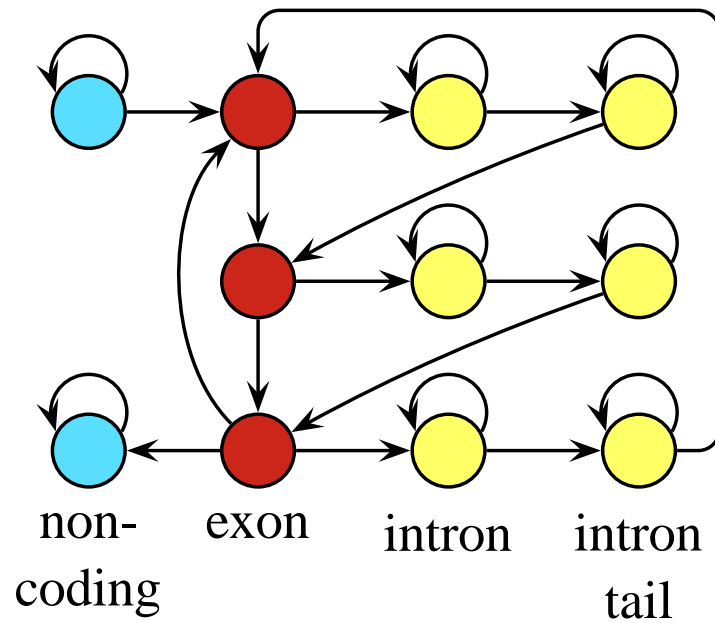
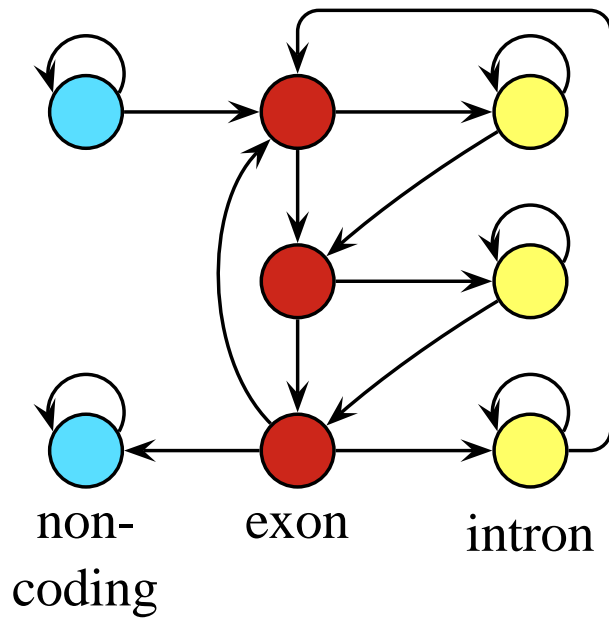


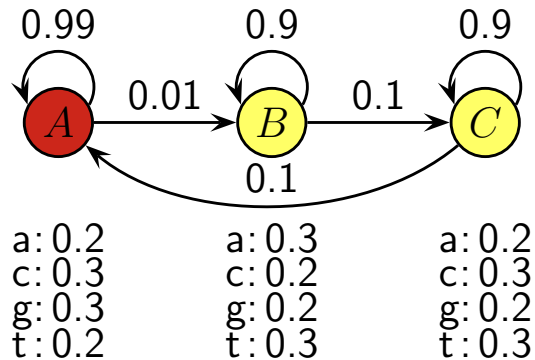
## Najpravdepodobnejšia anotácia na zložitejších HMM

Môžeme mať viacero stavov s tým istým významom.

**Príklad:** koniec intrónu je obohatený o C a T



## Cesta vs. anotácia



Sekvencia  $X = x_1 \dots x_n \in \{a, c, g, t\}^*$

Cesta  $\pi = \pi_1 \dots \pi_n \in \{A, B, C\}^*$

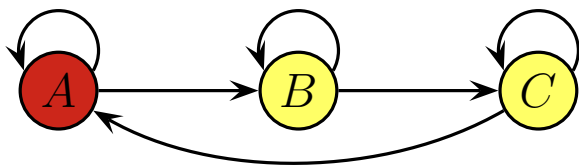
Anotácia  $\bar{A} = a_1 \dots a_n \in \{\blacksquare, \square\}^*$

AAAAAAAAAAAAAAAAAAAA
BBBBBBBBBBCCCCCC
AAAAA  
gttccgctgtgctgttttt
gtaggctcgcactaagg
tcta

$$P(X, \pi) = s(\pi_1) e(\pi_1, x_1) \prod_{i=2}^n t(\pi_{i-1}, \pi_i) e(\pi_i, x_i)$$

$$P(X, \bar{A}) = \sum_{\pi: \mathcal{A}(\pi) = \bar{A}} P(X, \pi)$$

## Cesta vs. anotácia



Sekvencia  $X = x_1 \dots x_n \in \{a, c, g, t\}^*$

Cesta  $\pi = \pi_1 \dots \pi_n \in \{A, B, C\}^*$

Anotácia  $\bar{A} = a_1 \dots a_n \in \{\text{red}, \text{yellow}\}^*$

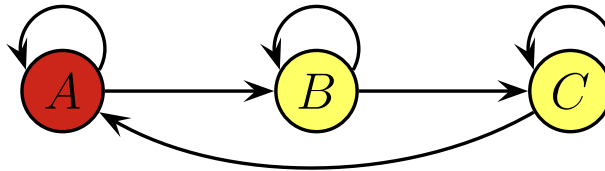
$$P(X, \bar{A}) = \sum_{\pi: \mathcal{A}(\pi) = \bar{A}} P(X, \pi)$$

Viterbiho algoritmus hľadá najpravdepodobnejšiu cestu

$$\pi^* = \arg \max_{\pi} P(X, \pi)$$

My chceme najpravdepodobnejšiu anotáciu  $\bar{A}^* = \arg \max_{\bar{A}} P(X, \bar{A})$

## Cesta vs. anotácia



$p_1$  : **tcggaagtcta ctggtggcaaggcgc cacgc aaacagttggcacta aggcagcccgc**at

$p_2$  : **tcgga**gtctact||ggtggcaaggcgc**ccacgcaa**acagttggcacta**aggcag**cccgc**at**

$p_3$  : **tcgga**gtctactggtggcaaggc||gccacgcaa**acagttggcactaaggcag**cccgc**at**

$p_4$  : **tcgga**gtctactggtggcaaggcgc**ccacg**||caa**acagttggcactaaggcag**cccgc**at**

...

$p_k$  : **tcgga**gtctactggtggcaaggcgc**ccacgcaa**acagttggcca||cta**aggcag**cccgc**at**

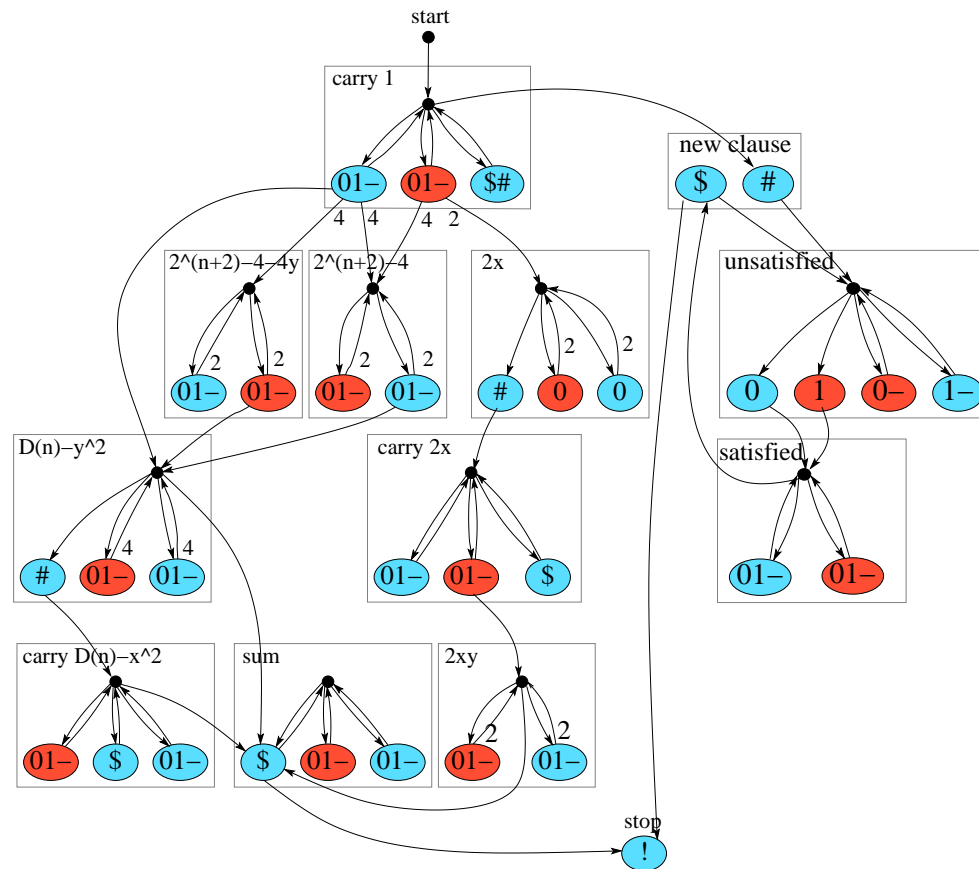
Viterbiho algoritmus:  $\arg \max\{p_1, p_2, \dots, p_k\}$

**Najpravdepodobnejšia anotácia:**  $\arg \max\{p_1, p_2 + p_3 + \dots + p_k\}$

# Problém nájdenia najpravdepodobnejšej anotácie je NP-t'ážký

[Lyngso and Pedersen 2002; B., Brown, Vinař 2007]

... nepodarí sa nám teda asi nájsť efektívny algoritmus



Logická formula

$$x_1 \wedge (x_2 \vee \neg x_3)$$

zakódovaná ako sekven-  
cia

$$X = 000\#1--\$-10\$!$$

najpravdepodobnejšia


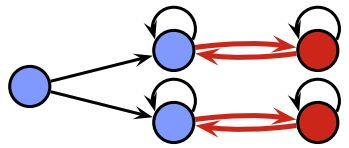
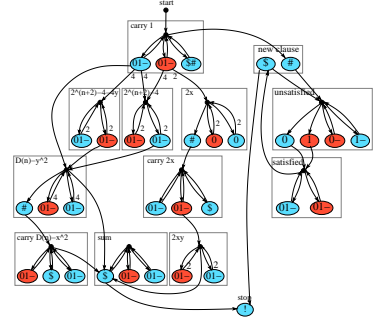

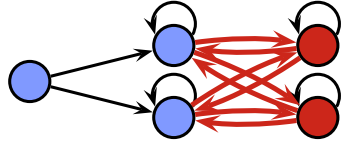
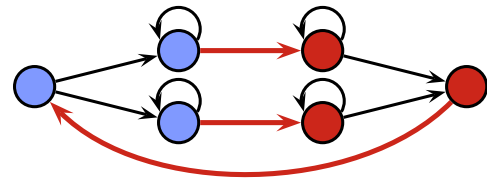
anotácia

⇒

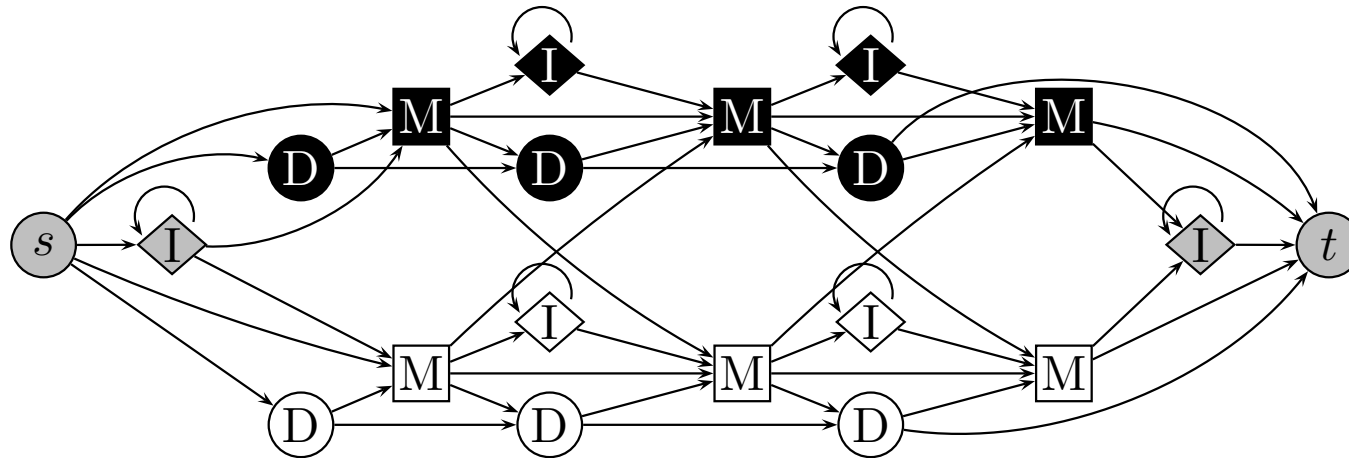
ohodnotenie premenných

**Pre niektoré HMM sa najpravdepodobnejšia anotácia dá spočítať v polynomiálnom čase**

Všetky HMMs

<p>Viterbiho alg. <math>O(nm^2)</math></p> 	<p>tiež <math>\in P</math></p> 	<p>NP-ťažké</p> 
<p><math>O(n^2m^3)</math></p> 	<p>?</p> 	
<p><math>O(n^3m^5)</math></p> 		

## Annotation issues in jumping HMMs



State path: alignment of sequence to subtype profiles

Annotation: segments of inputs emitted by subtype profiles

### Problems with most probable annotation:

- probably hard to decode
- many annotations with slightly shifted boundaries

### Change the objective function for decoding

## Gain function [Hamada et al. 2009]



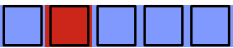

$G(A, A')$  measures accuracy of  $A$  wrt. correct annotation  $A'$

### Examples:

Identity: score 1 iff  $A$  completely correct, 0 otherwise

Pointwise: score +1 for every correct label in  $A$

Boundary: score +1 for every correct boundary,  $-\gamma$  for incorrect boundary

	Identity	Pointwise	Boundary
$A =$ 	1	5	4
$A' =$ 			
$A =$ 	0	4	$3 - \gamma$
$A' =$ 			



## Optimizing expected gain

**Goal:** find annotation  $\hat{A}$  that maximizes

$$E_{A'|X}[G(A, A')] = \sum_{A'} G(A, A')P(A'|X)$$

**Identity gain function:** Viterbi algorithm

**Pointwise gain function:** Posterior decoding (forward-backward)

**Boundary gain function:** [Gross et al. 2007]

The choice of gain function is application-dependent