

Text Algorithms

Lecture: PWM, HMM, SCFG, Transducers

Jaak Vilo
2008 fall

Jaak Vilo

MTAT.03.190 Text Algorithms

1

Probabilistic models and other extensions

- Position Weight Matrices PWM
 - Position specific scoring matrices
- Hidden Markov models (HMM)
- Stochastic Context Free Grammars (SCFG)
- Finite state transducers (FST)
- Attributed automata

- It is not possible to always define the problem using strings, regular expressions, or approximate matching with few mismatches.
- Sometimes probabilistic models are better suited
- And correspond better to real world. For example, the affinity between a protein and the DNA is the physical property.
- Biological systems are not discrete digital computers but rather analog ones

- Durbin, R., Eddy, S., Krogh, A., Mitchison Biological sequence analysis **Probabilistic models of proteins and nucleic acids** (Cambridge University Press, 1998)
- [Molecular Information Theory and the Theory of Molecular Machines](http://www.lecb.ncifcrf.gov/%7Eetoms/) (Tom Schneider)
- G. Z. Hertz, Gary D. Stormo: [Identifying DNA and protein patterns with statistically significant alignments of multiple sequences](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&list_uids=104878648&dopt=Abstract). Bioinformatics 15(7): 563577 (1999)

- A **position weight matrix (PWM)**, also called **position-specific weight matrix (PSWM)** or **position-specific scoring matrix (PSSM)**, is a commonly used representation of [motifs](#) (patterns) in biological sequences.

(Wikipedia)

- **See compression algorithms for Entropy and Information Content**
- Lets have 8 sequences with certain properties (e.g. a protein is known to bind to the respective DNA stretches)
- How to summarize that info? What is the true signal?
- $C_{i,j}$ nr of occurrences of character i in column j .

GATGAG
GATGGG
GCTGAG
GAAGAG
AATGAG
GATGAA
GTTGAG
GATGAT

• Counts: $C_{i,j}$

A	1	6	1	0	7	1
C	0	1	0	0	0	0
G	7	0	0	8	1	6
T	0	1	7	0	0	1

GATGAG
 GATGGG
 GCTGAG
 GAAGAG
 AATGAG
 GATGAA
 GTTGAG
 GATGAT

G A T G A G (consensus?)

Match

A	G	A	T	G	A	G	T	G	
1	0	1	0	1	1			= 4	
	7	6	7	8	7	6		= 41	
		1	0	0	0	1	1	= 3	
			0	0	1	8	0	6	= 15

A	1	6	1	0	7	1
C	0	1	0	0	0	0
G	7	0	0	8	1	6
T	0	1	7	0	0	1

GCTGAA = 7 1 7 8 7 1 = 31

A	1	6	1	0	7	1
C	0	1	0	0	0	0
G	7	0	0	8	1	6
T	0	1	7	0	0	1

G A T G A G (consensus)

Does consensus represent the information? No.

Consensus with one mismatch? no

But mismatch can occur at an unlikely position then.

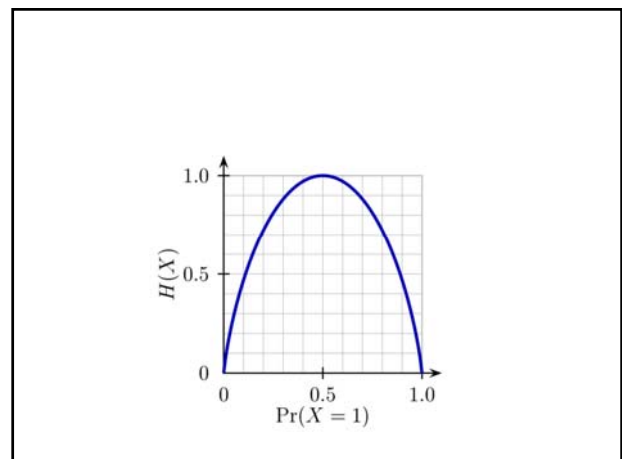
- $p_{i,j}$ probability of letter i at position j
- $p_{i,j} = f_{i,j} / N$ (N = nr of sequences)

A	0.125	0.75	0.125	0	0.875	0.125
C	0	0.125	0	0	0	0
G	0.875	0	0	1	0.125	0.75
T	0	0.125	0.875	0	0	0.125

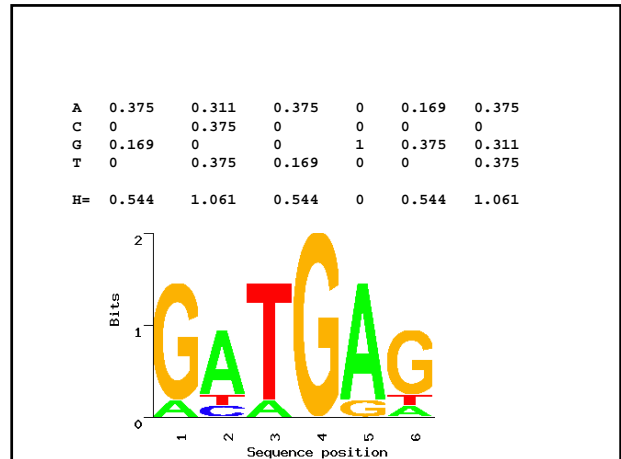
- Is 0 probability really 0?
- $p_{i,j} = (f_{i,j} + n_i) / (N+n)$ (n_i pseudocounts, sum n)
- Proportions are the same as in absolute numbers

Information content

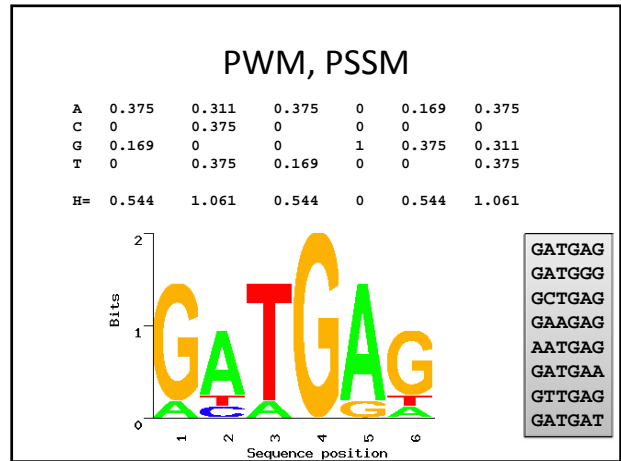
- At column i the sum of probabilities is 1: $\sum_{i=1..c} p_{i,j} = 1$
- Lets look at column j
- $\log_2 P_i$ Level of "surprise" in bits P_i
- If p is small, the more surprised we would be, the more information the position carries.
- Overall information content is counted over all possible characters

$$H = \sum_{i=1..c} p_i \log_2 p_i$$


A	0.375	0.311	0.375	0	0.169	0.375
C	0	0.375	0	0	0	0
G	0.169	0	0	1	0.375	0.311
T	0	0.375	0.169	0	0	0.375
H=	0.544	1.061	0.544	0	0.544	1.061



- Logo: Column height is 2 H
- Logo: Letter height is the proportion of the column: $p_i \cdot (2H)$
- Logo: Sort letters according to probability $p_i \cdot (2H)$



- [PATMATCH](http://www.bioinf.ebc.ee/EP/EP/)
- Try: 1:...TGAAAATTTT....
- -1:...TG.GATGAG.T....

PDF

http://bit.cs.ut.ee/~vilo/edu/2005-06/Text_Algorithms/10_Probabilistic/Articles/Hertz_Stormo_Bioinformatics_1999_Significance.pdf

G.Z. Hertz and C. Stormo

a) Alignment Matrix

A	A	T	T	C	A
A	C	C	T	C	C
A	C	A	T	C	
A	G	C	G	T	
1	2	3	4	5	6

consensus: A G C T G N

b) Weight Matrix

A	1.2	0	1.6	0	1.6	0
C	1.6	1.6	1.6	0	0	0
G	-1.6	1.6	1.6	-1.6	1.6	0
T	1.6	-1.6	0	1.6	0	0

test sequence: A G C T G C

$$\ln \frac{(m_{ij} + p_j) / (N + 1)}{p_j} \approx \ln \frac{f_{ij}}{p_j}$$

Fig. 1. Examples of the simple matrix model for summarizing a DNA alignment. (a) An alignment matrix describing the alignment of the four 6-mers on top. The matrix contains the number of times, m_{ij} , that letter i is observed at position j of this alignment. Below the matrix is the consensus sequence corresponding to the alignment (N indicates that there is no nucleotide preference). (b) A weight matrix derived from the alignment in (a). The formula used for transforming the alignment matrix to a weight matrix is shown above the arrow. In this formula, N is the total number of sequences (four in this example), p_j is the a priori probability of letter j (0.25 for all the bases in this example) and $f_{ij} = m_{ij}/N$ is the frequency of letter i at position j . The numbers enclosed in blocks are summed to give the overall score of the test sequence. The overall score is 4.5, which is also the maximum possible score with this weight matrix.

- Example from <http://www.biomedcentral.com/14712180/2/29>
- Observed frequencies (counted from the data)
- Odds score (ratio of observed over background frequencies)
- The odds score is the frequency observed divided by the theoretical frequency expected (i.e., the background frequency of the base, usually averaged over the genome ~0.25/base).
- In DNA a crude (and wrong) estimate is: $P_c = P_t = P_g = 1/4$
- If f is 0.79 and bkg prob. is 0.25, is odds score $0.79/0.25 = 3.16$.
- Finally, odds scores were converted to log odds scores by taking the logarithm base 2.
- $W_{ij} = \log_2(F_{ij}/P_j)$ W_{ij} is score for character i in column j .
- As the logarithm of zero is infinity, a zero occurrence of a particular base in the matrix creates a problem.
- A formula proposed by Hertz and Stormo can be used
- $W_{ij} = \log_2(C_{ij} + p_j) / (N + 1)P_j$
- (Which is approximately $= \log_2(F_{ij}/P_j)$)
- Where C_{ij} is nr of occurrences of i 'th character in position j and N is the nr. of sequences

- The statistical significance of such alignment can be calculated
- $I_{seqalign} = \sum_i \sum_j f_{ij} \log (f_{ij} / p_{ij})$
- The information content (aka normaliseeritud logodds ratio) Other names include KullbackLeibler information (Kullback, Leibler, 1951), relative entropy. Or, *large deviation rate function* (Bucklew, 1990).
- Formula has properties which satisfy the intuition about the information content
- It is for example a distance measure from the center of the distribution with $f_{ij} = p_{ij}$.
- Distance is minimal when $f_{ij} = p_{ij}$
- Distance is maximal when the least expected character is exclusive, or $f_{m,j} = 1$ ja $p_m \leq p_j$ for every i .
- Schneider et.al. (1986) noticed that $2^{I_{seqalign}}$ is approximately equal to frequency of the occurrences of DNA binding sites in *e.coli*
- Stormo et.al: $2^{I_{seqalign}}$ is the upper limit for the expectation with which the words from the alignment occur in the random sequence

- Log likelihood

Perceptron

- PWM is like a simple neural network, the 1node perceptron
- $F(X=x_1x_2..x_m) = \sum w_i x_i$
- DNA signals can also be represented in more complex neural networks (multilayer)

Matching of PWM

- Alphabet Σ and PWM ($W_{i,j}$) where for a character i at column j there is a weight $W_{i,j}$
- Like brute force exact match

Algorithm: Search with PWM
 Input: PWM ($W_{i,j}$), $i \in 1..|\Sigma|$, $j \in 1..m$, string S , threshold K
 Output: All occurrences of $W_{i,j}$ in S with quality at least K

```

for p=1 .. |S| m+1
  sum = 0 ;
  for j=1 .. m
    sum += W[ S[p+j-1] ][ j ]
  if sum ≥ K then report match at position p with score sum

```

Matching of PWM

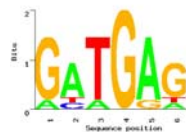
- Can it be improved when K is high? Stop when K cannot be reached
- Use some precalculated scores?
- Aho Corasick?

Some local research

- Marek Zäuram – BSc
- Jelena Zaitseva, Pavlos Pavlidis – PWM matching with
 - indels (e.g like edit distance – add/remove)
 - and closures (repeat the last column...)
- [PATMATCH](#)
- Still room for completing some of that research

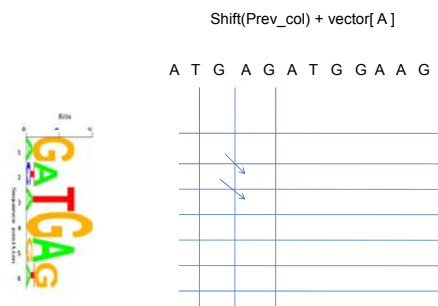
Matches above threshold t

1. Match scores of potentially most significant positions
2. If remaining best possible scores can not achieve the threshold, stop



E.g. You can afford only Max- t "mistakes"

Matching v2



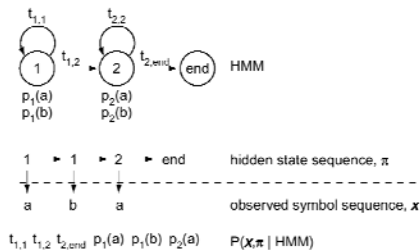
- Motif combinations
- Indels, closures
- What is the "best signal" ?

Research is going on...

- Cinzia Pizzi, Esko Ukkonen: **Fast profile matching algorithms — A survey** [Theoretical Computer Science Volume 395, Issues 23](#), 1 May 2008, Pages 137157
- GAPWM: a genetic algorithm method for optimizing a position weight matrix Li, Liang, Bass (2007) <http://bioinformatics.oxfordjournals.org/cgi/content/short/23/10/1188>
- A Feature-Based Approach to Modeling Protein–DNA Interactions
Eilon Sharon,^{#1} Shai Lubliner,^{#1} and Eran Segal^{1,2*}
PLoS Comput Biol. 2008 August; 4(8): e1000154. doi:10.1371/journal.pcbi.1000154.
<http://www.pubmedcentral.nih.gov/articlerender.fcgi?tool=pubmed&pubmedid=18725950>

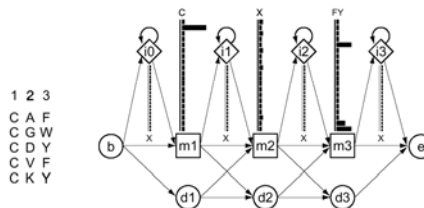
HMM (see also previous notes)

- **Introductory material on hidden Markov models**
- ["Profile hidden Markov models"](#) Sean R. Eddy, *Bioinformatics* 14(9):75563, 1998. (PDF)
 http://hlt.cs.ut.ee/~57Eviho/edu/2005-06/text_Algorithms/L10_Probabilistic/Articles/hmmreview-bioinformatics-98.pdf
- ["Hidden Markov Models and Protein Sequence Analysis"](#) Rachel Karchin, UCSC.



- Given an HMM, what is the p that the state has been passed through and respective sequence output?
- Given the sequences, generate the HMM that maximises the probability that these sequences are generated by that HMM
- Examples generate the model, model generates the examples
- Describe a family of sequences with a single HMM (e.g. protein family)
- For a new sequence what is the probability to be from the same family?

More complex example
 Proteins, alphabet of size 20
 Would produce 3/character long sequences
 Plus states for insertion and deletion



In state m1 the most probable is C, in state m2 any, and m3 F, Y and W.

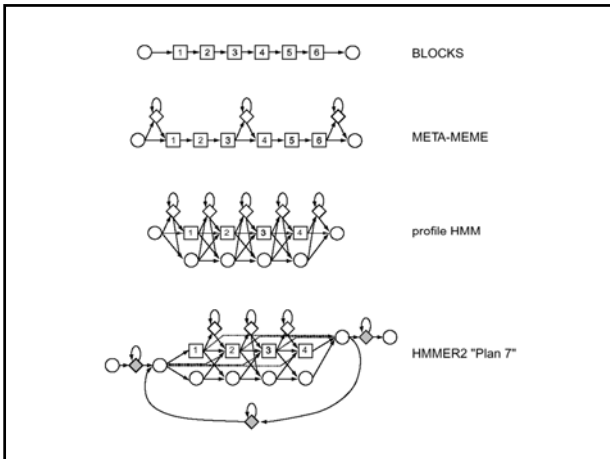
Using hidden Markov models

http://en.wikipedia.org/wiki/Hidden_Markov_model

- There are three canonical problems associated with HMM:**
- Given the parameters of the model, compute the probability of a particular output sequence, and the probabilities of the hidden state values given that output sequence. This problem is solved by the [forwardbackward algorithm](#).
 - Given the parameters of the model, find the most likely sequence of hidden states that could have generated a given output sequence. This problem is solved by the [Viterbi algorithm](#).
 - Given an output sequence or a set of such sequences, find the most likely set of state transition and output probabilities. In other words, discover the parameters of the HMM given a dataset of sequences. This problem is solved by the [BaumWelch algorithm](#).

HMM training

- Different representations, topologies
- Usually the training algorithm learns the parameters (weights), not the topology
- It is up to the designer to come up with good topology.

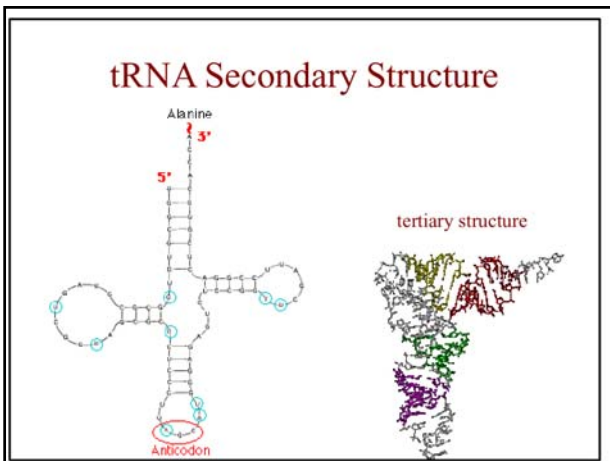
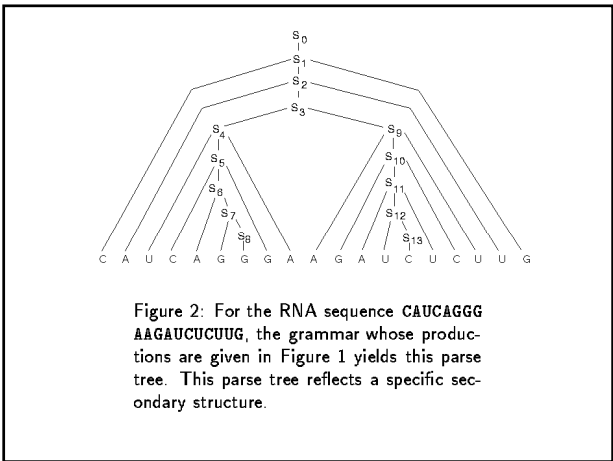


Stochastic Context Free Grammars

- Stohhastilised kontekstivabad grammatikad
- [Articles](http://bit.cs.ut.ee/%7Eviljo/edu/2005-06/Text_Algorithms/L10_Probabilistic/Articles/)
- Sakakibara et.al. 1993

$$P = \left\{ \begin{array}{ll} S_0 \rightarrow S_1, & S_7 \rightarrow G S_8, \\ S_1 \rightarrow C S_2 G, & S_8 \rightarrow G, \\ S_1 \rightarrow A S_2 U, & S_8 \rightarrow U, \\ S_2 \rightarrow A S_3 U, & S_9 \rightarrow A S_{10} U, \\ S_3 \rightarrow S_4 S_5, & S_{10} \rightarrow C S_{11} G, \\ S_4 \rightarrow U S_5 A, & S_{10} \rightarrow G S_{11} C, \\ S_5 \rightarrow C S_6 G, & S_{11} \rightarrow A S_{12} U, \\ S_6 \rightarrow A S_7, & S_{12} \rightarrow U S_{13}, \\ S_7 \rightarrow U S_7, & S_{13} \rightarrow C \end{array} \right\}$$

Figure 1: This set of productions P generates RNA sequences with a certain restricted structure, such as that in Figure 2. S_0, S_1, \dots, S_{13} are nonterminals; A, U, G and C are terminals representing the four nucleotides.

$$S_0 \Rightarrow S_1 \Rightarrow C S_2 G \Rightarrow C A S_3 U G \Rightarrow C A S_4 S_5 U G \Rightarrow C A U S_6 A S_7 U G \Rightarrow C A U C S_9 G A S_{10} U G \Rightarrow C A U C A S_7 G A S_9 U G \Rightarrow C A U C A G S_8 G A S_9 U G \Rightarrow C A U C A G G G A S_9 U G \Rightarrow C A U C A G G G A A S_{10} U U G \Rightarrow C A U C A G G G A A G S_{11} C U U G \Rightarrow C A U C A G G G A A G A S_{12} U C U U G \Rightarrow C A U C A G G G A A G A U S_{13} U C U U G \Rightarrow C A U C A G G G A A G A U C U C U U G.$$


$$\text{Prob}(s | G) = \sum_{\text{all derivations (or parse trees) } \delta} \text{Prob}(S_0 \xrightarrow{\delta} s | G)$$

$$= \sum_{\alpha_1, \dots, \alpha_n} \text{Prob}(S_0 \Rightarrow \alpha_1 | G) \cdot \text{Prob}(\alpha_1 \Rightarrow \alpha_2 | G) \cdot \dots \cdot \text{Prob}(\alpha_n \Rightarrow s | G)$$

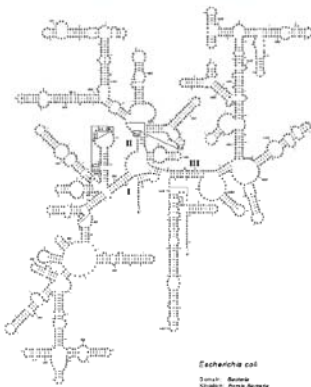
Efficiently computing this quantity, $\text{Prob}(s | G)$, presents a problem because the number of possible parse trees for s is exponential in the length of the sequence. However, a dynamic programming technique analogous to the Cocke-Kasami-Young or Earley methods [1] for non-stochastic CFGs can accomplish this task efficiently (in time proportional to the cube of the length of s). We define the negative logarithm of the probability of a sequence given by the grammar, $-\log(\text{Prob}(s | G))$, as the *negative log likelihood (NLL) score* of the sequence. This quantifies how well the sequence s fits the grammar.

Type I	$S \rightarrow SS$	1.0		
Type II	$S \rightarrow A S A$	0.05	$S \rightarrow A S C$	0.05
	$S \rightarrow A S G$	0.05	$S \rightarrow A S U$	0.1
	$S \rightarrow C S A$	0.05	$S \rightarrow C S C$	0.05
	$S \rightarrow C S G$	0.1	$S \rightarrow C S U$	0.05
	$S \rightarrow G S A$	0.05	$S \rightarrow G S C$	0.05
	$S \rightarrow G S G$	0.1	$S \rightarrow G S U$	0.05
Type III	$S \rightarrow A S$	0.2*	$S \rightarrow G S$	0.2*
	$S \rightarrow C S$	0.2*	$S \rightarrow U S$	0.2*
	$S \rightarrow S$	0.2*		
Type IV	$S \rightarrow A$	0.25	$S \rightarrow G$	0.25
	$S \rightarrow C$	0.25	$S \rightarrow U$	0.25

Figure 3: To obtain production probabilities for this initial grammar, we placed a uniform distribution over each set of same type productions, but weighted Watson-Crick base pairs twice as heavily. Starred values may differ if no skip productions are needed. For simplicity, this figure omits subscripts.

- Proc Int Conf Intell Syst Mol Biol. 2000;8:5766.Links
- **Small subunit ribosomal RNA modeling using stochastic contextfree grammars.**
- [Brown MP.](#)
- HNC, San Diego, CA 921213278, USA. mpsb@hnc.com
- We introduce a model based on stochastic contextfree grammars (SCFGs) that can construct small subunit ribosomal RNA (SSU rRNA) multiple alignments. The method takes into account both primary sequence and secondary structure basepairing interactions. We show that this method produces multiple alignments of quality close to hand edited ones and outperforms several other methods. We also introduce a method of SCFG constraints that dramatically reduces the required computer resources needed to effectively use SCFGs on large problems such as SSU rRNA. Without such constraints, the required computer resources are infeasible for most computers. This work has applications to fields such as phylogenetic tree construction.
- <http://www.ncbi.nlm.nih.gov/pubmed/10977066>

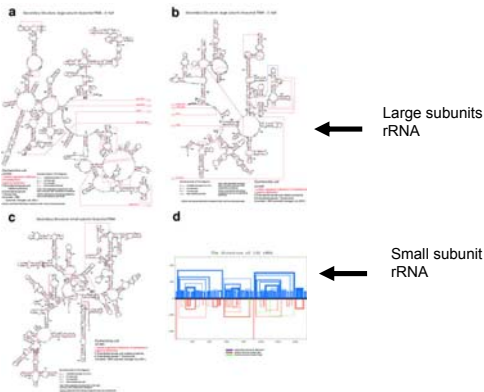
Small Subunit Ribosomal RNA



```

UGGAGAGUUGAUCCUGGCUACAGGUGAAGCUGGGCG  CUUGCCUAAAGACADCAAGUCUGUGGGCCCGGGGUA
CUCGUGUUCAGCGGACGAGGUGAGUAACGCGUGGG  UGACCCUACCCGGAAGAGGGGACACACCGGGGAAACUC
GGGCUAAUCCCCAUGUGAACCCGCCUUCUGGGUGUG  UCCAAAGGGCUUUGCCCGCUUCGGAUGGGCCCGGUC
CCACGACUAGUUGGUGGUAUUGCCACCAAGGCG  ACACAGGGUAGCCGUCUGAGAGGAUGCCGCTCACAG
GGGCAUUGAGCAACCGCCACUUCUACGGAGGCGAG  CAGUUAUGAUCUUCUCCAUUGGGGAAACCTUGGCG
AGCAAGCCGCUUGGAGGAAAGAGCCUUCGAGGUGUA  AACUCUGAACCCGGAGCAAAACCCCGACGAGGGGAC
UGACGGUACCGGGUUAUAGCCCGCCAAUCGCGUUC  CAGCAGCCCGGUAUACGGAGGGGCGGAGCGUUAACC
GGUAAAGGCGUUGAGGCGCCUGGGGCGUCCAGUG  GAAGAACCAGGCUAACCGUGGGGAGCGUGGUAUC
GUCAGCCUAGACGGUGGAGGGGUGGUAUUCUCC  GGGUACCGGUGAUAUCGCGUUAACCGGAGGACCC
CGAUGGGGAGGCGCCACUGGUCACCCGUGAGCCU  GAGGCGGAAAGCGUGGGGACAAACCGGUAUAACCG
GGUAGUCCAGCCUAAAGUAGCGCGUAGGUCUCUG  GGUCUCUGGGGGCCGAGCTUAAACGCUUAAGCGCC
GCTUGGGAGUACGGCCGACAGGCGUAAACUAAAGGA  AUUGACGGGGCCCGCACAGCGGGGAGCAUUGGUGU
UAATUCGCGAAGAACUUAACCGCCUUCACAGUUA  GGAACCCGGUGAAGAGCCUGGGGCGCCGCAAGGGAC
CCTAGCACAGGUCUGGGCCUUCGUCGUGGCGCU  GAGUGUUGGGUUAAGUCCCGAAGAGCGCACCCCC
GCCUUAUGUCCAGCGUUCGGCCGGCACUUAACG  GGACUCCCGGAAAGCGGGAGGAGGGAGGAGCG
UCUGGUCAGCAGGCCUUAAGCGCCUGGGGCGACACG  UGUCACAAUCCGCUAACAGGCGAAGCCCGGCAAG
GGGAGCUAUCGCAAAAGGUGGGCCGAGUUGAGU  GGUCUCACACCCGCAUAGAGCCGGAUUCGCUAG
UAUUCGCGAUCAGACCGCCGUGAUAUGUUCUCC  GGCCUUGUACACACCCCGUCUACCGCCUGGGAGCG
GUCUACCCGAGUUGCCGGAGGAGCUACGGGCAAGCC  CGAGGUAAGGGCCGUGACUGGGGGAAGUCGUAACA
GGUAGCUGUACCGAAGGUGCGUGGUAUACUCC
    
```

RNA secondary Structure representation



30S Small subunit of Ribosomal RNA

The **16S rRNA** is a 1542 nt long component of the small prokaryotic ribosomal subunit (**30S**) and has several functions

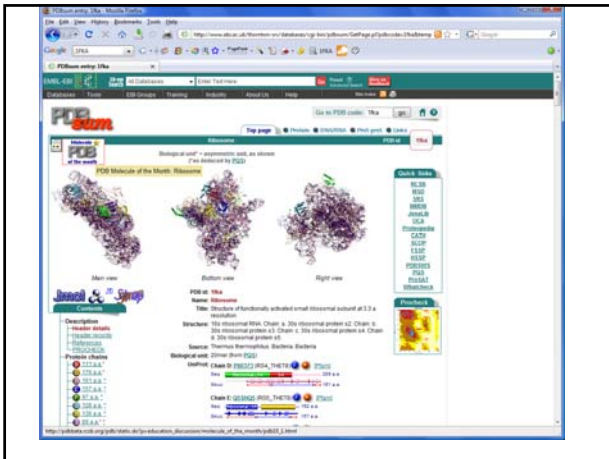
30S is the smaller subunit of the 70S **ribosome** of **prokaryotes**. It is a complex of ribosomal RNA and ribonucleoproteins which functions in **mRNA translation**. It includes the subunit **16S**, which is composed entirely of ribosomal RNA.

The 30S subunit is the site of inhibition for antibiotics such as **tetracycline**.

<http://en.wikipedia.org/wiki/30S>



Atomic structure of the 30S Subunit from *Thermus thermophilus*. Proteins are shown in blue and the single RNA strand in orange.^[1]



Examples of known interactions of RNA secondary structural elements

Pseudoknot

These patterns are excluded from the prediction schemes as their computation is too intensive.

Kissing hairpins

Hairpinbulge contact

- **Stochastic Context Free Grammars (SCFG)**
- [Loengumaterial, Mark Craven](#)
- **Main questions**
- How probable is the sequence for this grammar? the Inside algorithm
- What is the most probable parse tree? the CockeYoungerKasami (CYK) algorithm
- How to learn the parameters from the examples? the InsideOutside algorithm

Finite state machines

	Acceptors	Transducers
Unweighted		
Weighted		

JHU 600.405 - Finite-State Methods in NLP - Fall 2000 - J. Eisner 6

Definition of FST

A FST is $(Q, \Sigma, \Gamma, I, F, \delta)$

- Q: a finite set of states
- Σ : a finite set of input symbols
- Γ : a finite set of output symbols
- I: the set of initial states
- F: the set of final states
- $\delta \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times (\Gamma \cup \{\epsilon\}) \times Q$: the transition relation between states.

→ FSA can be seen as a special case of FST

Fei Xia; courses.washington.edu/ling570/fei_fall07/10_10_FST.ppt

An example of composition operation

Fei Xia; courses.washington.edu/ling570/fei_fall07/10_10_FST.ppt

FST Algorithms

- **Recognition:** Is a given pair of strings accepted by an FST?
 - $(x,y) \rightarrow$ yes/no
- **Composition:** Given two FSTs T_1 and T_2 defining regular relations R_1 and R_2 , create the FST that computes the composition of R_1 and R_2 .
 - $R_1=\{(x,y)\}, R_2=\{(y,z)\} \rightarrow \{(x,z) \mid (x,y) \in R_1, (y,z) \in R_2\}$
- **Transduction:** given an input string and an FST, provide the output as defined by the regular relation?
 - $x \rightarrow y$

Fei Xia; courses.washington.edu/ling570/fei_fall07/10_10_FST.ppt

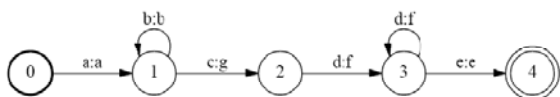
Finite-state transducers

A (2-way) finite-state transducer is a quintuple $M = (K, s, F, \Sigma \times \Sigma, \delta)$ where:

1. K is a finite set of states
2. s is a designated initial state
3. F is a designated set of final states
4. Σ is an alphabet of symbols, and
5. δ is a transition relation from $K \times (\Sigma \cup \epsilon \times \Sigma \cup \epsilon)$ to K

Richard Sproat: compling.ai.uiuc.edu/rws/newindex/acm.ppt⁵⁵

An FST



Richard Sproat: compling.ai.uiuc.edu/rws/newindex/acm.ppt⁵⁷

Composition

- In addition to *union*, *concatenation* and *Kleene closure*, regular relations are closed under *composition*
- Composition is to be understood here the same way as composition in algebra:
 - $R_1 \circ R_2$ means take the output of R_1 and feed it to the input of R_2

Richard Sproat: compling.ai.uiuc.edu/rws/newindex/acm.ppt⁵⁸

Composition: an illustration

Consider the following simple example using rewrite rules.

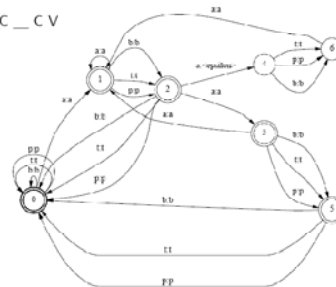
R1: $V \rightarrow \epsilon / V C _ C V$
 R2: $b \rightarrow p / _ t$

tabata (via R1) \rightarrow tabta (via R2) \rightarrow tapta
 As it happens, rewrite rules such as the above are implementable using FST's.

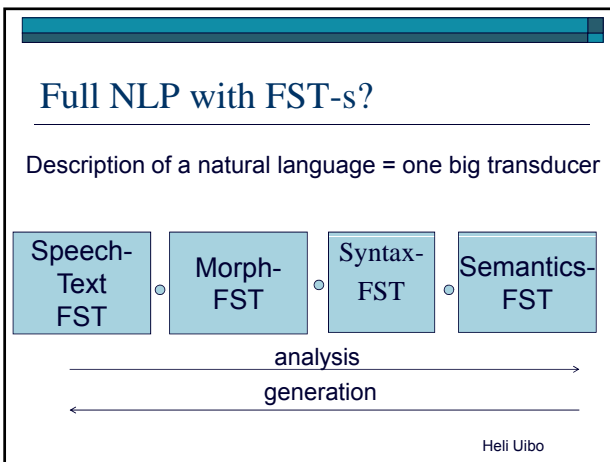
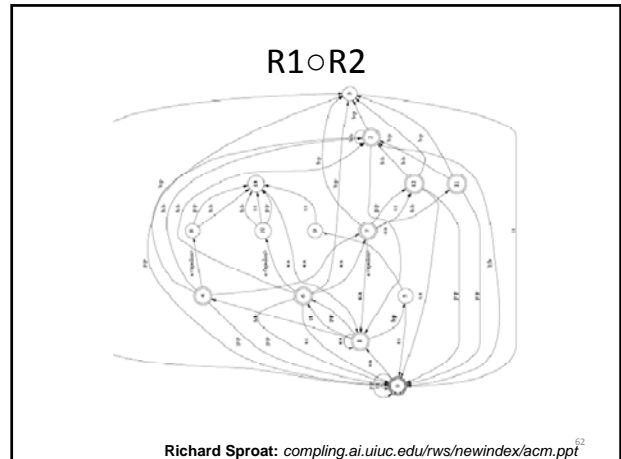
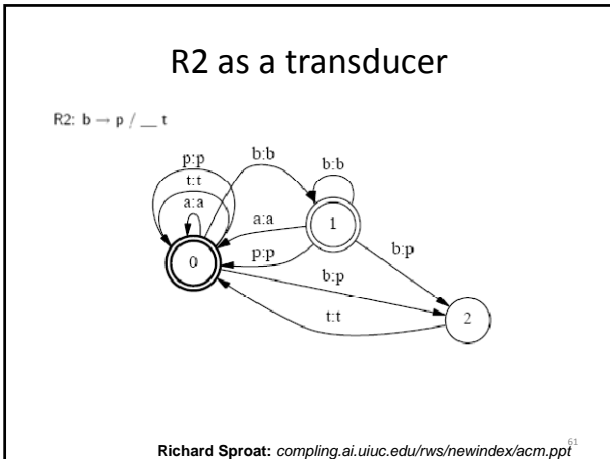
Richard Sproat: compling.ai.uiuc.edu/rws/newindex/acm.ppt⁵⁹

R1 as a transducer

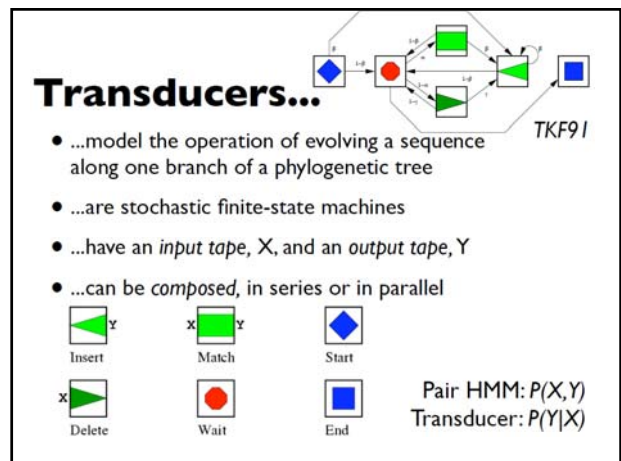
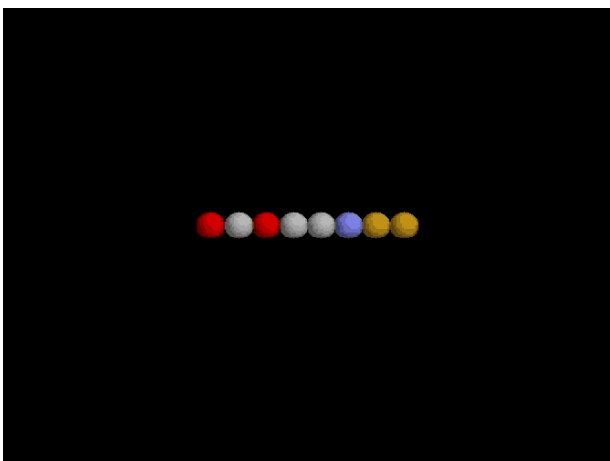
R1: $V \rightarrow \epsilon / V C _ C V$

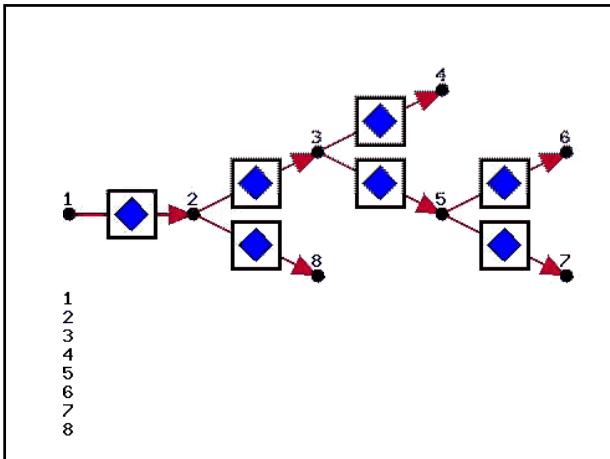


Richard Sproat: compling.ai.uiuc.edu/rws/newindex/acm.ppt⁶⁰



- ### Transducers
- Help to model the evolution
 - “Transform sequence into another”
 - Compose on different branches of evolutionary tree
 - <http://biowiki.org/PhyloFilm>
 - <http://biowiki.org/~yam/fruitfly/TKFmovie/>





“Every good work of software starts by scratching a developer’s personal itch” - Eric Raymond

2001 (Holmes & Bruno)
MCMC based on TKF91. **Poor performance due to affine gaps, rate variation**

2002 (Holmes & Rubin)
EM algorithm for estimating substitution rates & in particular rate variation

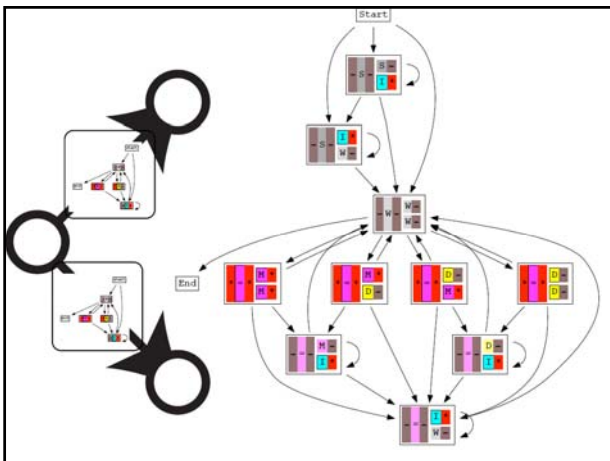
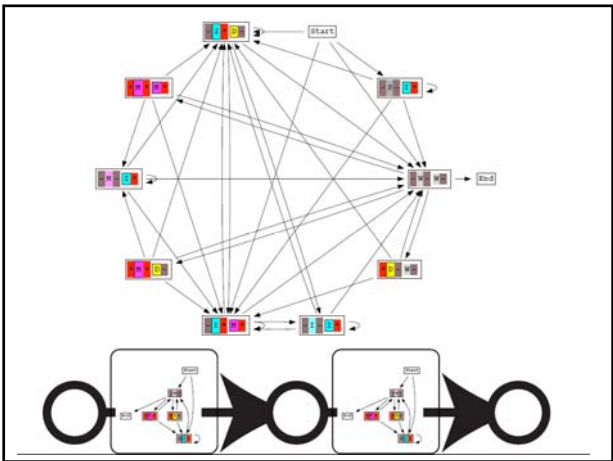
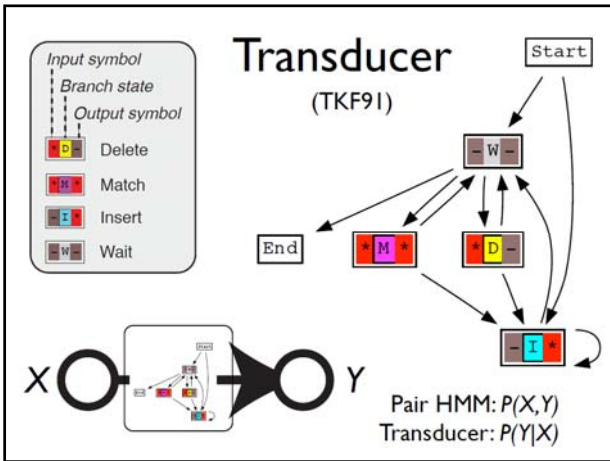
2003 (Holmes)
General algorithm for transducer composition on phylogenetic trees

2004 (Miklós, Lunter & Holmes)
“Long Indel” model: affine-gap Pair HMMs from evolutionary models

2005-2006 (Holmes; Klosterman et al)
(Started extending transducer theory to RNA sequence analysis & SCFGs)
Enormous difficulties debugging transducer composition & DP algorithms!

2007 (Holmes)
Phylocomposer: a transducer compositor, parser-generator & debugger

Ian Holmes (UC Berkeley)

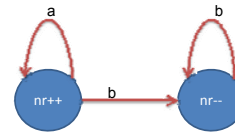


Video lectures

- http://videolectures.net/aop07_bertinoro/
- Mehryar Mohri, **Weighted Transducers and Rational Kernels**
– http://videolectures.net/aop07_mohri_wtt/
- **Similarity and differences by finite automata**
[Tamás Gaál](http://videolectures.net/aop07_gaal_sdf/), XEROX Research Centre Europe
– http://videolectures.net/aop07_gaal_sdf/



Attributed automata



$a^n b^n$ (if $nr=0$)

Graph Edit Distance

- Measuring the distance (or similarity) of graphs is an important task in pr and related areas.
- Graph edit distance (ged) is one of the most general graph distance measures.
- ged measures the distance of a pair of graphs in terms of the minimum number of edit operations required to transform one graph into the other one.

A Quadratic Programming Approach to the Graph Edit Distance Problem
June, 2007
2