

Introduction to Bioinformatics

JTMS-19

Marc-Thorsten Hütt

mhuett@constructor.university

Felix Jonas

fjonas@constructor.university

What is this session about?

Hidden Markov models (HMMs) are introduced. Viterbi algorithm of HMMs is discussed. First ideas of posterior decoding are introduced.

How can you revise the material after the session?

Read Durbin et al. chapters 3.2, 3.3, 3.4

Read Baxevanis/Oullette pages 208 – 210

alternative reading: Hütt/Dehnert chapters 2.8.1 – 2.8.4

Markov chains as a tool for studying CpG islands

+	A	C	G	T	-	A	C	G	T
A	0.180	0.274	0.426	0.120	A	0.300	0.205	0.285	0.210
C	0.171	0.368	0.274	0.188	C	0.322	0.298	0.078	0.302
G	0.161	0.339	0.375	0.125	G	0.248	0.246	0.298	0.208
T	0.079	0.355	0.384	0.182	T	0.177	0.239	0.292	0.292

$$S(x) = \log \left(\frac{P(x | \text{model } +)}{P(x | \text{model } -)} \right) = \log \left(\frac{P(B) \prod_{i=1}^L a_{x_{i-1}x_i}^+}{P(B) \prod_{i=1}^L a_{x_{i-1}x_i}^-} \right) = \sum_{i=1}^L \log \left(\frac{a_{x_{i-1}x_i}^+}{a_{x_{i-1}x_i}^-} \right) = \sum_{i=1}^L \beta_{x_{i-1}x_i}$$

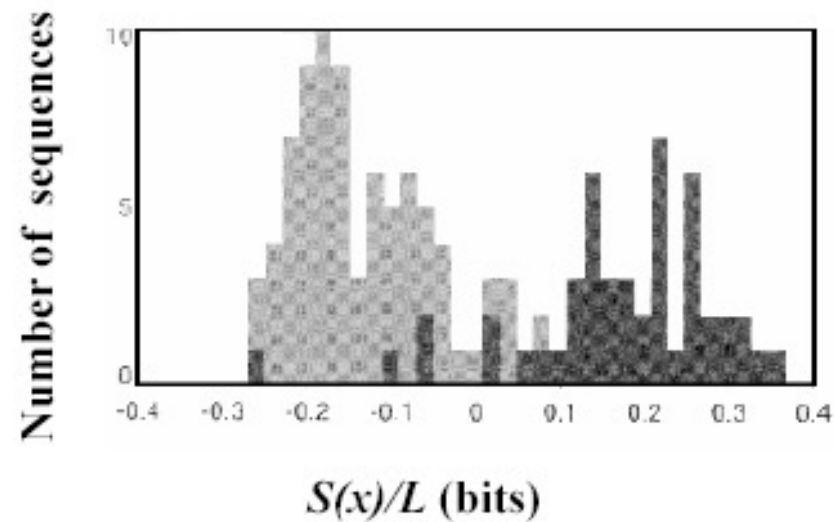
a number for each sequence x
 → histogram of score values S(x)
 for many sequences x

a number for each dinucleotide
 → table of "log-likelihoods"

Markov chains as a tool for studying CpG islands

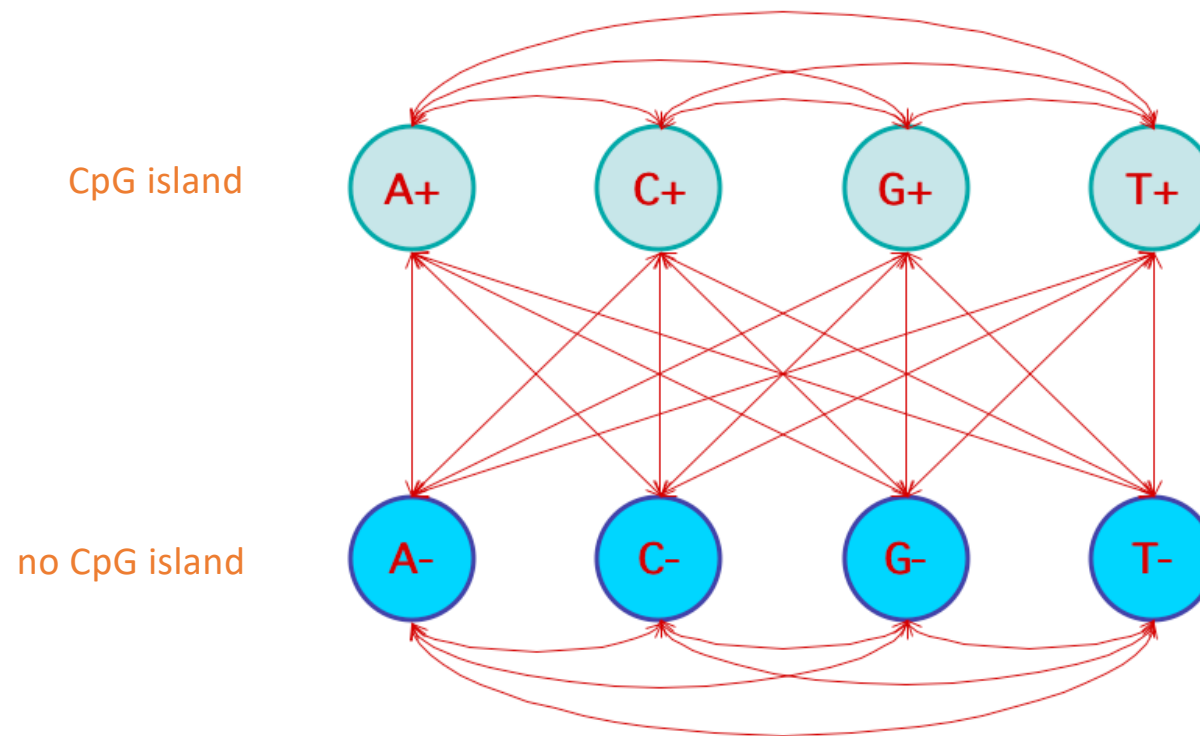
$\beta(\log_2)$	A	C	G	T
A	-0.740	0.419	0.580	-0.803
C	-0.913	0.302	1.812	-0.0685
G	-0.624	0.461	0.331	-0.730
T	-1.169	0.573	0.393	-0.679

table of
"log-likelihoods"



histogram of scores

Hidden Markov model (HMM) for CpG islands

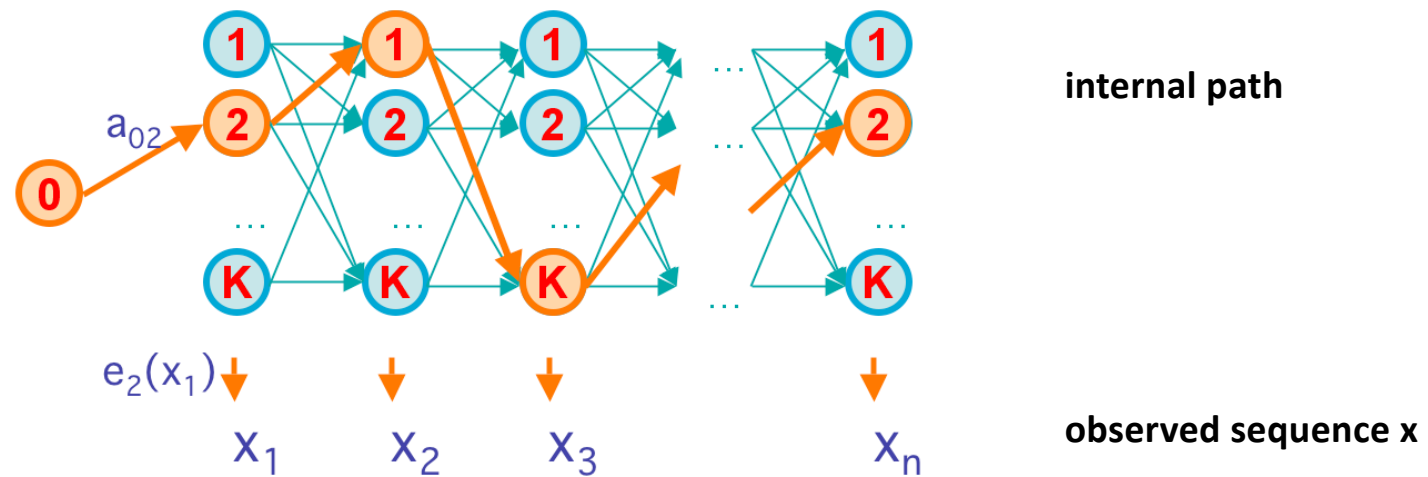


internal state: $A+, C+, G+, T+, A-, C-, G-, T-$

emitted state: A C G T A C G T

general properties of HMM

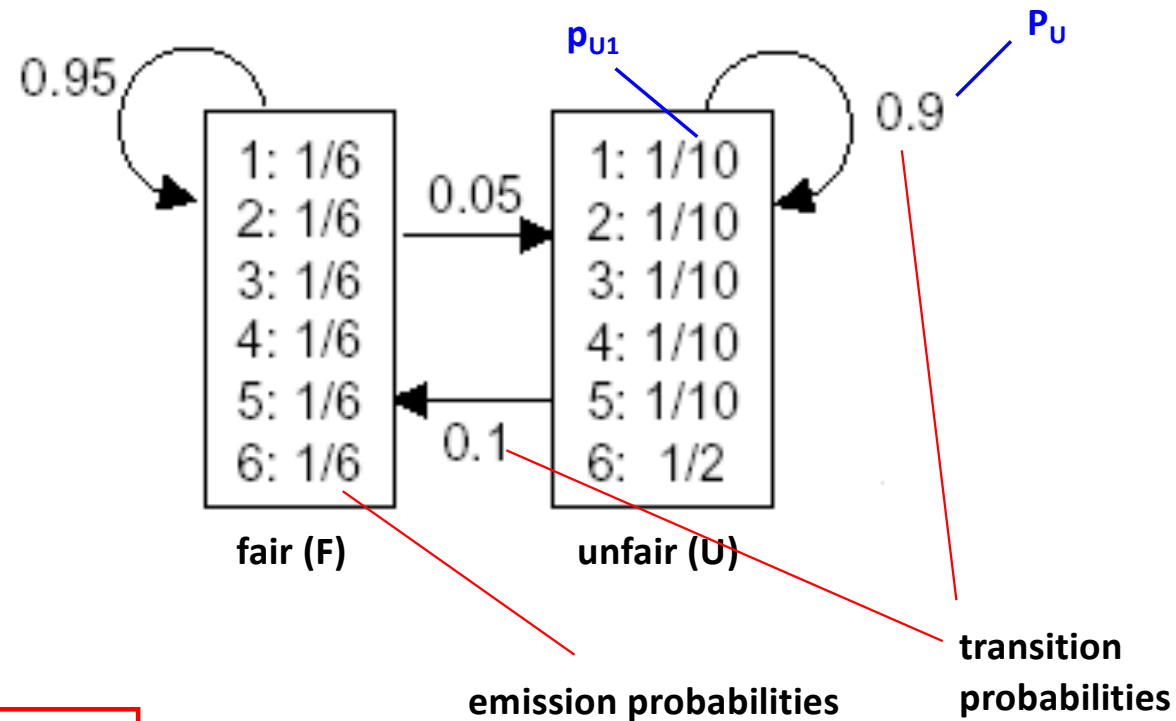
simulation of a sequence



joint probability of the path and the sequence

$$P(x, \pi) = a_{\pi_0 \pi_1} \prod_{i=1}^L e_{\pi_i}(x_i) a_{\pi_i \pi_{i+1}}$$

elementary example: casino with two dice



$$S_{\text{HMM}} = \{ F, U \}$$

$$S_{\text{HMM}} = \{ 1F, 2F, \dots, 6F, 1U, 2U, \dots, 6U \}$$

- emission prob's are 0 or 1
- transition prob's are products (e.g. $a_{2U,1U} = P_U p_{U1}$)

Training:

parameter estimation from an ensemble of sequences with a given internal structure

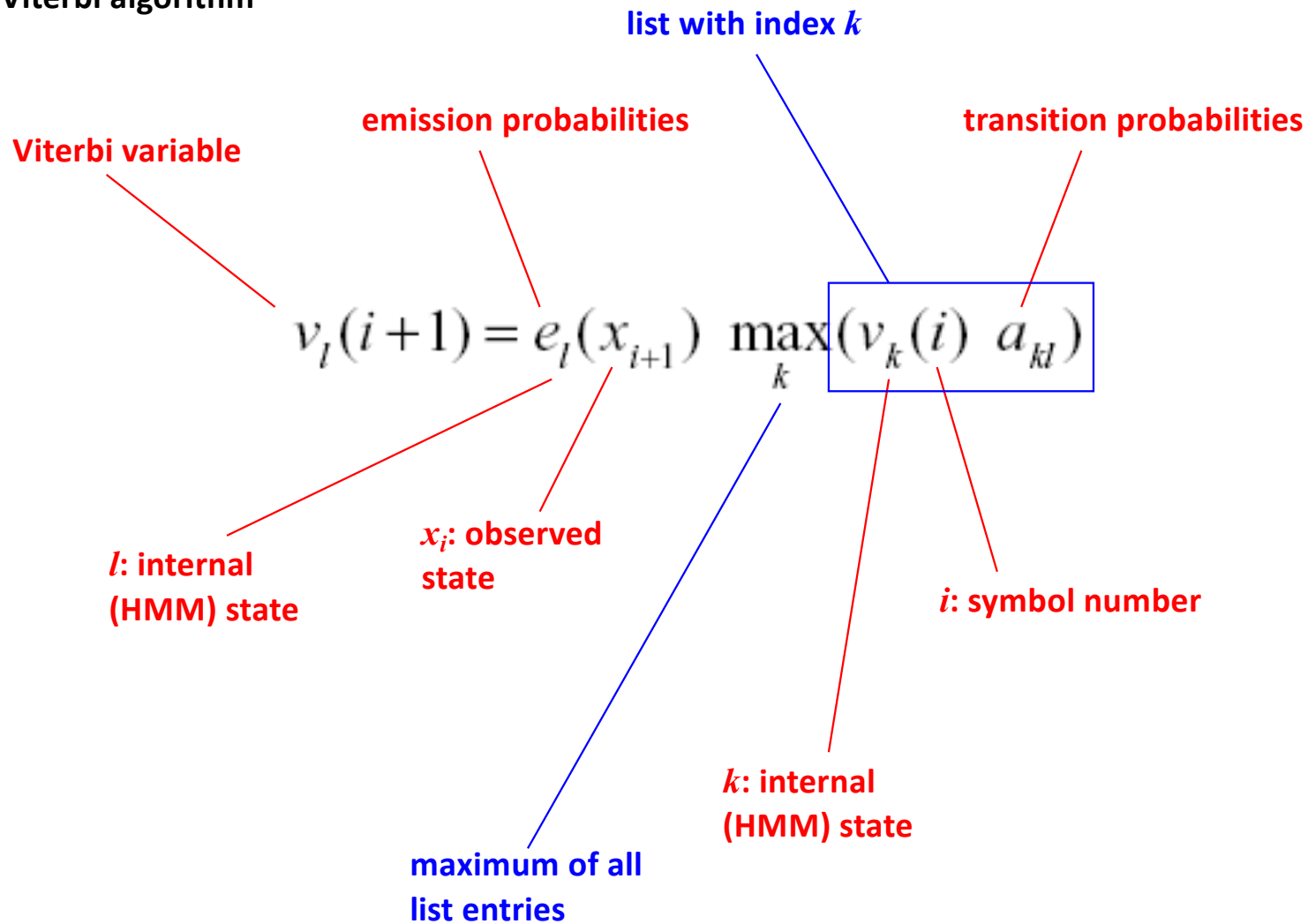
Evaluation:

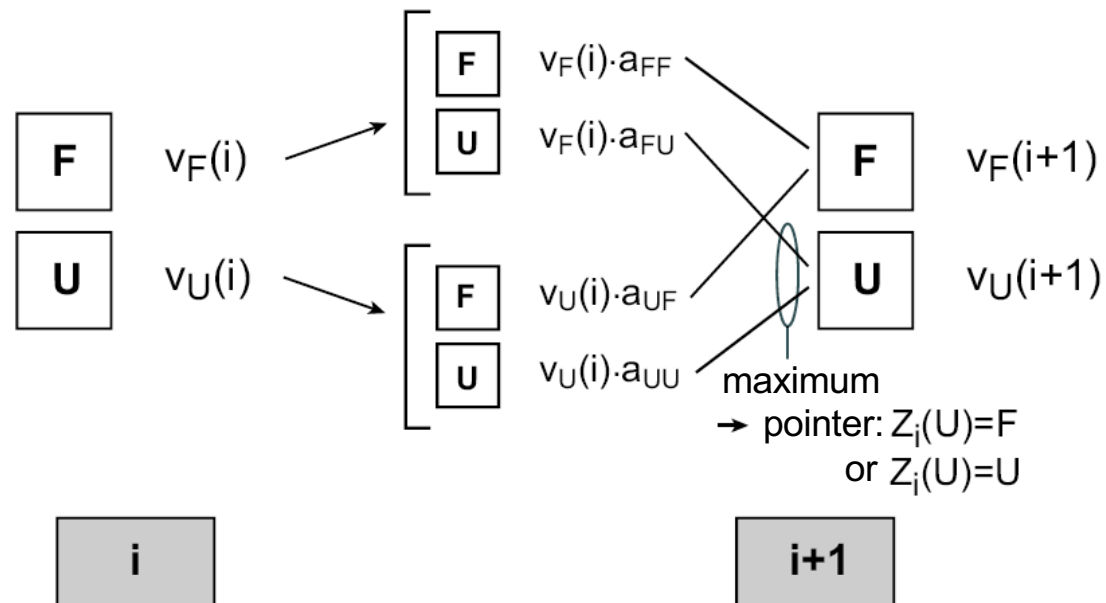
obtaining the sequence x out of a given internal sequence (path) p

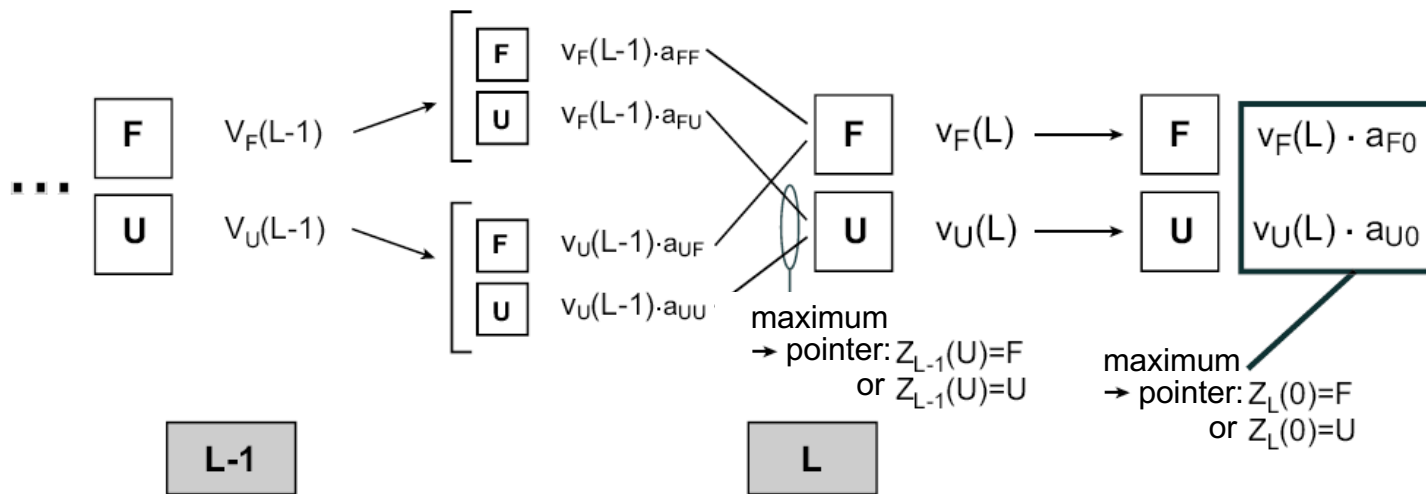
Decoding:

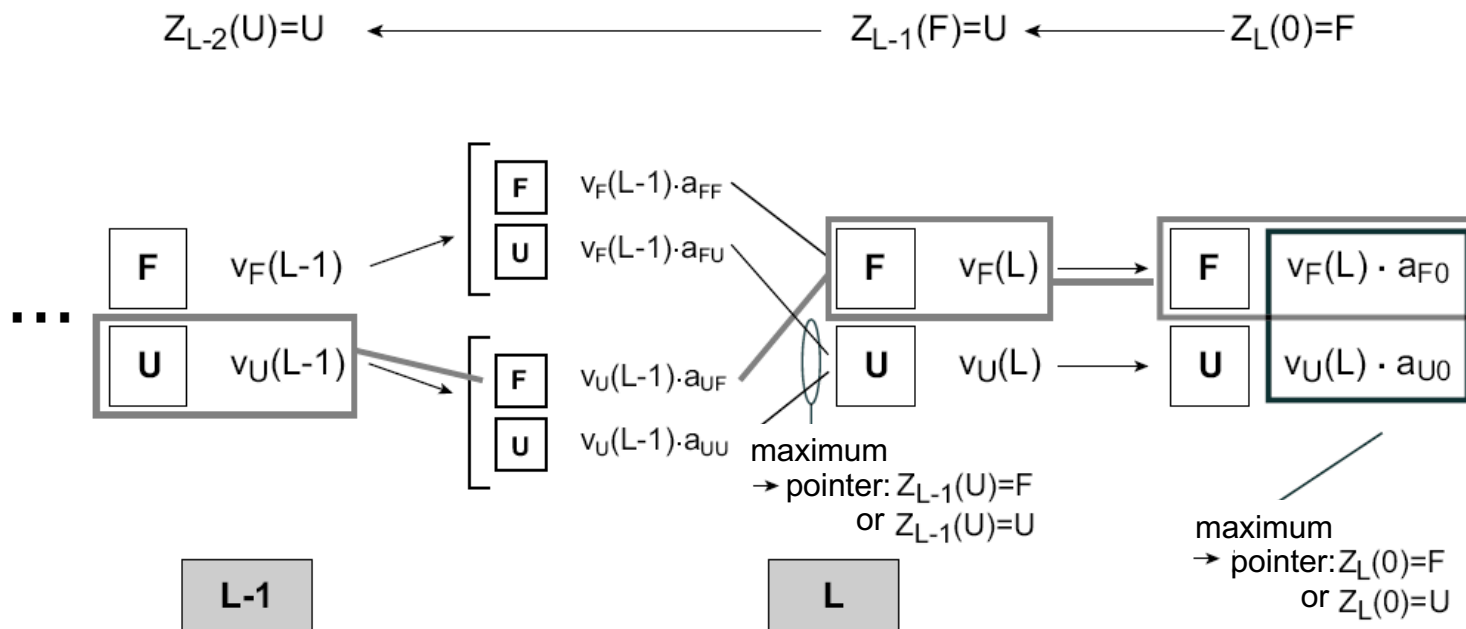
finding the internal path p behind an observed sequence x

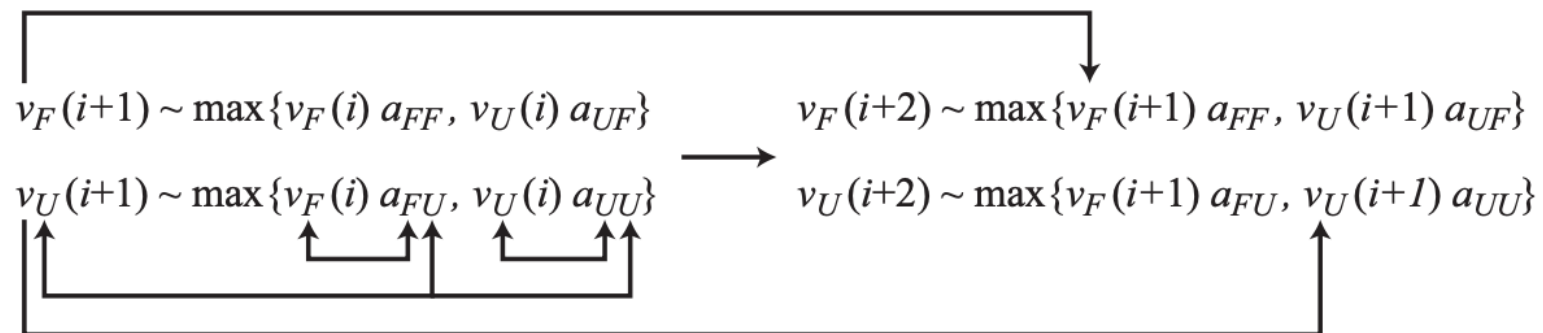
Viterbi algorithm











Casino: results

[illegible]

Finding the genes in genomic DNA

Christopher B Burge* and Samuel Karlin†

Addresses

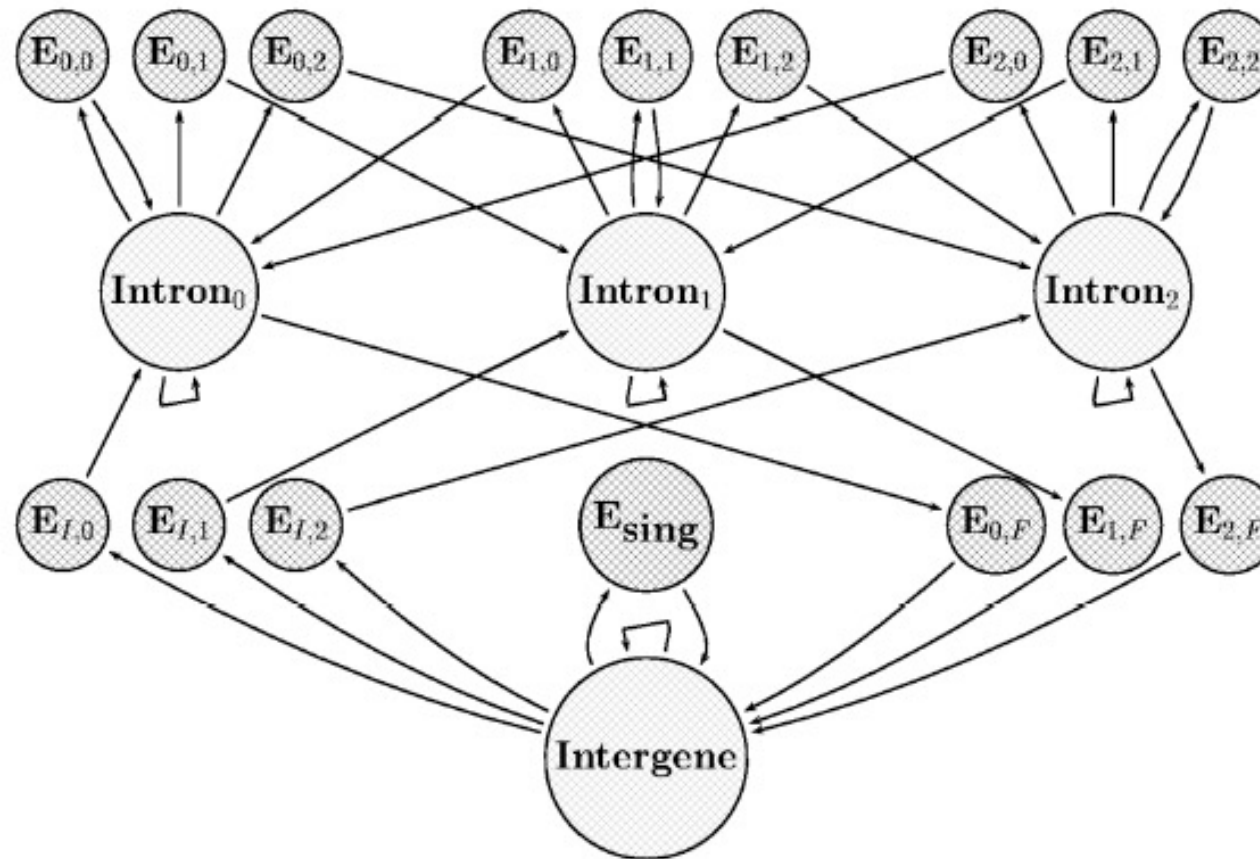
*Center for Cancer Research and Department of Biology,
Massachusetts Institute of Technology, 40 Ames Street, E17-526
Cambridge, MA 02139, USA; e-mail: cburge@mit.edu

†Department of Mathematics, Stanford University, 450 Serra Mall,
Stanford, CA 94305, USA; e-mail: sam@galois.stanford.edu

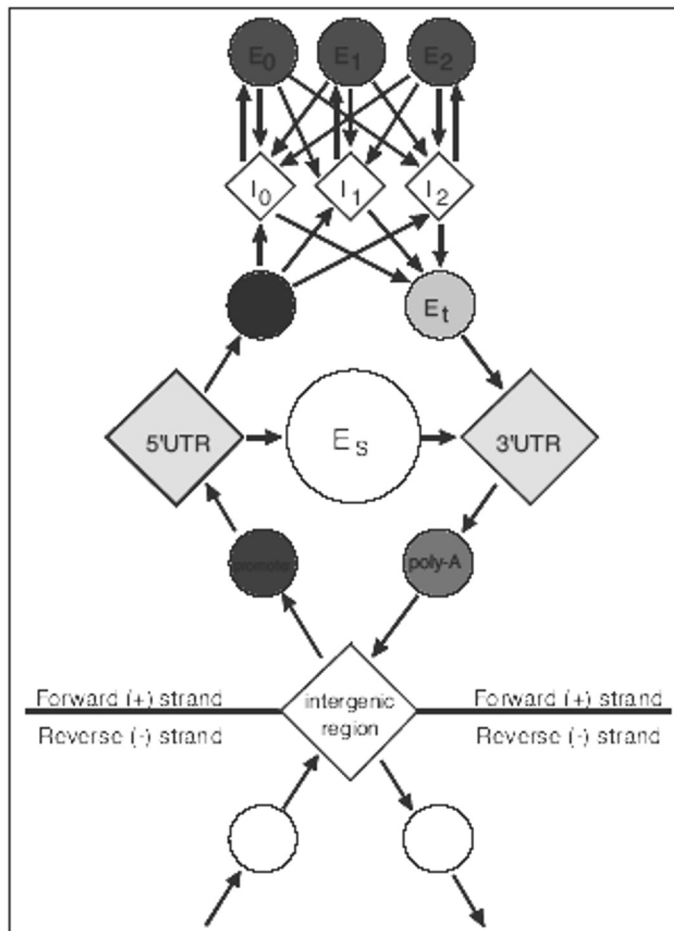
Correspondence: Samuel Karlin

Current Opinion In Structural Biology 1998, 8:346–354

The GenScan HMM

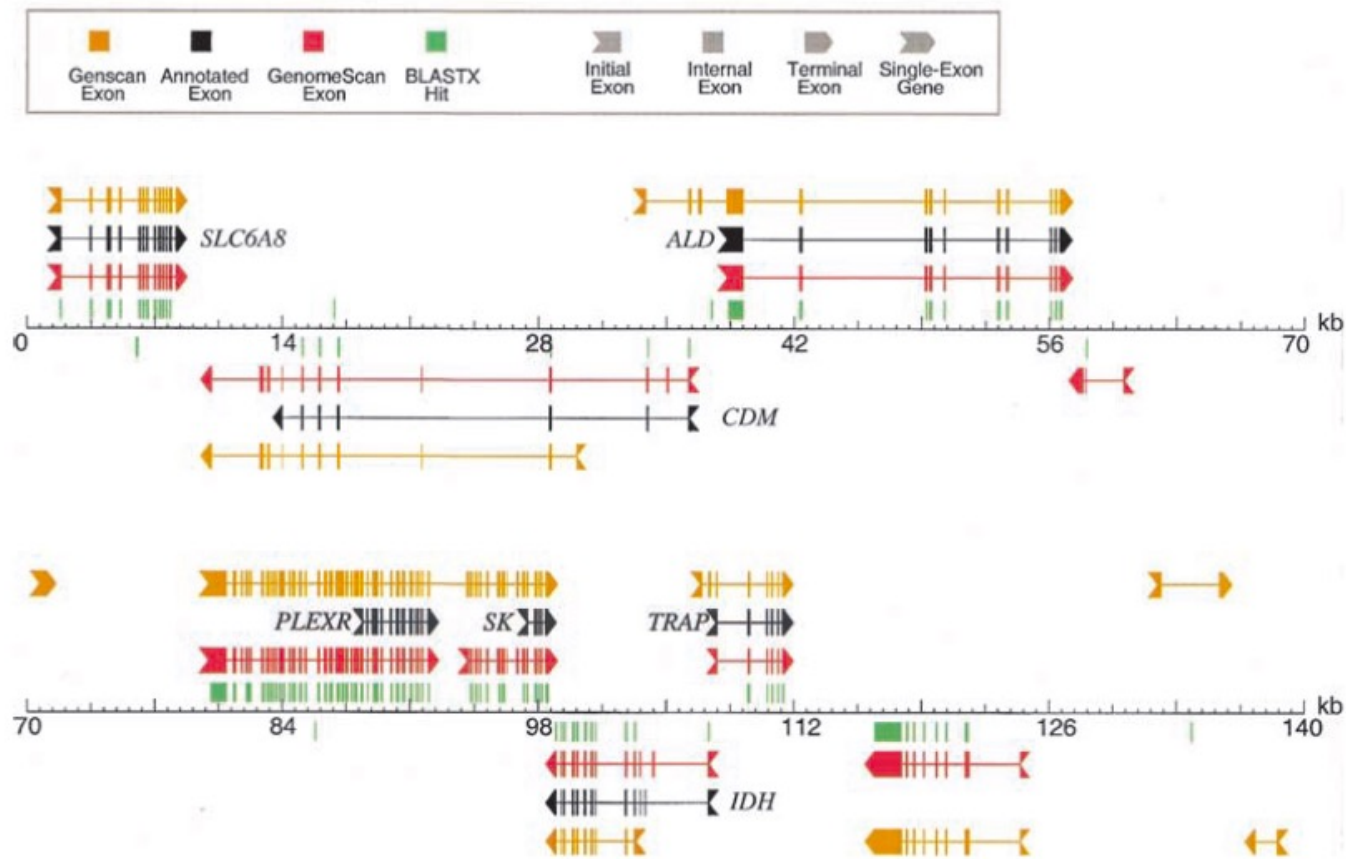


GENSCAN (Burge & Karlin)



```

62001 AGGACAGCTA GGGCTGTCTT CACTTACACC TCAGGCTGTG GAGCCAGACC
62051 CTAGGGTTGG GCAATCTACT CCAGGAGCA GGGAGGGCAG GAGCCAGGGC
62101 TGGGCATAAA AGTCAGGCA GAGGCATCTA TTGCTTACTT TTGCTTCTGA
62151 CACACCTGTG TTCACTAGCA ACCTCAACCA GACAC
62201
62251 CT TGTATCAG GTTACAGAC
62301 AGCTTANGG AGCCAAATG AACTGGCCA TGTGAGACA GAGAGACTC
62351 TTGGTTTCT GATAGGACT GACTCTCTCT GCTATTGCT CTATTTTCC
62401 ACCTTACG TCTGGTGT CTAGGCTTGG ACCCAGAGT TCTTGAAGT
62451 CTTGGGAT CTTCCACT CAGTGTCT TGTGGCAG CTTAGGTT
62501 AGGCTCAGG CAGAAAGTG CTGGTGGCT TTATGATGG CTTGGCTCAG
62551 CTGACAGCC TCAGGGGCG CTTGGCACA CTGATGAGC TGGACTGTA
62601 CAGCTGAC GTGATCTG AGAATTCAG CTGATCTA TGGAGCCTT
62651 GATGTTTCT TTGCTTCT TTCTATGCT TATGTCAG TCTAGGAG
62701 GGCAGACTA AGGCTTACA GTTAGATG GAAACAGAC GATGATTC
62751 ATCAGTGG AGTCTCAG ATGCTTTAG TTCTTTTAT TTGCTTCA
62801 TAAATCTT TTCTTTGT TAACTCTG CTCTTTTT TTCTTCTC
62851 GCAATTTT ACTATATC TAAAGCCT AACTTGTG ATACAAAG
62901 GAAATCTC TCAGATAC TACTTACTT AAAAAAGC TTACACAGT
62951 CTGCTAGTA CATTACTAT TGAATATAT GTGCTCTT TTGATATC
63001 ATATCTCC TACTTATT TCTTTATT TAAATGAT CAAATCAT
63051 ATATATTT AGGCTTAA GTATATGT TAAATGAT TACATATT
63101 GACCAATCA GGTATTTT GATTTGTA TTTAAGAA TCTTTCTC
63151 TTTAATTA CTTTFTGT TATCTATT CAAATCTT GCTAATCT
63201 TTCTTTCA GCAATATG ATACATGTA TCATGCTCT TTGACCAT
63251 CTAAAGATA ACAGTATTA TTCTGGCT AGGCAATG CATTATCT
63301 GCAATATAT ATTTCTGT ATAAATGTA ACTGATGA GAGTTTCA
63351 ATGCTATA GCACTACA TCAGCTACC ATCTGCTT TATTTATG
63401 TTGGATAG GCTGATAT TCTGATCA AGTAGGCG TTGCTAT
63451 CATCTTCA CTTCTATCT TCTGCCCA CTCTGGGC AAGTGCTG
63501 TCTGTGCT GGCCTATC TTGCAAGG ATTCACCC ACCAGTCA
63551 GCTGCTATC AGAAGTGT GCTGTGTG GCTAATGCC TGGCCACAA
63601 GTATCTTA GCTGCTTC TTGCTGCA ATTTATTA AAGTTCTT
63651 TGTTCCTA GTCCACTC TAACTGGG GATTTATGA AGGCTTGA
63701 GATCTGGAT TCTGCC AT AAAAACTT TTTTCTAT CCAATGAT
    
```

Yeh, Lim, and Burge, *Genome Research* 11:803-816, 2001.

Casino: results

$$\left\{ \begin{array}{c} \boxed{\mathbf{x}} \\ \boxed{\pi} \end{array} \begin{array}{c} x_1 \ x_2 \ \dots \ x_i \ x_{i+1} \\ \updownarrow \\ \vdots \end{array} \right\} = \sum_{k \in \Sigma_{\text{HMM}}} \left\{ \begin{array}{c} \boxed{\mathbf{x}} \\ \boxed{\pi} \end{array} \begin{array}{c} x_1 \ x_2 \ \dots \ x_i \quad x_{i+1} \\ \quad \uparrow \\ \quad e_l(x_{i+1}) \\ k \xrightarrow{a_{kl}} \vdots \end{array} \right\}$$

► Summary of the Viterbi algorithm

initialization → recursion → termination → traceback

$$v_0(0) = 1, v_k(0) = 0 \quad \forall k \in \Sigma_{HMM} \quad \text{initialization}$$

$$v_l(i+1) = e_l(x_{i+1}) \max_k \{v_k(i) a_{kl}\} \quad \text{recursion}$$

$$\max_k \{v_k(L) a_{k0}\} = P(x, \pi^*) \quad \text{termination}$$

$$\downarrow \quad \swarrow \quad a_{k0} = \frac{1}{|\Sigma_{HMM}|} \quad \forall k \in \Sigma_{HMM}$$

$$\pi_L^*$$

$$\pi_{L-1}^* = Z_{L-1}(\pi_L^*)$$

⋮

$$\pi_{i-1}^* = Z_{i-1}(\pi_i^*)$$

⋮

traceback