

## **Conditional Random Fields**

### **a ich použitie na anotáciu biologických sekvencií**

Broňa Brejová

Katedra informatiky FMFI UK

## Conditional Random Fields

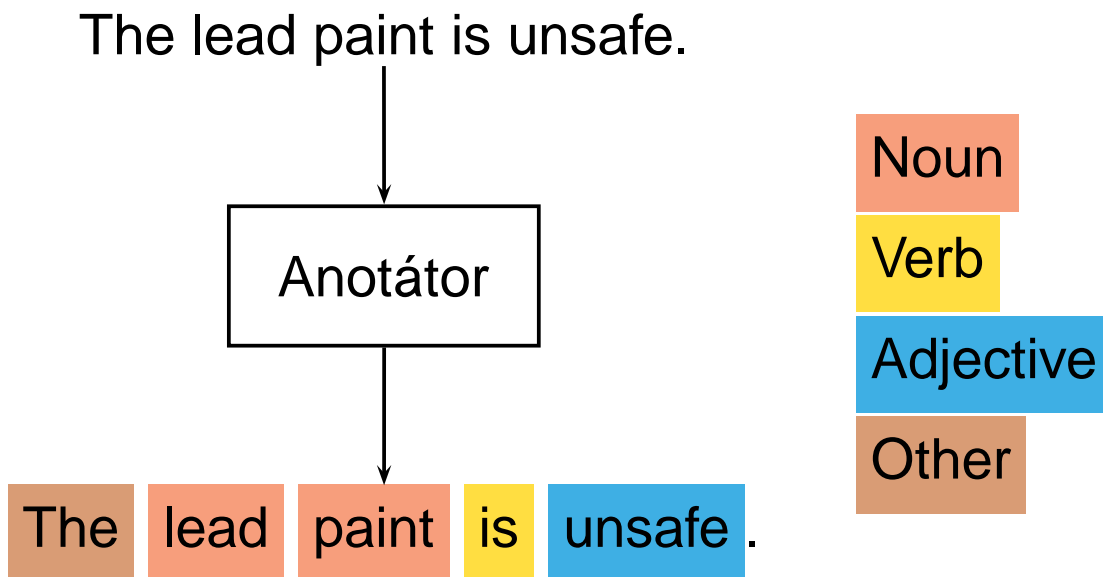
### a ich použitie na anotáciu biologických sekvencií

- Anotácia sekvencií
- Skryté Markovove modely
- Conditional random fields (CRF) [Lafferty, McCallum, Pereira 2001]
- Použitie CRF na hľadanie génov [DeCaprio et al 2007]
- Náš výskum

## Anotácia sekvencií/ret'azcov

Označ (ofarbi) každý symbol na vstupe podľa jeho funkcie/významu.

### Príklad: určovanie slovných druhov (part-of-speech tagging)



Zdroj: Andrew McCallum

## Použitie v biológii: hľadanie génov

**Vstup:** DNA sekvencia  $X = x_1, \dots, x_n$

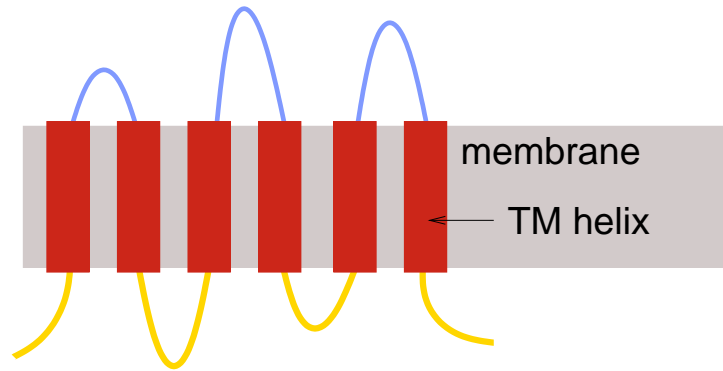
```
tgggcgtat t t t t g c g c t a g t g t t g g g t g t t c c g c t g t g c t g t t t t t c c g t c a t g g c t c g c a  
c t a a g c a a a c t g c t c g g a a g t c t a c t g g t g g c a a g g c g c c a c g c a a a c a g t t g g c c a c t a  
a g g c a g c c c g c a a a a g c g c t c c g g c c a c c g g c g g c g t g a a a a a g c c c c a c c g c t a c c g g c  
g t a a a c t a c c t t t c c a g c g c c t g g t g c g c g a g a t t g c g c a g g a c t t t a a a a c a g a c c t g c  
a c a t c c a g c t c g c c c g c c g c a t c c g c g g a g a g a g g g c g t g a t t a c t g t g g t c t c t c t g a c
```

**Výstup:** anotácia (postupnosť farieb)  $\hat{A} = a_1, a_2, \dots, a_n$

Červená: kóduje proteín, žltá: intrón, modrá: medzigénová oblasť

```
tgggcgtat t t t t g c g c t a g t g t t g g g t g t t c c g c t g t g c t g t t t t t c c g t c a t g g c t c g c a  
c t a a g c a a a c t g c t c g g a a g t c t a c t g g t g g c a a g g c g c c a c g c a a a c a g t t g g c c a c t a  
a g g c a g c c c g c a a a a g c g c t c c g g c c a c c g g c g g c g t g a a a a a g c c c c a c c g c t a c c g g c  
g t a a a c t a c c t t t c c a g c g c c t g g t g c g c g a g a t t g c g c a g g a c t t t a a a a c a g a c c t g c  
a c a t c c a g c t c g c c c g c c g c a t c c g c g g a g a g a g g g c g t g a t t a c t g t g g t c t c t c t g a c
```

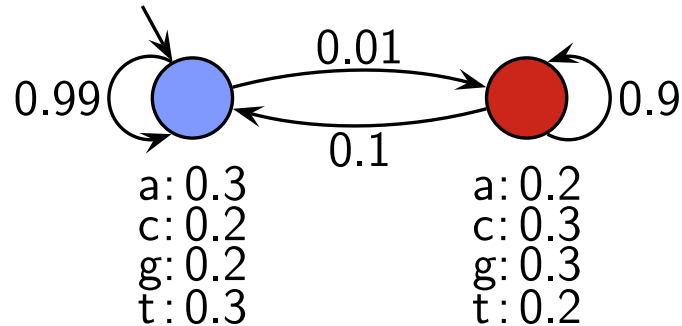
## Použitie v biológii: topológia transmembránových proteínov



Červená: TM hélix, žltá: vnútri bunky, modrá: vonku

MNYSLHLAFVCLSLFTERMCIQGSQFNVEVGRSDKLSLPGFENLTAGYNKFLRPNFGGEP  
VQIALTLDIASISSISESNMDYTATIYLRQRWMDQRLVFEGNKSFTLDARLVEFLWVPDT  
YIVESKKSFLHEVTVGNRLIRLFSNGTVLYALRITTTVACNMDLSKYPMDTQTCKLQLES  
WGYDGNDVEFTWLRGNDSVRGLEHLRLAQYTIERYFTLVTRSQQETGNYTRLVLQFELRR  
NVLYFILETYVPSTFLVVLVSWVSEWISLDSVPARTCIGVTTVLSMTTLMIGSRTSLPNTN  
CFIKAIDVYLGICFSFVFGALLEYAVAHYSSLQQMAAKDRGTTKEVEEVSITNIINSSIS  
SFKRKISFASIEISSDNVDYSDLTMKTSDFKFKFVFREKMGRIVDYFTIQNPSNVDHYSKL  
LFPLIFMLANVFYWAYYMYF

## Skryté Markovove modely (hidden Markov models, HMMs)



Sekvencia:  $\{a, c, g, t\}^*$

$X = x_1, x_2, \dots, x_n$

Anotácia:  $\{\square, \blacksquare\}^*$

$A = a_1, a_2, \dots, a_n$

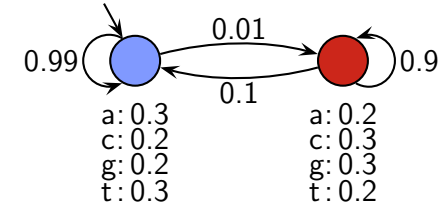
tatttagcgtcttctatcatccaatcactgcactttacacactataaatagagcagctca  
 tgggcgtatttgcgctagtggtgggtggtccgctgtgctgtttttccgctc**atggctcgca**  
**ctaagcaaactgctcggaa**gtctactggtggcaaggcgccacgcaaacagttggccacta

HMM definuje  $P(X, A)$  pre sekvencie  $X$  a anotácie  $A$ :

$$P(X, A) = s(a_1)e(a_1, x_1)t(a_1, a_2)e(a_2, x_2) \cdots t(a_{n-1}, a_n)e(a_n, x_n)$$

## Viterbiho algoritmus [Forney 1973]

Vstup: reťazec  $X = x_1, \dots, x_n$  a HMM



tatttagcgtcttctatcatccaatcactgcactttacacactataaatagagcagctca  
tgggcgtat ttgcgctagtgttgggtgttccgctgtgctgtttttccgctcatggctcgca  
ctaagcaaactgctcggaaagtctactggtggcaaggcgccacgcaaacagttggccacta

Cieľ: nájdí **najpravdepodobnejšiu anotáciu**  $A^* = \arg \max_A P(X, A)$

tatttagcgtcttctatcatccaatcactgcactttacacactataaatagagcagctca  
tgggcgtat ttgcgctagtgttgggtgttccgctgtgctgtttttccgctc**atggctcgca**  
**ctaagcaaactgctcggaa**gtctactggtggcaaggcgccacgcaaacagttggccacta

Jednoduché dynamické programovanie

$O(n|E|)$  čas ( $n$  = dĺžka sekvencie,  $|E|$  = počet prechodov)

## Trénovanie HMM z anotovaných dát

**Kritérium maximálnej vierohodnosti** (maximum likelihood)

Poznáme  $(X, A)$

Chceme parametre  $\theta$  také, že  $P_{\theta}(X, A)$  je maximálna

$$\theta^* = \arg \max_{\theta} P_{\theta}(X, A)$$

**Jednoduché počítanie frekvencií**, napr.

$$t(\blacksquare, \blacksquare) = \frac{|\{i \mid a_i = \blacksquare \wedge a_{i+1} = \blacksquare\}|}{n}$$



## Trénovanie HMM z neanotovaných dát

Poznáme  $X$

Chceme parametre  $\theta^*$  také, že

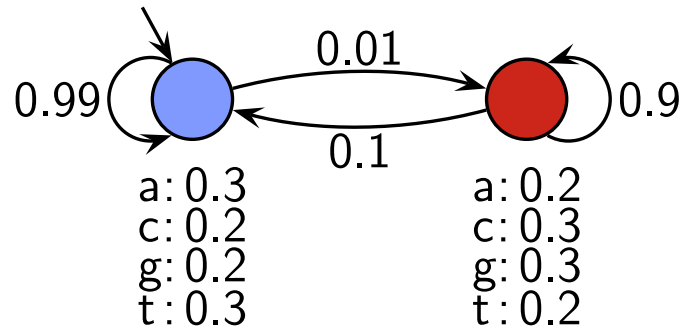
$$\theta^* = \arg \max_{\theta} P_{\theta}(X) = \arg \max_{\theta} \sum_A P_{\theta}(X, A)$$

**Iteratívny Baum-Welchov algoritmus** (verzia EM)

Každá iterácia dynamické programovanie  $O(n|E|)$ , zvýši  $P_{\theta}(X)$

Nemusí konvergovať do globálneho optima

## Zhrnutie HMM



Sekvencia:  $\{a, c, g, t\}^*$

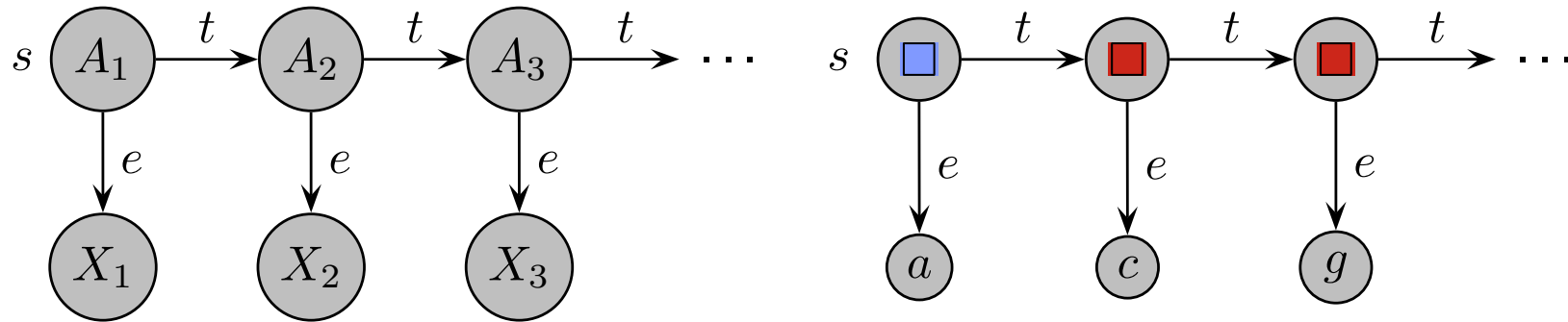
$X = x_1, x_2, \dots, x_n$

Anotácia:  $\{\square, \square\}^*$

$\tilde{A} = a_1, a_2, \dots, a_n$

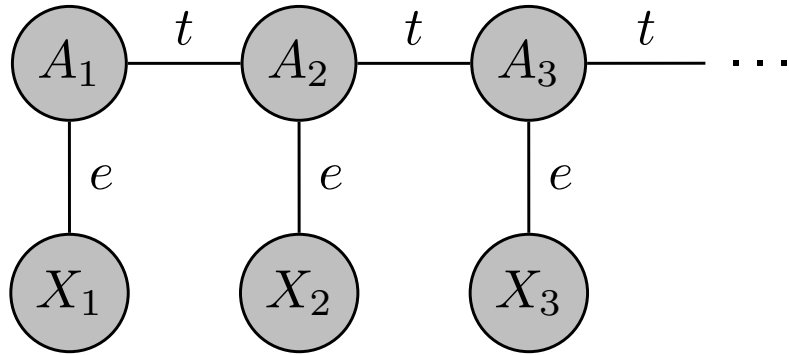
- Určíme stavový automat HMM
- Pravdepodobnosti trénujeme na anotovaných alebo neanotovaných dátach
- Efektívny algoritmus na hľadanie najpravdepodobnejšej anotácie

## HMM ako bayesovská sieť



$$P(X, A) = s(a_1)e(a_1, x_1)t(a_1, a_2)e(a_2, x_2) \cdots \\ t(a_{n-1}, a_n)e(a_n, x_n)$$

## Prvý krok k CRF: zmažeme šípky



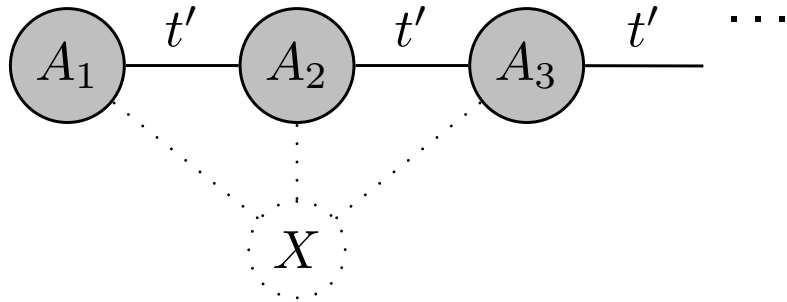
$$P(X, A) = \frac{1}{Z} e(a_1, x_1) t(a_1, a_2) e(a_2, x_2) \cdots t(a_{n-1}, a_n) e(a_n, x_n)$$

$$Z = \sum_{X', A'} P(X', A')$$

**Markovská sieť** (neorientovaný grafový model)

$e(a, x)$  a  $t(a, a')$  nenormalizované, ľubovoľné nezáporné čísla

## Druhý krok k CRF: “zabudneme” na $X$



$$P(A|X) = \frac{1}{Z(X)} t'(1, a_1, a_2, X) \cdots t'(n-1, a_{n-1}, a_n, X)$$

$$Z(X) = \sum_{A'} P(A'|X)$$

Sekvencia  $X$  je daná, nepotrebuje modelovať jej pravdepodobnosť.

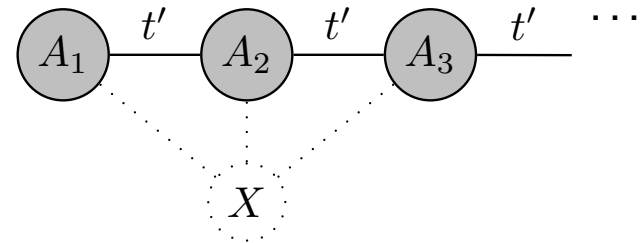
Ovplyvňuje však pravdepodobnosť anotácie  $A$ .

Funkcia  $t'(j, a, a', X)$  závisí od dvoch susedných stavov a celého  $X$ .

## Typický zápis CRF

$$P(A|X) = \frac{1}{Z(X)} t'(j, a_1, a_2, X) \cdots t'(j, a_{n-1}, a_n, X)$$

$$Z(X) = \sum_{A'} P(A'|X)$$



## Vážený součet lokálních atribútov

$$P(A|X) = \frac{1}{Z(X)} \exp \left( \sum_{i=1}^k \lambda_i \sum_{j=1}^{n-1} f_i(j, a_j, a_{j+1}, X) \right)$$

**Atribút** (feature)  $f_i(j, a, a', X)$

**Váha** atribútu  $\lambda_i$

$$t'(j, a, a', X) = \prod_{i=1}^k (e^{f_i(j, a, a', X)})^{\lambda_i}$$

## HMM sú špeciálnym prípadom CRF

$$P(A|X) = \frac{1}{Z(X)} \exp \left( \sum_{i=1}^k \lambda_i \sum_{j=1}^{n-1} f_i(j, a_j, a_{j+1}, X) \right)$$

$$Z(X) = \sum_{A'} P(A'|X)$$

### Parametre HMM ako atribúty, jednotkové váhy:

$$f_1(j, a, a', X) = \ln e(a, x_j)$$

$$f_2(j, a, a', X) = \ln t(a, a')$$

$$f_3(j, a, a', X) = \ln s(a) \text{ ak } j = 1; \text{ inak } 0$$

$$\lambda_1 = \lambda_2 = \lambda_3 = 1$$

## HMM sú špeciálnym prípadom CRF

$$P(A|X) = \frac{1}{Z(X)} \exp \left( \sum_{i=1}^k \lambda_i \sum_{j=1}^{n-1} f_i(j, a_j, a_{j+1}, X) \right)$$

$$Z(X) = \sum_{A'} P(A'|X)$$

### Parametre HMM ako váhy, indikátorové atribúty:

$$f_{e,b,y}(j, a, a', X) = \delta(b, a) \delta(y, x_j)$$

$$f_{t,b,b'}(j, a, a', X) = \delta(b, a) \delta(b', a')$$

$$f_{s,b}(j, a, a', X) = \delta(b, a) \delta(j, 1)$$

$$\lambda_{e,b,y} = \ln e(b, y)$$

$$\lambda_{t,b,b'} = \ln t(b, b')$$

$$\lambda_{s,b} = \ln s(b)$$



## Tretí krok k CRF: zmeníme tréningové kritérium

Atribúty volíme ručne, trénujeme váhy  $\lambda_i$

Pri HMM maximalizujeme vierohodnosť  $\arg \max_{\theta} P(A, X)$

Pri CRF **maximalizujeme podmienenú vierohodnosť**  $\arg \max_{\lambda} P(A|X)$

### Tréning CRF:

- Nepoznáme uzavretý stav na výpočet váh z anotovaných dát
- Iteratívne algoritmy (iteratívne škálovanie, gradientové metódy)
- V každej iterácii dynamické programovanie na tréningových dátach
- Logaritmus podmienenej vierohodnosti je konkávny
  - konvergujeme ku globálnemu maximu

## Zhrnutie CRF

$$P(A|X) = \frac{1}{Z(X)} \exp \left( \sum_{i=1}^k \lambda_i \sum_{j=1}^{n-1} f_i(j, a_j, a_{j+1}, X) \right)$$

- Nemodelujú pravdepodobnosť známych dát  $X$
- Môžeme si voliť ľubovoľné atribúty závisiace od celej sekvencie, prípadne ďalších dostupných dát
- Pri tréovaní maximalizujeme podmienenú pravdepodobnosť anotácie
- Tréovanie náročnejšie  
(iteratívna optimalizácia namiesto počítania frekvencií)
- Parametre nemajú intuitívny význam
- Použitie natrénovaného modelu stále Viterbiho algoritmom

## Použitie CRF na hľadanie génov

**Vstup:** DNA + ďalšie zdroje informácií

**Výstup:** označenie oblastí kódujúcich proteíny

Bohatá história používania zovšeobecnených HMM

Program Conrad [DeCaprio et al 2007] založený na CRF

Používa semi-Markov CRF

- atribúty závisia od celého úseku s rovnakým stavom
- čas výpočtu potom kvadratický od  $n$
- umožňuje modelovať dĺžky úsekov

## Atribúty v Conrade

- Parametre HMM rozdelené do niekoľkých skupín, každá zvláštu váhu  
Např. emisia v stave  $b$ :  $f_b(j, a, a', X) = \delta(a, b) \ln e(b, x_j)$   
Emisie v rôznych stavov môžu mať rôznu váhu
- Použitie ďalších dát: porovnanie s genómami príbuzných druhov,  
experimentálne overenie transkripcie (indikátorové atribúty)

## Zlepšenie presnosti prechodom na CRF

Training method	No. of species	Additional features	Nucl. sens.	Nucl. spec.	Exon sens.	Exon spec.	Gene sens.	Gene spec.
MEA-splice	5	Gap/Foot/EST	98.8	99.4	95.0	96.7	86.0	86.0
<b>CML</b>	<b>2</b>	—	<b>97.9</b>	<b>98.2</b>	<b>85.6</b>	<b>90.2</b>	<b>61.3</b>	<b>61.8</b>
MEA-Nucleotide	2	—	98.0	98.1	86.0	90.0	61.1	61.0
MEA-splice	2	—	97.8	98.1	85.8	89.9	60.8	60.8
GHMM	2	—	94.7	98.4	82.3	87.0	52.7	54.1

# Adaptácia Conradu na mitochondriálne genómy kvasiniek

## Diplomová práca Juraja Meštánka 2010

Mierná modifikácia základného modelu

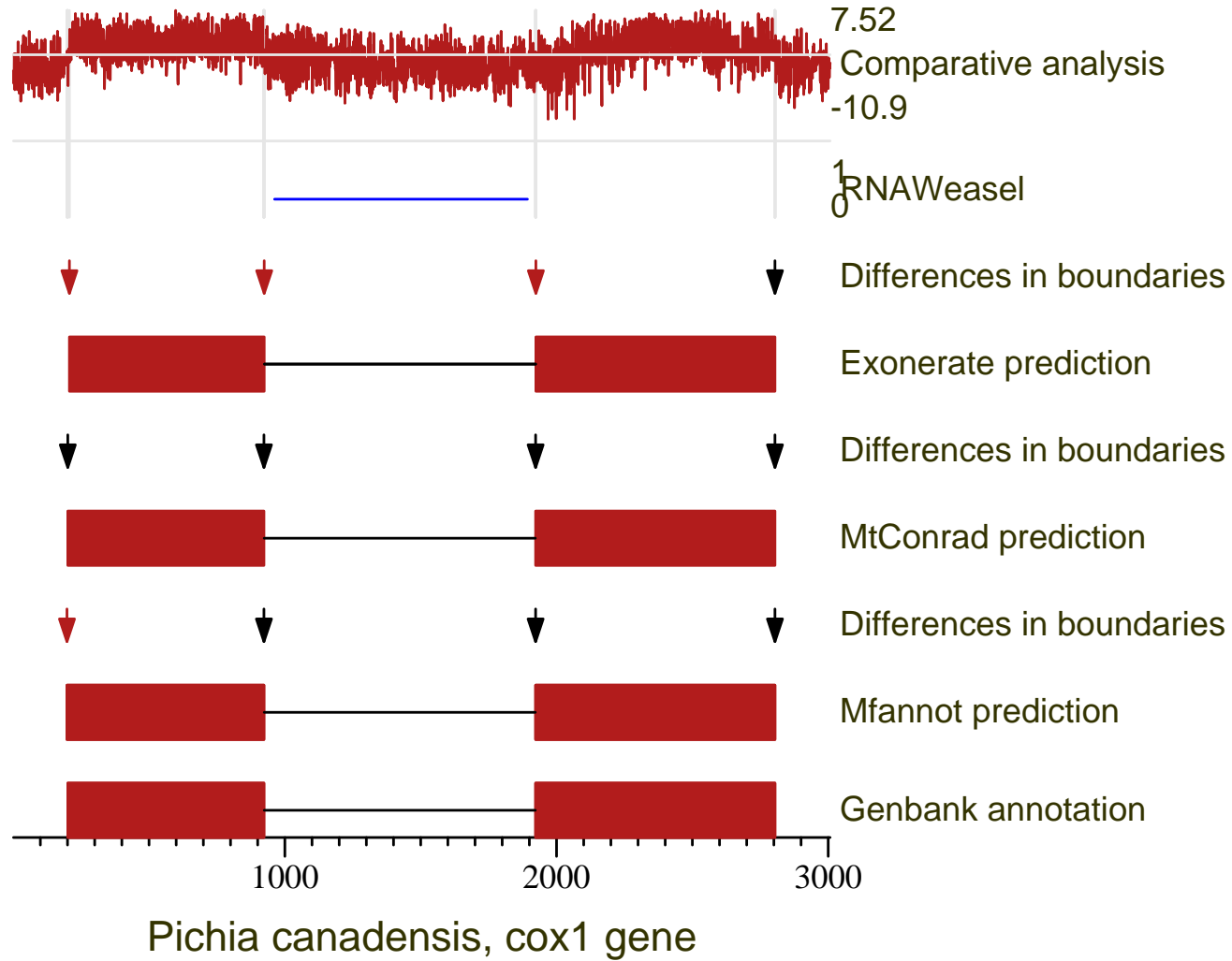
Využitie viacerých druhov prídavnej informácie

Zostavenie tréningových a testovacích dát

### Porovnanie presnosti

Tool	Perfect	Exon		Intron		Coding nucl.	
		sens	spec	sens	spec	sens	spec
HMM	58%	45%	27%			96%	99%
<b>MtConrad</b>	<b>78%</b>	70%	74%	55%	59%	<b>99%</b>	99%
MFannot	74%	<b>77%</b>	<b>81%</b>	<b>82%</b>	<b>82%</b>	94%	<b>99.8%</b>

# Ukážka jedného génu



Červené šípky nesprávne okraje.

## Ďalšie plány

- Implementácia CRF tréovania do nášho hľadača génov ExonHunter (prototyp diplomant Lukáš Bača)
- Oproti Conradu ExonHunter lepšie spracováva dĺžky úsekov, má systém na tréovanie na nových organizmoch.
- Rozšírenie HMM resp. CRF na zložitejšie typy prídavnej informácie (diplomant Marcel Kucharík)
- Rôzne kritériá odvodzovania v HMM resp. CRF (doktorand Michal Nánási)