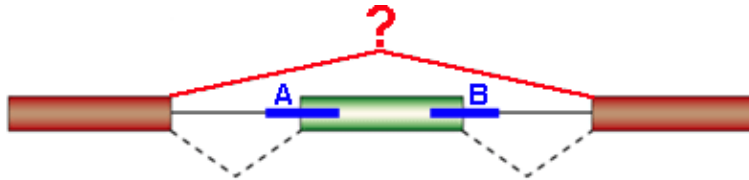


Solution to Task 2

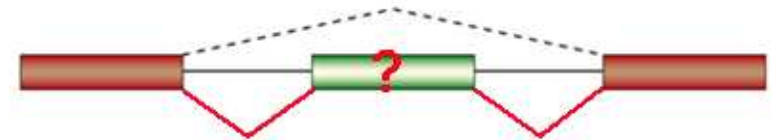
- *Problem:* We do not know yet the exon boundaries!
- *Solution:* Consider all possible exons within the intron.

PREDICTION OF ALT. SPLICED EXONS

- Two different Tasks:
 - Exon is known
 - Can it be skipped?



- Intron is known
- Does it contain an exon?



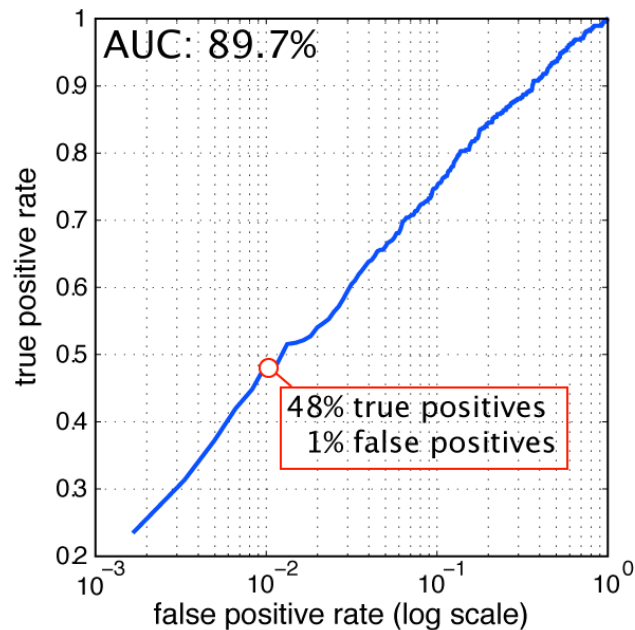
Solution to Task 2

- *Problem:* We do not know yet the exon boundaries!
- *Solution:* Consider all possible exons within the intron.
- Classify true exons vs. wrong exons
- ⇒ Use SVM-like algorithm using the
 - * sequences A & B (± 100 nt of splice sites)
 - * exon & intron lengths and
 - * *splice site scores* (SVM based)

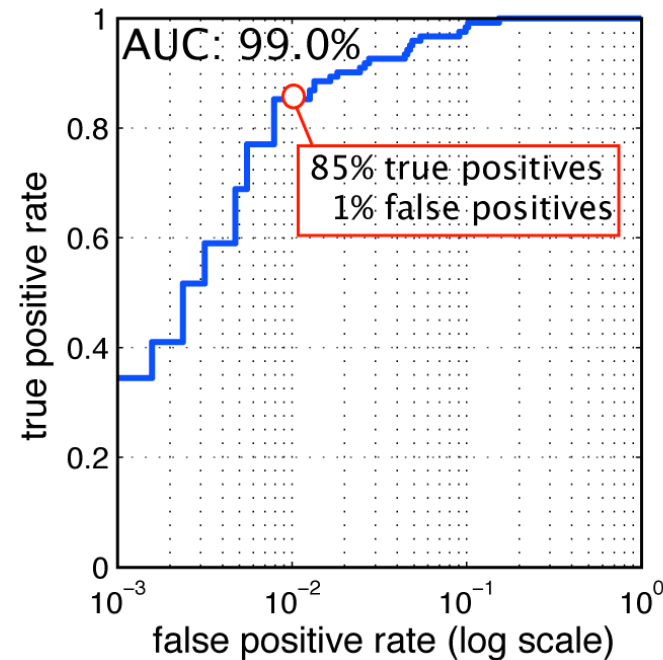
COMPUTATIONAL RESULTS

- 487 alternatively and 2531 constitutively spliced exons

Task 1 (exon known)



Task 2 (intron known)



As accurate as for conserved alternatively spliced exons in human.

⇒ Conservation not needed!

⇒ presented at ISMB 2005

TFBS MODULE DISCOVERY

project just started. . .

- given *genome wide* binding sites locations for MEF2 binding genes
- task1: find binding sites that form modules with MEF2 associated BS
- task2: genes are active in different development stages (early,late,always)
- which are the discriminating TFBS ?

FUTURE

- Genefinding
- Alternative Promoters
- Quantitative Predictions of Alternative Splicing
- Mass Spectrography
- . . .

Questions ?