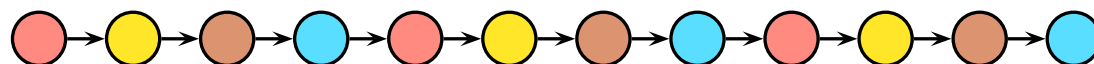


PROBLÉM NAJPRAVDEPODOBNEJŠEJ ANOTÁCIE
POMOCOU HMM
(SKRYTÉHO MARKOVSKÉHO MODELU)

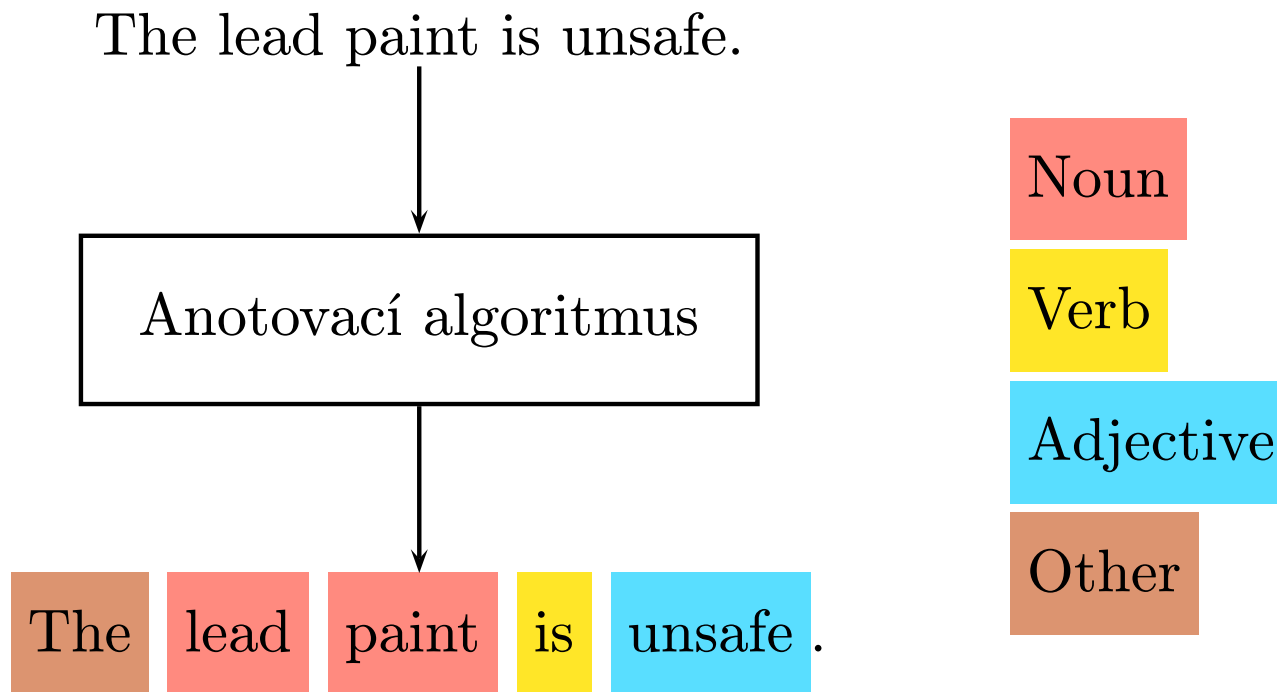


BROŇA BREJOVÁ
SPOLUATORI DAN BROWN A TOMÁŠ VINAŘ

ANOTÁCIA REŤAZCA

Ofarbi jednotlivé znaky podľa významu (závisí od kontextu)

PRÍKLAD 1: URČOVANIE SLOVNÝCH DRUHOV (TAGGING)



(Zdroj: Andrew McCallum)

HĽADANIE GÉNOV

cggtgaaactgcacgattgttgctggcttaaagatagaccaatcagagtgtgtaacgtca
tatttagcgtcttctatcatccaatcactgcactttacacactataaatagagcagctca
tgggcgtatattgcgctagtgttgggtgttccgctgtgctgtttttccgtcatggctcgca
ctaagcaaactgctcggaagtctactggtggcaaggcgccacgcaaacagttggccacta
aggcagcccgcaaaagcgctccggccaccggcggcgtgaaaaagcccaccgctaccggc
cgggcaccgtggctctgcgcgagatccgccgttatcagaagtcactgaactgcttattc
gtaaactacctttccagcgcctggtgcgcgagattgcgcaggactttaaacagacctgc
gtttccagagctccgctgtgatggctctgcaggaggcgtgcgaggcctacttggtagggc
tatttgaggacactaacctgtgcgccatccacgccaagcgcgtcactatcatgcccaagg
acatccagctcgcccgcgcgatccgcggagagagggcgtgattactgtggtctctctgac
ggtccaagcaaaggctcttttcagagccaccaccttttcaagtaaagtagctgtaagaaa
ccaatttaagacaaaagggaatgcattgggagcacttttcgttttaatgctactgaaggc

Gén

Medzigénový úsek

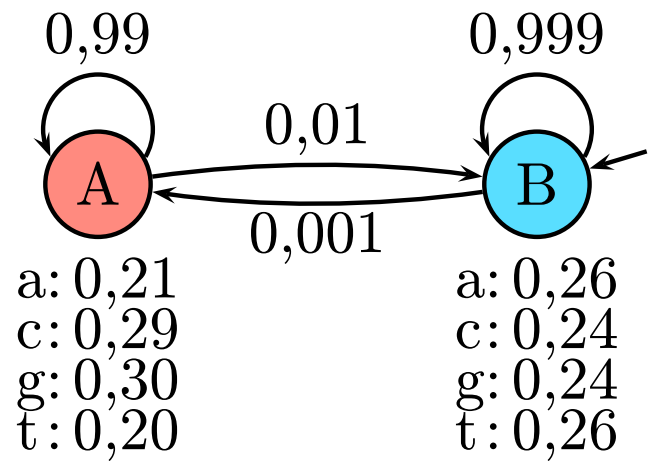
SKRYTÝ MARKOVOVSKÝ MODEL (HIDDEN MARKOV MODEL, HMM)



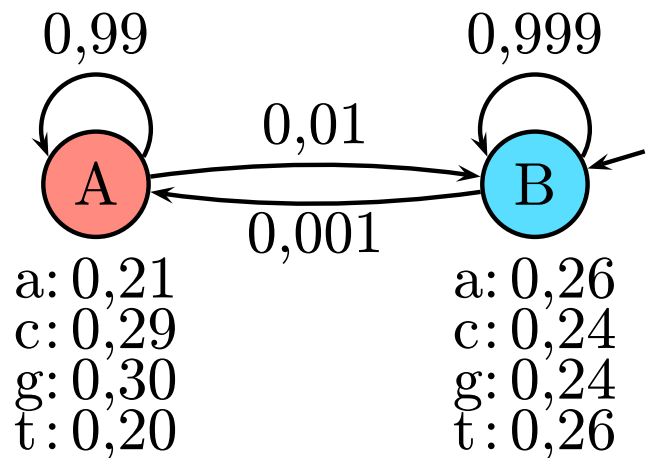
Definuje pravdepodobnosť $\Pr(X, A)$, že vygenerujeme reťazec X s anotáciou A .

“Najlepšia” anotácia reťazca X : $\arg \max_A \Pr(X, A)$

ŠTRUKTÚRA HMM



ŠTRUKTÚRA HMM



Reťazec: $\{a, c, g, t\}^*$

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

Postupnosť stavov: $\{A, B\}^*$

$S = s_1, s_2, \dots, s_n$

Anotácia: $\{\text{red}, \text{blue}\}^*$

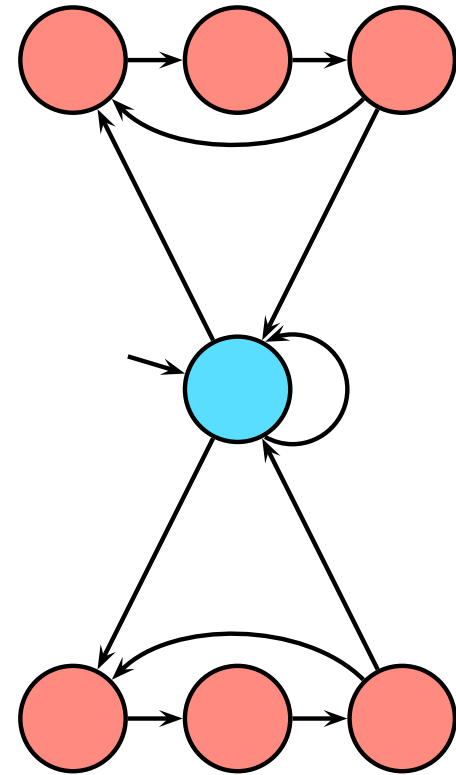
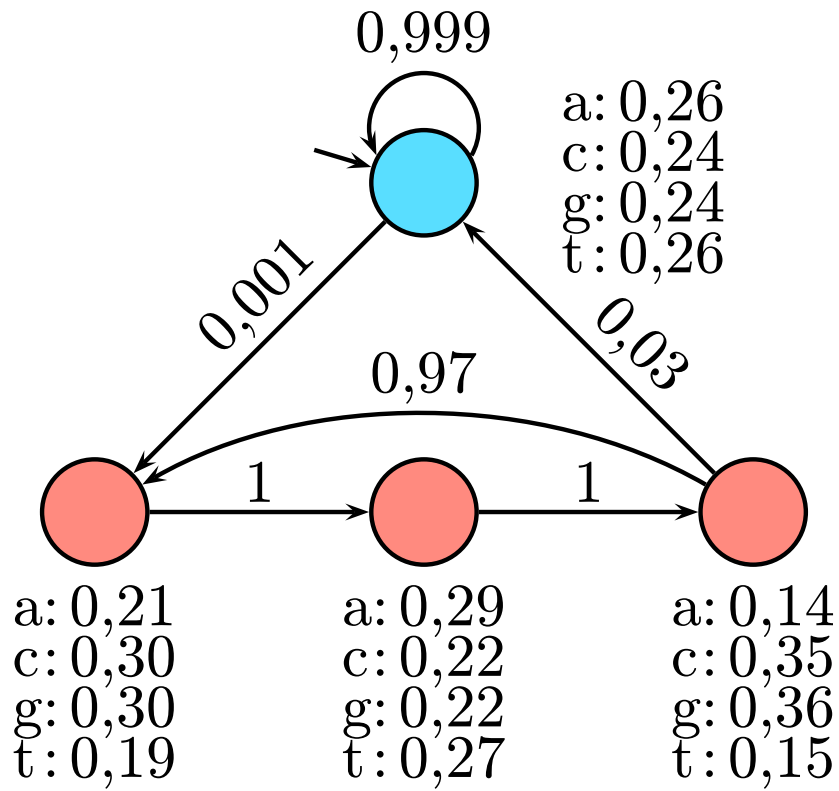
$A = a_1, a_2, \dots, a_n$

$a_i = \text{farba } s_i$

$$\Pr(X, S) = \Pr(x_1|s_1) \Pr(s_1 \rightarrow s_2) \Pr(x_2|s_2) \cdots \Pr(s_{n-1} \rightarrow s_n) \Pr(x_n|s_n)$$

$$\Pr(X, A) = \sum_{S: S \Rightarrow A} \Pr(X, S)$$

ZLOŽITEJŠIE HMM



PROBLÉM NAJPRAVDEPODOBNEJŠEJ ANOTÁCIE

Dané: HMM (graf, pravdepodobnosti, farby), reťazec X

Nájdí: $A = \arg \max_A \Pr(X, A) = \arg \max_A \sum_{S \Rightarrow A} \Pr(X, S)$

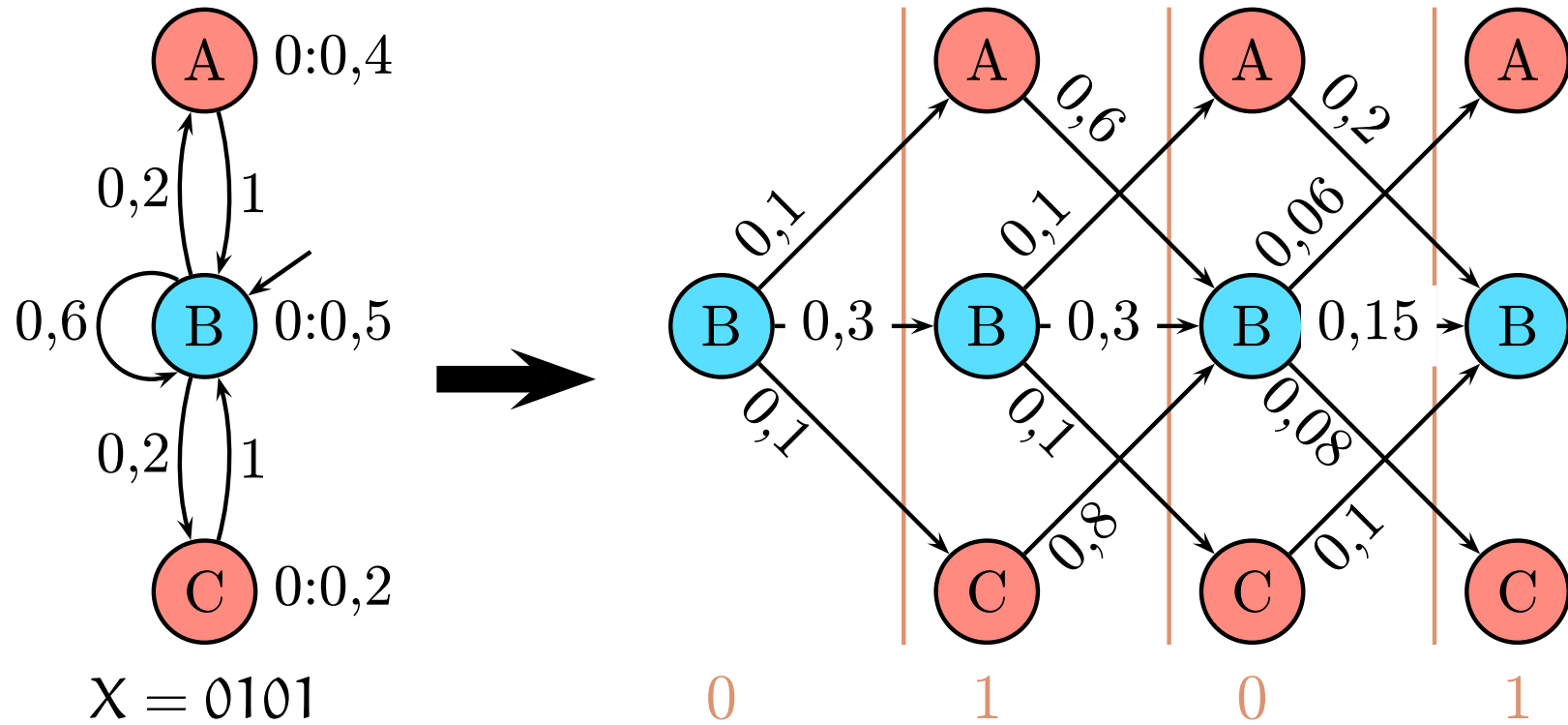
ŤAŽKÝ PROBLÉM:

- NP-ťažké pre vstup HMM a X [Lyngsø, Pedersen 2002]
- NP-ťažké pre jeden pevný HMM a vstup X [BBV 2004]

PRE NIEKTORÉ HMM POLYNOMIÁLNE:

- Ak pre každé A jediné S [Viterbi 1967]
- Všeobecnejšia trieda HMM [BBV 2004]

ZJEDNODUŠME PROBLÉM



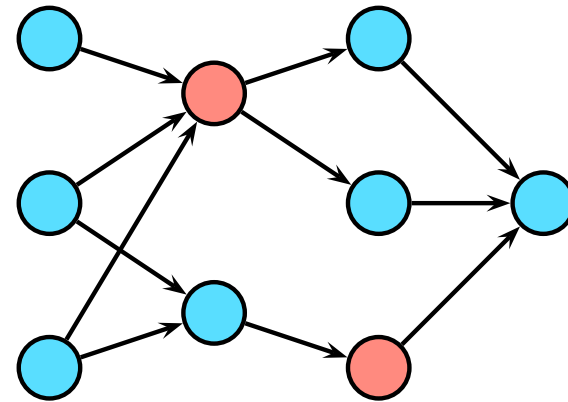
- S cesta
- $\Pr(X, S)$ súčin cien hrán
- A ofarbenie stĺpcov
- $\Pr(X, A)$ súčet zodpovedajúcich ciest

ZJEDNODUŠME PROBLÉM EŠTE VIAC

...zrušme ceny hrán

PROBLÉM OFARBENIA VRSTIEV

- Orientovaný graf
- Rozdelený do vrstiev s konšt. počtom vrcholov
- Vrcholy ofarbené 2 farbami
- Nájdi ofarbenie vrstiev, ktorému zodpovedá najviac ciest (resp. $\geq T$ ciest)



●	●	●	●	0
●	●	●	●	6
●	●	●	●	2

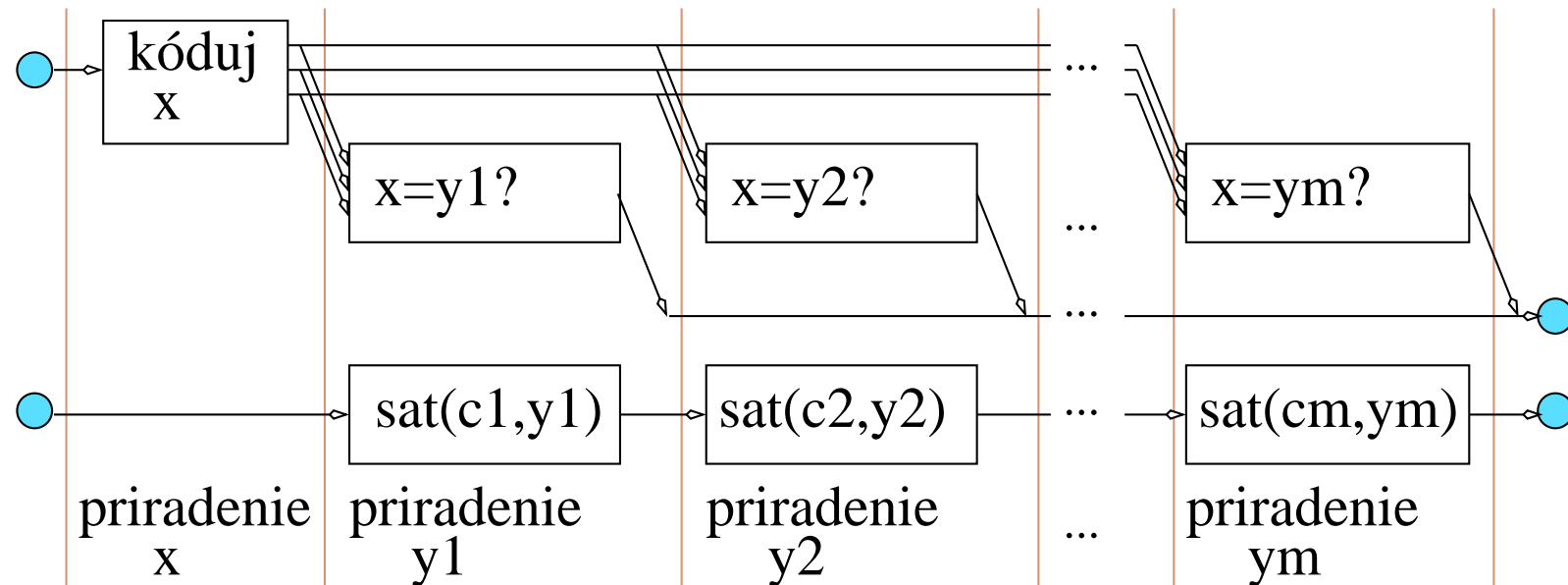
REDUKCIA ZO SAT

Formula v CNF s m klauzulami a n premennými

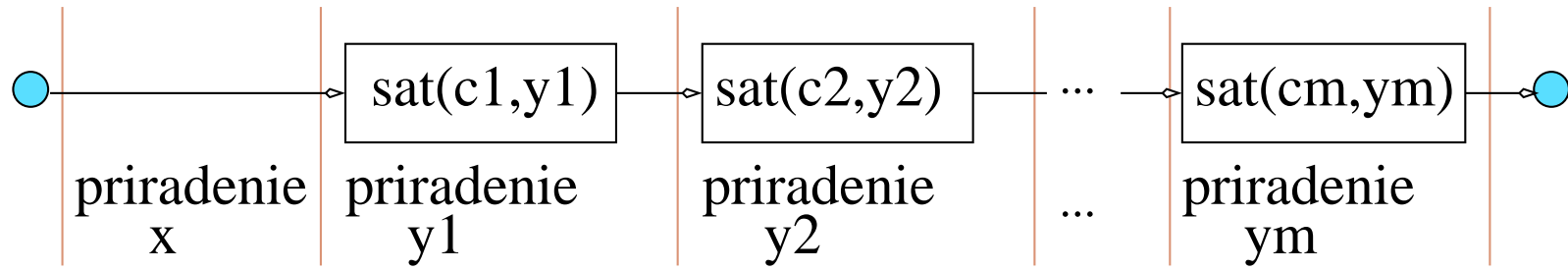
$$c_1 \wedge c_2 \wedge \dots \wedge c_m.$$

Pravdivostné priradenie zodpovedá ofarbeniu vrstiev

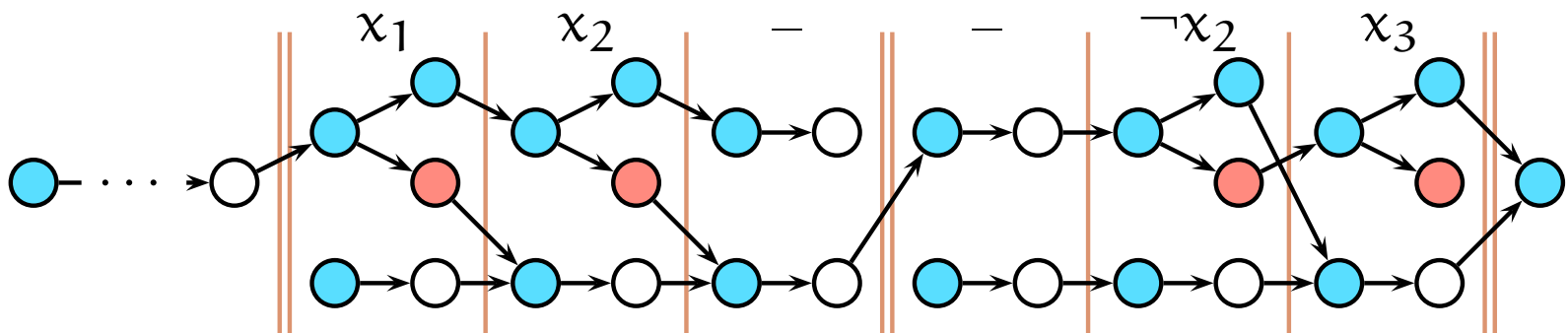
●● 0 ●● 1



JE KLAUZULA c_i SPLNENÁ PRE PRIRADENIE y_i ?



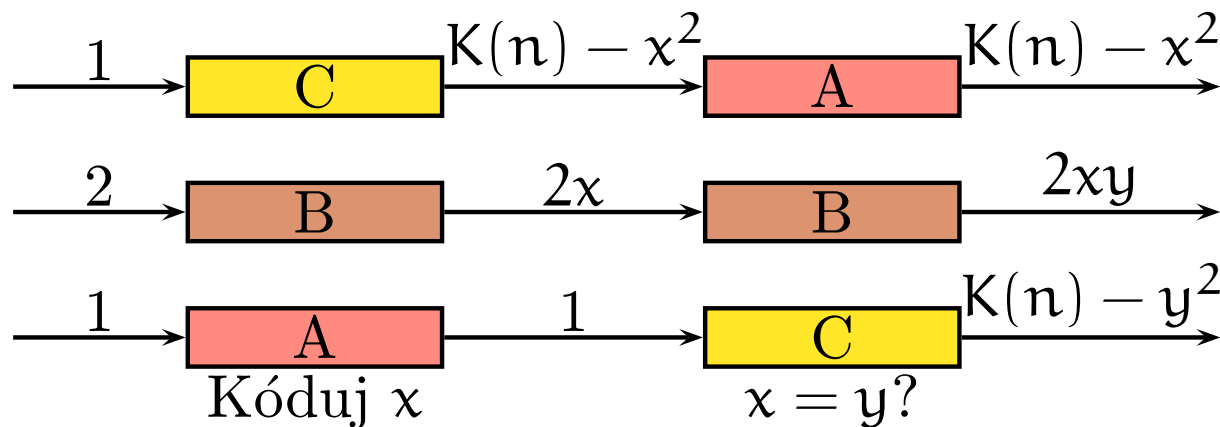
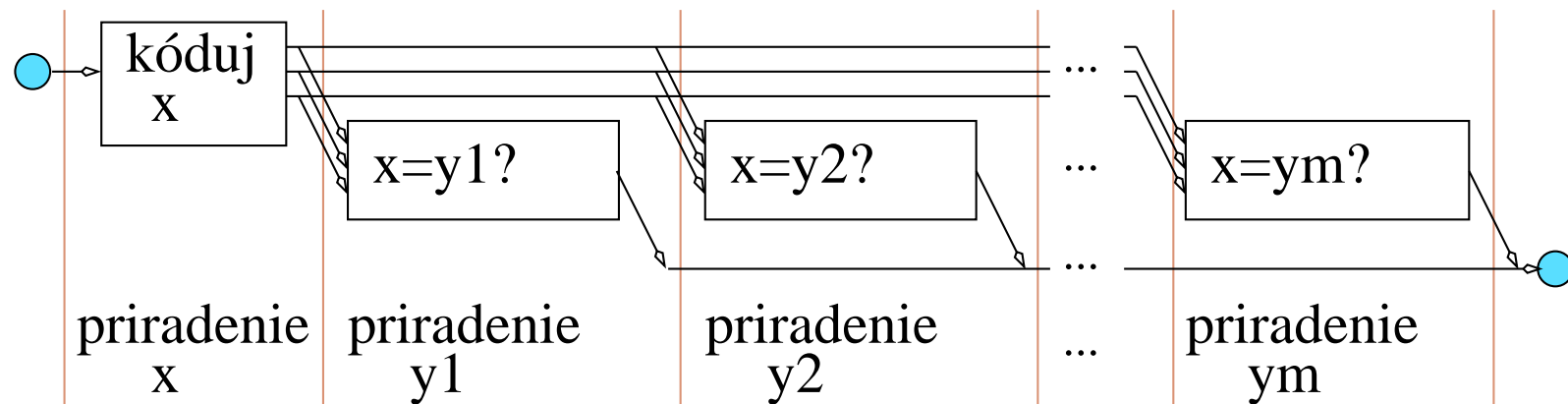
$$(x_1 \vee x_2) \wedge (\neg x_2 \vee x_3)$$



TESTOVANIE ROVNOSTI

Ofarbenie vrstiev = pravdivostné priradenie = binárne číslo

●● 0 ●● 1



Spolu:

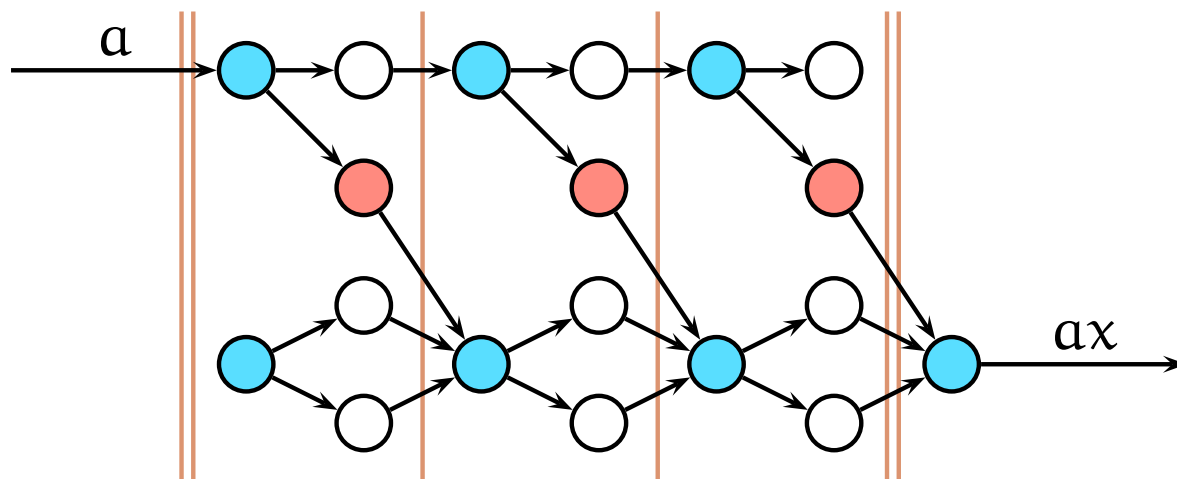
$$2K(n) - (x - y)^2$$

VYTVORENIE ČÍSLA V BINÁRNEJ SÚSTAVE

$2n$ vrstiev – n bitov: $\bullet\bullet 0$ $\bullet\bullet 1$



$n = 3$:



VÝPOČET $K(n) - x^2$

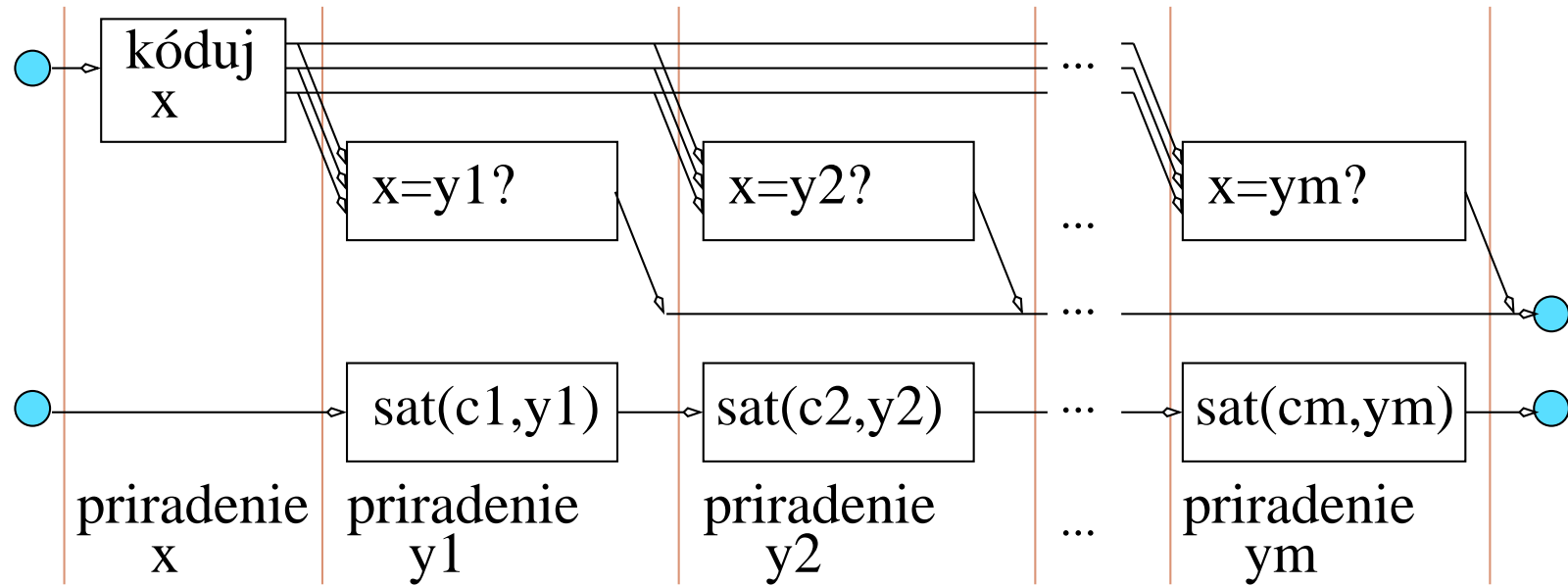


REKURENCIE PRE PRIDANIE k -TEHO BITU

Prvých k bitov tvorí číslo $z = 2y + t$ (t je k -ty bit)

$$\begin{aligned} U_k &= 2^{k+2} - 4 = 2U_{k-1} + 4 \\ V_{z,k} &= U_k - 4z = \begin{cases} 2V_{y,k-1} + 4, & \text{if } t = 0 \\ 2V_{y,k-1}, & \text{if } t = 1 \end{cases} \\ W_{z,k} &= \underbrace{(2^k - 1)^2}_{K(k)} - z^2 = \begin{cases} 4W_{y,k-1} + U_{k-1} + 1, & \text{if } t = 0 \\ 4W_{y,k-1} + V_{y,k-1}, & \text{if } t = 1 \end{cases} \end{aligned}$$

PROBLÉM OFARBENIA VRSTIEV JE NP-ŤAŽKÝ



- Ak $x = y_1 = \dots = y_m$:
 - ak x spĺňa formulu: $2^m K(n) + 1$ ciest
 - ak x nespĺňa formulu: $2^m K(n)$ ciest
- Inak: $\leq 2^m K(n)$ ciest

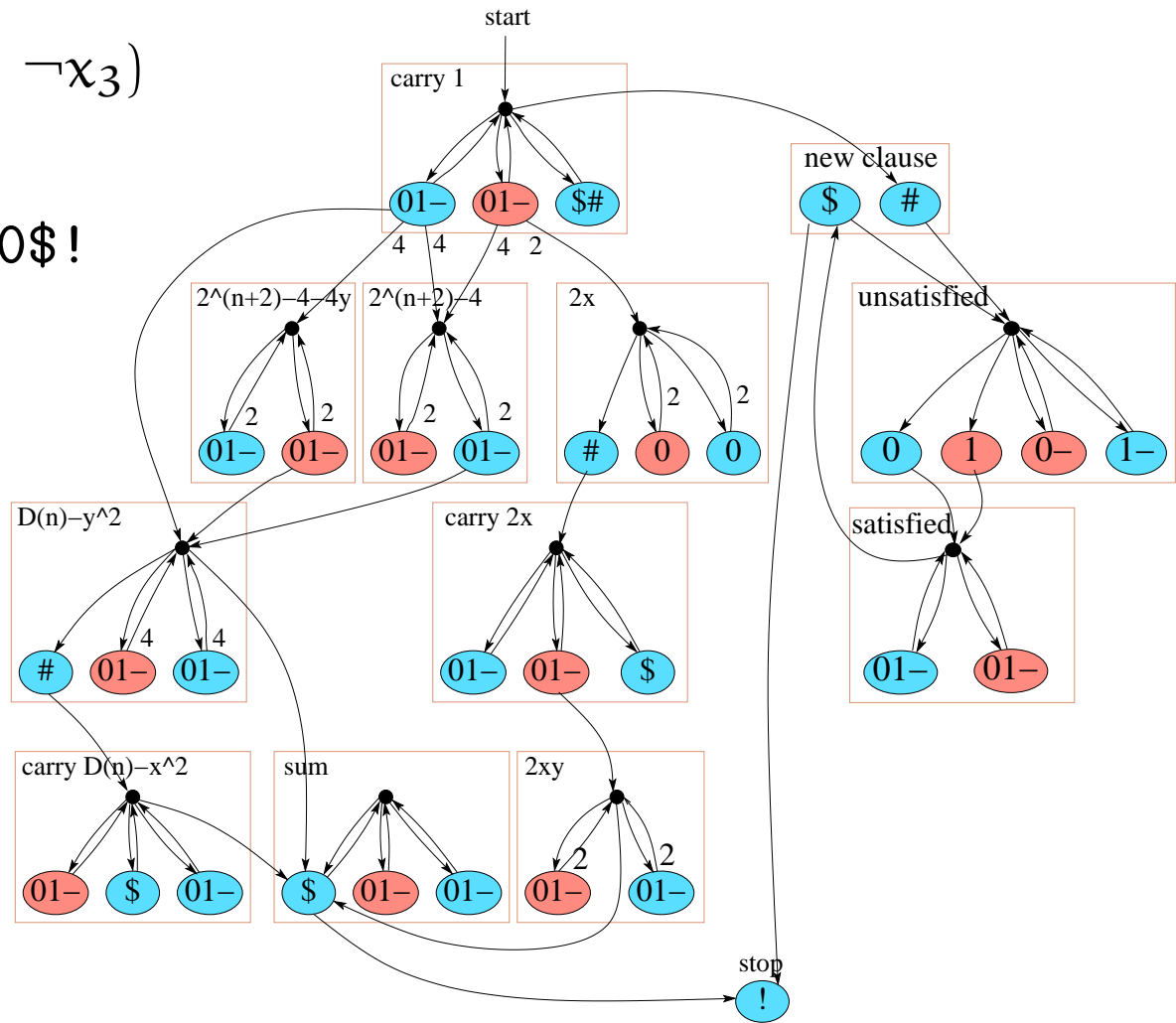
Každá vrstva má najviac 29 vrcholov

PROBLÉM NAJPRÁVDEPODOBNEJŠEJ ANOTÁCIE JE NP-ŤAŽKÝ PRE TENTO HMM

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

zodpovedá

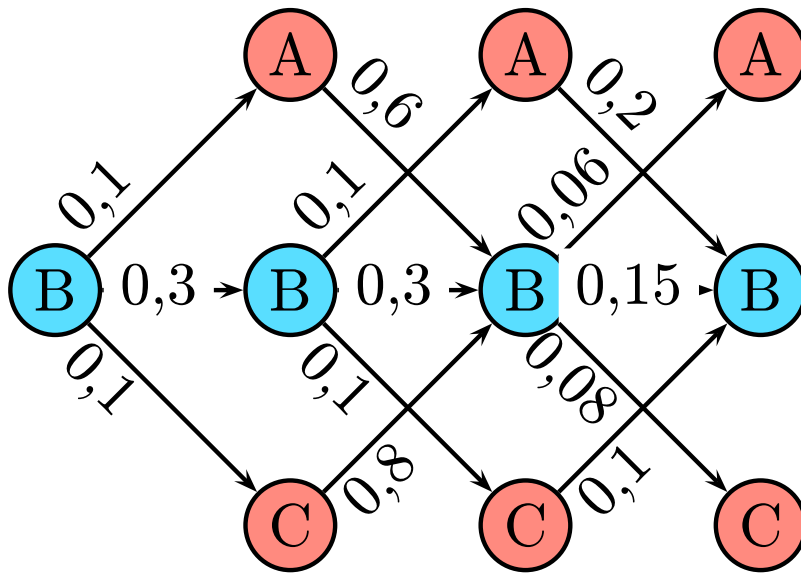
$X = 000\#1--\$-10\$!$



ALGORITMICKÁ ČASŤ

Efektívne algoritmy pre najpravdepodobnejšiu anotáciu v špeciálnych typoch HMM

VITERBIHO ALGORITMUS: nájde $\arg \max_S \Pr(X, S)$



Dyn. programovanie:

$D[v]$ = najlepšia cesta do v

$D[v] = \max_{u \rightarrow v} D[u] \cdot c(u, v)$

Čas $O(nm^2)$

$n = |X|$

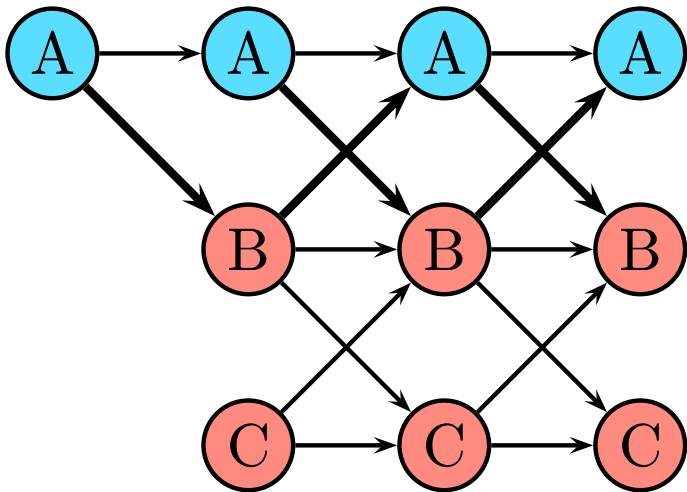
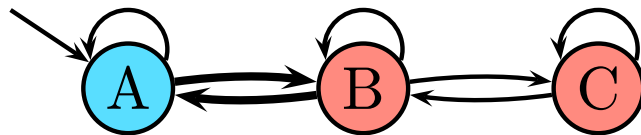
$m = \text{počet stavov}$

Ak každá anotácia A má iné S , nájde aj $\arg \max_A \Pr(X, A)$

ROZŠÍRENÝ VITERBIHO ALGORITMUS

KRITICKÁ HRANA: spája dve rôzne farby

ROZŠÍRENÁ ANOTÁCIA R: postupnosť farieb stavov a kritických hrán



Postupnosť stavov:

$S = A, B, C, B$

Rozšírená anotácia:

$R = \blacksquare, \blacksquare, \blacksquare, \blacksquare, (A \rightarrow B)$

Anotácia:

$A = \blacksquare, \blacksquare, \blacksquare, \blacksquare$

CIEĽ: najpravdepodobnejšia rozšírená anotácia R

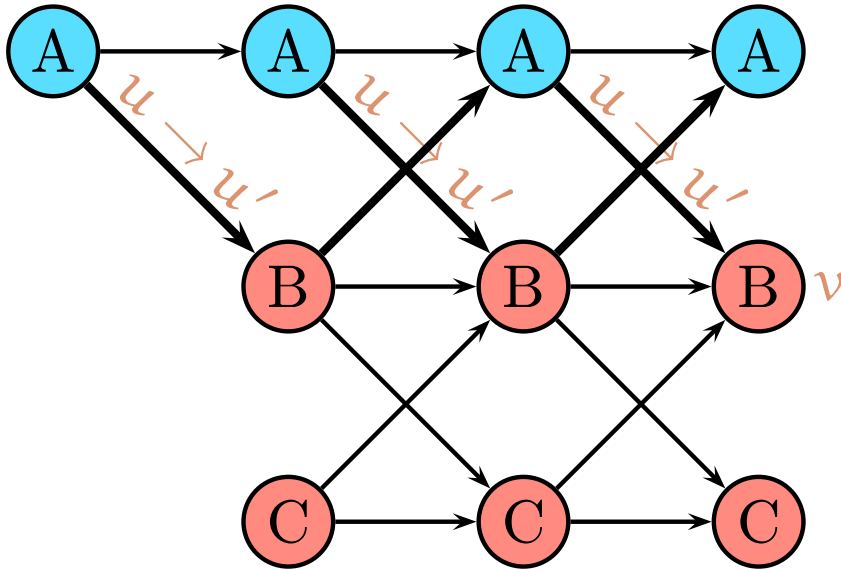
ROZŠÍRENÝ VITERBIHO ALGORITMUS

$D[v]$ = najlepšia rozšírená anotácia končiaci vo vrchole v

$$D[v] = \max_{u \rightarrow u', f(u')=f(v), f(u) \neq f(v)} D[u] \cdot c(u, u') \cdot D'[u', v]$$

$D'[s, t]$ = súčet cien jednofarebných ciest z s do t

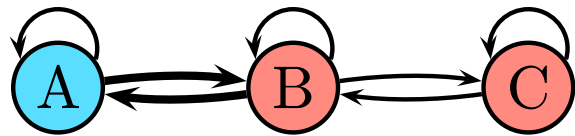
$$D'[s, t] = \sum_{s \rightarrow s', f(s)=f(s')} c(s, s') \cdot D'[s', t]$$



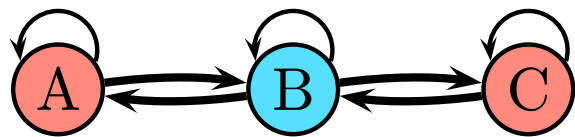
Čas $O(n^2m^3)$

ROZŠÍRENÝ VITERBIHO ALGORITMUS (RVA)

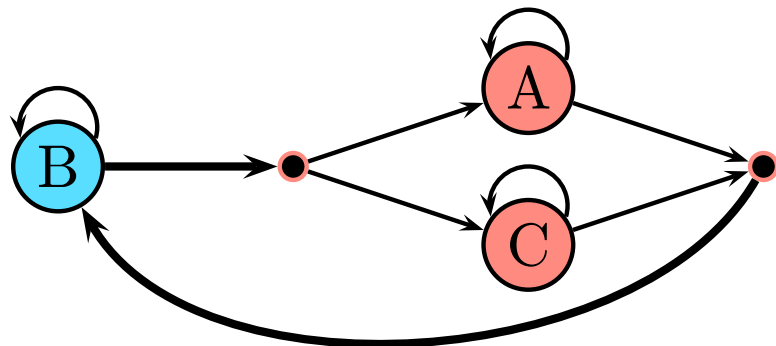
Kedy nám RVA nájde najlepšiu anotáciu A?



Pre každé A jediné R
 \Rightarrow môžeme použiť RVA.

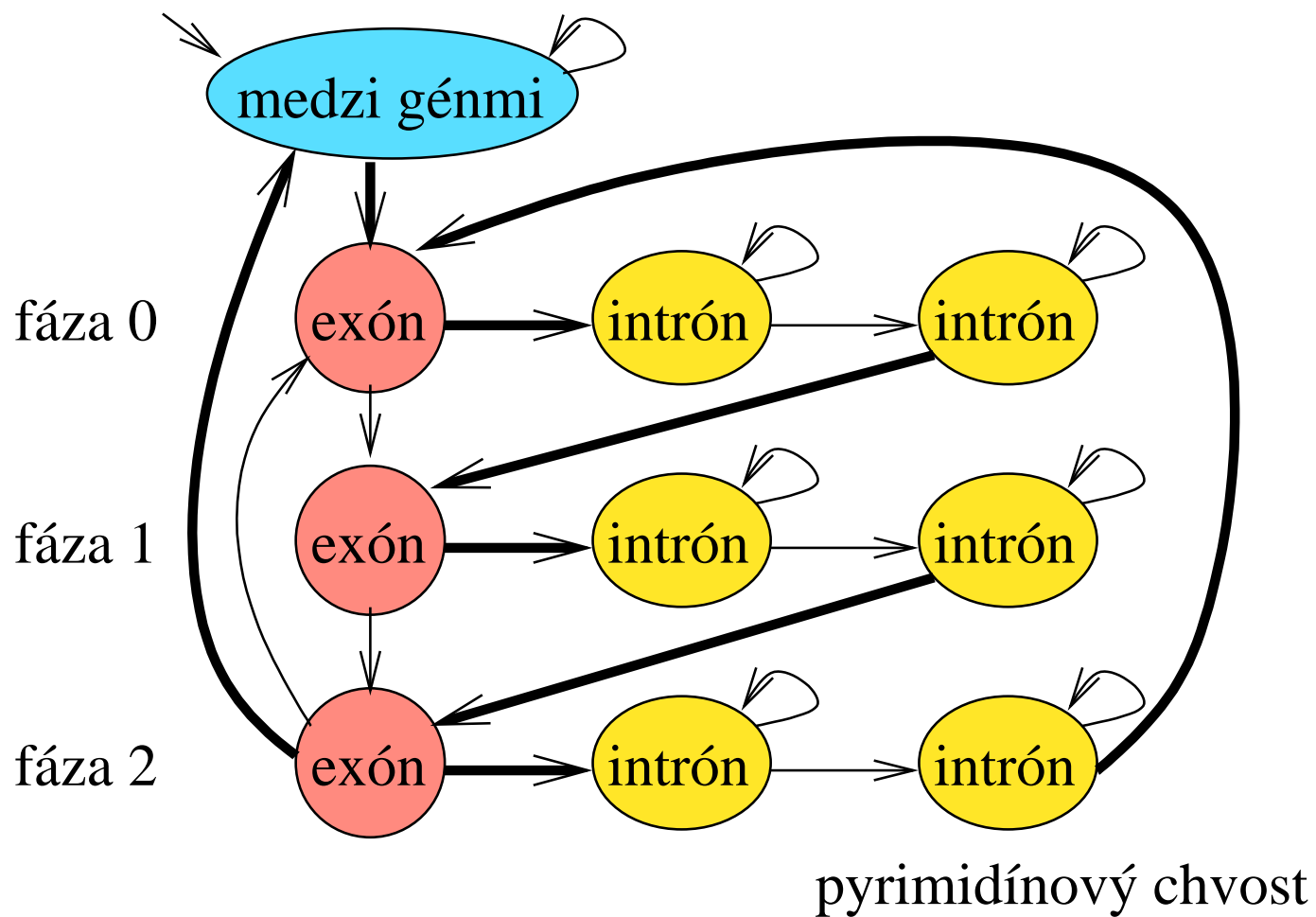


A = ■, ■, ■ má dve možné R
 \Rightarrow nemôžeme použiť RVA



Pridáme stavy, ktoré nič
neregnerujú
Môžeme použiť RVA

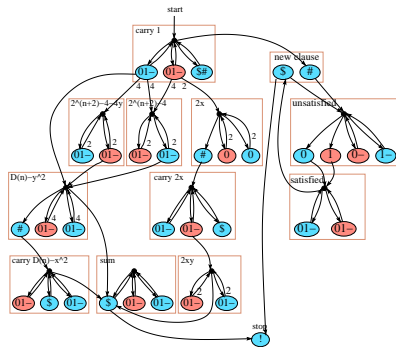
NAŠA PÔVODNÁ MOTIVÁCIA



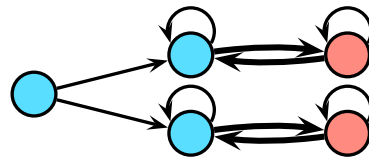
ZHRNUTIE

Všetky HMM

NP-ťažké



?



Viterbi alg. $O(nm^2)$



RVA $O(n^2m^3)$



ešte rozšír $O(n^3m^5)$

