

#### Alan Moses

Associate Professor and Canada Research Chair in Computational Biology Departments of Cell & Systems Biology, Computer Science, and Ecology & Evolutionary Biology Director, Collaborative Graduate Program in Genome Biology and Bioinformatics Center for Analysis of Genome Evolution and Function University of Toronto

### **Classic Bioinformatics topics!**

- consensus sequences
- regular expressions
- PWMs / PSSMs / matrix models
  - Motif scanning / matrix matching
  - *De novo* motif finding
- profile HMMs



Beginner textbook



#### Chapter 7 Regulatory Motif Analysis

Alan Moses and Saurabh Sinha

Advanced book chapter



#### Intermediate textbook

### Biology of *cis*-regulation

- Gene regulation at the level of transcription is a major factor in the 1) response to cellular environment
  - 2) diversity of cell types in complex organisms
  - 3) evolutionary diversity of morphology



Major Mechanism:

Sequence specific DNA binding proteins



Berman et al.: The Protein Data Bank. Nucleic Acids Research, 28 pp. 235-242 (2000).

#### AG**CGGAAAACTGTCCTCCG**TGCTTAA

What bases does this protein prefer to bind?

Bioinformatics. 2000 Jan;16(1):16-23. DNA binding sites: representation and discovery. Stormo GD YIR032CGATAAGYIR032CGGTAAGYIR032CGATAAGYJL110CGATAATYKR034WGATAGAYKR034WGATAACYKR039WGATAAGYKR039WGATAAC

8 GATA sites from SCPD database



## Probabilistic models of regulatory motifs (sequence families, consensus patterns)

Represent the sequence pattern quantitatively using a statistical model

YIR	032C	GATA	AAG			
YIR	032C	GGT	AAG			
YIR	032C	GATA	AAG			
YJL	110C	GATA	AAT	-		
YKF	R034W	GATAGA				
YKF	R034W	GATA	AAC			
YKF	R039W	GATAAG				
YKF	R039W	GATA	AAC			
	8 GATA s from SC	sites PD				

	A	С	G	Т					
1	0	0	8	0					
2	7	0	1	0					
3	0	0	0	8					
÷	8	0	0	0					
	7	0	1	0					
W	1	2	4	1					
data matrix, X									
	(и	/ x  /	<b>۱</b> )						



#### Probabilistic motif model



# Probabilistic models of regulatory motifs (sequence families, consensus patterns)

- Represent the sequence pattern quantitatively using a statistical model
- Identify "matches to the motif" or instances using a likelihood ratio score (Naïve Bayes classification)

$$S(X) = \log \frac{p(X \mid \text{motif})}{p(X \mid bg)} = \log \frac{\prod_{i=1}^{w} \prod_{b \in ACGT} f_{ib}^{X_{ib}}}{\prod_{i=1}^{w} \prod_{b \in ACGT} g_{b}^{X_{ib}}} = \sum_{i=1}^{w} \sum_{b \in ACGT} X_{ib} \log \left(\frac{f_{ib}}{g_{b}}\right)$$
  
Sequence of length w  
Likelihood  
ratio score

# Probabilistic models of regulatory motifs (sequence families, consensus patterns)







#### ChIP-chip and ChIP-seq



Can be used to measure binding of transcription factors to the entire genome

#### Allelic imbalance in ChIP-seq

HNF6 (from Lannoy et al. JBC 1998) r.hnf-3b -126/-139 AAAAATCAATATCG -105/-92CTAAGTCAATAATC m.ttr r.pfk-2 prom. -211/-198 TGAAATCAATTTCA r.pfk-2 GRUc 1st intron AAAAATCCATAACT AGAAGTCAATGATC m.hnf-4 -389/-376 r.pepck -258/-245 TTTAGTCAATCAAA r.a2-uq TTTTATCAATAATA -183/-196 -41/-54r.cyp2c13 CAGAATCAATATTT r.cyp2c12 AAAAATCAATATTT -36/-49TCTAATCAATAAAT m.mup -132/-119 -208/-221 ATAAATCAATAGAC r.tog CAAAGTCAATAAAG r.a-fp -6103/6090

Look at peaks in liver of a Heterozygous Mouse

data from Mike Wilson c. 2011



#### Between Slc16a12 (monocarboxylic acid transporter) and Pank1 (liver CoA biosynthesis) chr19:34758655-34758994

HNF6 peak score=162.35

rs31029568 ref=A var=G 1/1:99:52:0:255,157,0:1

Alig	nment block 1 of	2 in window, 34758754 - 34758798, 42 bps
ВD	Mouse	gctca-tttgatgaggacagaa-a <mark>aagtcaatag</mark> aaccaagtga
ΒD	Rat	ggtcg-tttgatgagaacagaa-a <mark>aagtc<b>gg</b>tag</mark> aacaaagtga
ВD	Kangaroo rat	aaccattttgatgagaagagag-a <mark>aaatcaatag</mark> gacaaagtga
ΒD	Naked mole-rat	catcattttgatgagatcagaa-a <mark>atcaatag</mark> gacaaagtga
ВD	Guinea pig	agtcatctttgtgagggcagaa-a <mark>atcaatag</mark> gacaaagcaa
ВD	Squirrel	agccatttcaatgagaactggg-t <mark>aaatcaatag</mark> gacaaagaga
ВD	Rabbit	agtcattttgatgagaacagag-a <mark>aaa<b>c</b>ca<b></b>ta</mark> gataaagtga
ВD	Pika	agtcattttgatgagaatagag-a <mark>aac<b>c</b>caataa</mark> gataaagtga
ВD	Human	attgatgagaatggag-a <mark>aaatcaatag</mark> gacaaagtga
ВD	Chimp	agtcattttgatgagaatggag-a <mark>aaatcaatag</mark> gacaaagtga
ВD	Gorilla	agtcattttgatgagaatggag-a <mark>aaatcaatag</mark> gacaaagtga
ВD	Orangutan	agtcattttgatgagaatggag-a <mark>aaatcaatag</mark> gacaaagtga
ВD	Gibbon	agtcattttgatgagaatggag-a <mark>aaatcaatag</mark> gacaaggtga
ВD	Rhesus	agtcattttgatgagaatggag-a <mark>aaatcaatag</mark> gacaaagtga
ВD	Baboon	agtcattttgatgagaatggag-a <mark>aaatcaatag</mark> gacaaagtga
ВD	Marmoset	agtcattttgatgagaatggag-a <mark>aaatcaatcg</mark> gacaaagtga
ВD	Squirrel monkey	agtcattttgatgagaatggag-a <mark>aaatcaattg</mark> gacaaagtga
ВD	Mouse lemur	catcattctgataagaatggaa-a <mark>aaatcaatag</mark> gacaaagtga
ВD	Bushbaby	tgtcattttgataaagatggag-a <mark>aaatcaatag</mark> gccaaagtga
ВD	Pig	agtcgttttgacgagaacaggg-a <mark>aaatcaatag</mark> gacaaagtga
ВD	Alpaca	agtcattttgatgagaacagag-a <mark>aaatcaatag</mark> gacaaaggga
ВD	Dolphin	agtcattttgatgagaacaggg-a <mark>aaatcaatag</mark> gacaaagtga
ΒD	Sheep	agtetttttgatgagaacaggg-a <mark>aaateaatag</mark> gacaaagtga
ВD	Cow	agtetttttgatgagaacaggg-a <mark>aaateaatag</mark> gacaaagtga
ВD	Cat	agtcattttgaggagaacagagaa <mark>aaatcaatag</mark> gacaaagtgc
ΒD	Dog	agtcattttgatgagggcagag-a <mark>aaatcaatag</mark> gatagagtgc
ВD	Panda	agtccttttgatgagagcagag-a <mark>aaatcaatag</mark> gacaaagtgc
ΒD	Horse	agtcattttgatgaaaacaga <mark>aaatcaatag</mark> gacaaagtga
ΒD	Microbat	atacgagaacagag-g <mark>aaatcaatag</mark> aa-aaagtga
ВD	Megabat	agtcattttgacgagaacagag-a <mark>aaatcaatag</mark> gacaaaatga
ΒD	Shrew	agtccttttgaggagaacagagag <mark>aaatc<b>gg</b>tgc</mark> gacaaaatga
ΒD	Elephant	-gcaattttgaagagaacagag-g <mark>aaatcaatag</mark> gacaaagtga
ВD	Tenrec	-gccgttttgaagagaacagag-g <mark>aaatcaatag</mark> gacagtgtga
ВD	Manatee	-gccattttgaaaagaacagag-g <mark>aaatcaatag</mark> gacaaaatga
ВD	Armadillo	-gtcattttgatgagaacagag-a <mark>aaatcaatag</mark> gacaaagtga
ВD	Sloth	-ttcaatttgatgagaacagag-a <mark>aaatcaatag</mark> gacaaagtga

HNF6: AATCAATAA

쿱1-AAGTCAATAG S=7.56 ref: var: AAGTCAGTAG S=5.40 ↑ rs31029568 ∆S=2.16 In do40: 70 reads with 'A' Q>20 10 reads with 'G' **Binomial tests:**  $Z_{f=0.5} = 6.7$ Heterozygote  $Z_{f=0.99} = -10.3$ 

Homozygote, 1% error

Strong evidence for imbalance in the direction predicted by binding affinity chr1:39813171-39813514

Not near any genes of known function

Т

HNF6 peak score=301.85

rs47325601 ref=C var=T 1/1:99:32:0:255,96,0:1

Alig	nment block 1 of	2 in window, 39813380 - 39813 <b>7</b> 12, 33 bps
ΒD	Mouse	cgtgagcaaata <mark>taatcaacca</mark> aagca-catgag
ВD	Rat	tgtgagcaaata <mark>caatcaatca</mark> gagca-caagagag
ΒD	Naked mole-rat	cataagcaaaca <mark>cgatgaatcaatca</mark> atgca-taggtg
ΒD	Guinea pig	cataagcaaatacggtgaacc <mark>aatca</mark> atgtg-tagcca
ВD	Squirrel	cagaagcaaaggcaatg <mark>a<b>c</b>tcaatgc</mark> a-cagcag
ΒD	Rabbit	cacaagcaaatacaatg <mark>aatcaatgc</mark> a-tgagtg
ΒD	Human	cacaagcaaatataatg <mark>aatcaatgc</mark> a-tagccg
ΒD	Chimp	cacaagcaaatataatg <mark>aatcaatgc</mark> a-tagccg
ВD	Gorilla	cacaagcaaatataatg <mark>aatcaatgc</mark> a-tagccg
ВD	Gibbon	cacaagcaaatataatg <mark>aatcaatgc</mark> a-tagctg
ВD	Rhesus	cataagcaaatataatg <mark>aatcaatgc</mark> a-tagcgg
ВD	Baboon	cataagcaaatataatg <mark>aatcaatgc</mark> a-tagcgg
ΒD	Marmoset	caagagcaaatacaatg <mark>aatcaatgc</mark> a-tagccg
ВD	Squirrel monkey	cacgagcaaatacaatg <mark>aatcaatgc</mark> a-tagccg
ВD	Tarsier	cacaagtaaatacaatt <mark>aatcaatgc</mark> a-tggcag
ВD	Mouse lemur	cacaagcaaatattatgcaatg <mark>aatcaatgc</mark> a-tagcaa
ВD	Bushbaby	cacaagcaaataccatg <mark>agtcaatgc</mark> a-cagcag
ΒD	Tree shrew	cat-agcaaatataagg <mark>aatcaatgc</mark> a-tcgcag
ВD	Pig	cac-agcaaacacaatg <mark>aatcaatgc</mark> a-tagaag
ВD	Alpaca	cacaagcagacgtagct <mark>g</mark> agca <b>gga</b> cgctatgag
ВD	Dolphin	cacaagcaaatacaata <mark>aatcaa<b>a</b>gc</mark> a-tagaag
ВD	Cat	cataagcaaatacaatg <mark>aatcaa<b>c</b>ac</mark> a-tagaag
ВD	Dog	cacaagcaaacacaatg <mark>aatcaatgc</mark> a-ttgccc
ВD	Panda	cacaagcaaatacaatg <mark>aatcaa<b>a</b>gc</mark> a-tagaag
ΒD	Horse	cacaagcaactgcaatg <mark></mark> a-ttgaag
ВD	Microbat	cacgcaaatacaata <mark>aatcaatgc</mark> c-tagaag
ΒD	Megabat	caaaagcaaacacaatg <mark>aatcaatgc</mark> c-tagaag
ВD	Elephant	cacaagcaaatacactg <mark>aatcaa<b>c</b>gc</mark> a-tagaag
ВD	Manatee	cgcaagcaaatacaatg <mark>aatca<b>g</b>tgc</mark> a-tagaag
ВD	Armadillo	cacaagcaaatagaatg <mark>aatcaa<b>g</b>gc</mark> a-cagaag
		HNF6:

TAATCAACCA S=4.54 ref: var: TAATCAATCA S=6.71 ↑ rs47325601  $\Delta S = -2.16$ In do40:

21 reads with 'C' 59 reads with 'T'  $^{Q>20}$ Binomial tests:  $Z_{f=0.5} = 4.2$  Heterozygote

 $Z_{f=0.99} = -65.4$  Homozygote, 1% error

Strong evidence for imbalance in the direction predicted by binding affinity



Unpublished data

# De novo motif-finding with the probabilistic model

• Search for statistically enriched sequence families in non-coding DNA



Motif finding



• DNA sequences contain some w-mers from the motif model, and some from the background model

#### Mixture model





#### A two-component Mixture Model

We assume that...



Each position in the sequence is drawn from a multinomial representing the background or a position in the motif

#### A two-component Mixture Model



Each position in the sequence is drawn from a multinomial representing the background or a position in the motif

#### A two-component Mixture Model

Unobserved indicator variables determine whether each position is drawn from motif or background component



E.g., MEME

Bailey TL, Elkan C *ISMB*, pp. 28-36, AAAI Press, Menlo Park, California, (1994.)

#### How to estimate parameters?

Even though there are "hidden" variables, it's still possible to write the likelihood and find Maximum-Likelihood estimates for



Using various optimization schemes :

Stormo GD, Hartzell GW 3rd. Proc Natl Acad Sci U SA. (1989) Feb;86(4):1183-7.

Lawrence CE, Reilly AA. Proteins. 1990;7(1):41-51.

Lawrence CE et al. Science. (1993) Oct 8;262(5131):208-14.



#### EM notes

- Guaranteed to increase the likelihood at each iteration
- Can get stuck in local maxima
- Non-linear gradient ascent can be very fast
- Usually not possible to derive analytic updates for all parameters
- General formulation based on graph theory allows application to complex models



### Unsupervised classification

- Classification when you don't have a training set is called "unsupervised".
- Finding motifs "*de novo*" is a form of unsupervised classification (or clustering)
- Fundamentally hard problem because motifs appear by chance in long enough sequences

See e.g., Zia & Moses 2012



# Motif models so far assume all fixed width, w

- Many biological patterns have variable lengths
- Especially important for protein families and domains (but also for some TFs, e.g., p53)
- Like the residues, insertions and deletions have position-specific preferences
  EVI1\_HUMAN/131-154 TRA1\_CAPEL/337-362 H9Z/04\_DROME/370-392



EVII_HUMAN/131-154
TRA1_CAEEL/337-362
H9ZJM4_DROME/370-392
ZSC22_HUMAN/352-374
ZN239_MOUSE/6-28
A0A024RC04_HUMAN/488-510
G2HH24_PANTR/517-539
ZNF17_HUMAN/442-464
ZFP59_MOUSE/493-515
ZFP60_MOUSE/428-450
A0A023ZFK3_YEASX/50-73
ZFP60_MOUSE/344-366
ZFP59_MOUSE/326-348
ZFP60_MOUSE/484-506
ZFP59_MOUSE/270-292
TRA1 CAEEL/306-331

-	0	F		3.1	c	7	145	7 2	TD	D	CN	т	OF	u	т	рс		L.
ľ	~~	-	• •	• 14	-	н.	. r. v	-	11	5	21	-	×г	u	÷	n	· • ¥	-
2	C	Q	Ι.	PQ	С	т.	. KS	ŝΥ	ΤD	P	SS	L	Rŀ	ίH	Ι	K7	. V	Η
k	c	N		.1	С	G.	.KA	F	SR	P	WI	L	Q	Η	I	RI	·	Н
F	(C	G		.E	С	G.	. KI	F	SR	s	ΤН	L	тс	)H	o	RV	<i>.</i>	Н
ŀ	c	D		.K	С	G.	. K	F	TR	S	SS	L	Ľ	- 7H	Ĥ	sv	<i>.</i>	Н
E	c	D		.E	С	G.	. KH	E	SH	A	GA	L	FI	н	ĸ	MV	<i>.</i>	Н
ŀ	c	N		.0	С	G.	.11	F	so	N	SE	F	IV	ТH	0	IA		Н
Æ	c	N		. K	c	G.	. KE	F	RY	С	FΤ	L	NF	RH	õ	RV	<i>.</i>	Н
F	c	ĸ		.v	С	G.	. KS	Ē	KR	E	SN	L	IC	DH	Ĝ	AV	<i>.</i>	Н
c	c	K		. D	С	w.	.EF	F	RR	R	SN	F	IĒ	ΪH	o	SI		Н
Ĉ	c	N		Ι.	с	L.	. KE	Ē	SR	Ι	DN	L	RC	DН	õ	SS	v	Н
Ē	c	ĸ		.0	c	G.	. КІ	F	SN	G	SY	L	LF	RH	Ŷ	DT		H
F	c	N		.v	c	G.	. SZ	F	RL	0	LY	L	SF	ιH	0	КП	•	Н
F	c	ĸ		. E	c	G.	. KZ	F	HF	ŝ	so	L	NN	тн	Ŕ	TS		Н
c		ĸ		. D	c	G.	. K	F	IV	T.	AH	L	TF	RH	0	SS		H
k	c	E	F.	AD	c	Ε.	.KZ	F	SN	A	SD	R	AF	CH	õ	NF	.т	Н

### Profile Hidden-Markov Models

• Probabilistic model with three types of states



• The states are 'hidden', but they explain the sequences that we observe





### Profile Hidden-Markov Models

• How to get the parameters?



For HMMer (Pfam) this is usually done starting with a (manual) alignment and priors



– GLAM2 offers unsupervised estimation of HMMs, also relies heavily on priors



#### What are priors?

## Why are priors so important for profile HMMs ?

#### CABIOS

Vol. 12 no. 4 1996 Pages 327-345

#### Dirichlet mixtures: a method for improved detection of weak but significant protein sequence homology

Kimmen Sjölander<sup>3</sup>, Kevin Karplus, Michael Brown, Richard Hughey, Anders Krogh<sup>1</sup>, I.Saira Mian<sup>2</sup> and David Haussler

With accurate prior information about which kinds of amino acid distributions are reasonable in columns of alignments, it is possible with only a few sequences to identify which prototypical distribution may have generated the amino acids observed in a particular column. Using this informed guess, we adjust the expected amino acid probabilities to include the possibility of amino acids that may not have been seen but are consistent with observed amino acid distributions. The statistical models produced are more effective at generalizing to previously unseen data, and are often superior at database search and discrimination experiments (Brown et al., 1993; Tatusov et al., 1994; Bailey and Elkan, 1995; Karplus, 1995a; Hughey and Krogh, 1996; Henikoff and Henikoff, 1996; Wang et al., 1996).

 For DNA/RNA, it's not as important what priors you use, but you still need something that gives psuedocounts

$$S(X) = \log \frac{p(X \mid \text{motif})}{p(X \mid bg)} = \log \frac{\prod_{i=1}^{w} \prod_{b \in ACGT} f_{ib}^{X_{ib}}}{\prod_{i=1}^{w} \prod_{b \in ACGT} g_{b}^{X_{ib}}} = \sum_{i=1}^{w} \sum_{b \in ACGT} X_{ib} \log \left(\frac{f_{ib}}{g_{b}}\right)$$

 For proteins, without priors, models of highly variable protein motifs/domains with few known examples will always be overfit



**Dirichlet distribution** 

(3 dimensional distribution is illustrated, for DNA we need 4 dimensions and for protein we need 20 dimensions)

 $A \ C \ G \ T$   $1 \ 0 \ 0 \ 8 \ 0$   $2 \ 7 \ 0 \ 1 \ 0$   $3 \ 0 \ 0 \ 0 \ 8$   $\vdots \ 8 \ 0 \ 0 \ 0$   $7 \ 0 \ 1 \ 0$   $w \ 1 \ 2 \ 4 \ 1$  X = (7,0,1,0)  $p(A) = 7 \ / \ 8$ 



7/8 is the 'maximum likelihood' estimate...



With priors, we have a "Bayesian" estimator: mean posterior estimate

E.g., If  $\vec{\alpha} = (3,2,2,3)$  Where did I get these parameters? p(A) = 10 / 18 = 5 / 9 $\equiv f_A$ 

- For proteins, a single column of parameters in not enough
  - Biochemical similarity
  - Surface vs. core
  - Secondary structure vs. disordered
  - And more
- Amino acid probabilities in proteins have a very complicated distribution, and these are not uniform across the protein sequence
  Use a mixture of

Use a mixture of Dirichlet distributions!

 Sjolander et al. 1996 trained a 9-component mixture of Dirichelet priors (Blocks9)

#### The mixture of Dirichlet prior

Component	Ratio (	atio (r) of amino acid frequency relative to background frequency										
	$\frac{1}{8 \leq r}$	$4 \le r \le 8$	$2 \leq r \leq 4$	$1 \le r \le 2$	$1/2 \le r < 1$	$1/4 \le r < 1/2$	$1/8 \le r < 1/4$	r < 1/8				
1			SAT	CGP	NVM	OHRIKFLDW	EY					
2	Y	FW	н		LM	NOICVSR	TPAKDGE					
3			QE	KNRSHDTA	MPYG	VLIWCF						
4		KR	ò	Н	NETMS	PWYALGVCI	DF					
5		LM	Ĩ	FV		WYCTQ	APHR	KSENDG				
6		IV		LM	CTA	F	YSPWN	EQKRDGH				
7		D	EN	QHS	KGPTA	RY	MVLFWIC					
8			М	IVLFTYCA	WSHQRNK	PEG	D					
9			PGW	CHRDE	NQKFYTLAM	SVI						

Table II. Preferred amino acids of Blocks9

The function used to compute the ratio of the frequency of amino acid *i* in component *j* relative to the background frequency predicted by the mixture as a whole is  $(\alpha_{l,i}/|\vec{\alpha}|)/(\sum_{k} q_k \alpha_{k,i}/|\vec{\alpha}_k|)$ .

An analysis of the amino acids favored by each component reveals the following:

Component 1 favors small neutral residues.

Component 2 favors the aromatics.

Component 3 gives high probability to most of the polar residues (except for C, Y and W).

Component 4 gives high probability to positively charged amino acids and residues with NH<sub>2</sub> groups.

Component 5 gives high probability to residues that are aliphatic or large and non-polar.

Component 6 prefers I and V (aliphatic residues commonly found in  $\beta$  sheets), and allows substitutions with L and M.

Component 7 gives high probability to negatively charged residues, allowing substitutions with certain of the hydrophilic polar residues.

Component 8 gives high probability to uncharged hydrophobics, with the exception of glycine.

Component 9 gives high probability to distributions peaked around individual amino acids (especially P, G, W and C).

### Profile Hidden-Markov Models

• How to get the parameters?



 For HMMer (Pfam) this is usually done starting with a (manual) alignment and priors



– GLAM2 offers unsupervised estimation of HMMs, also relies heavily on priors



- How to scan DNA/RNA/Protein sequences?
  - HMMer uses forward algorithms (after some speedup tricks)



- GLAM2 uses an alignment based approach

## Scoring DNA/proteins sequences with an HMM

- We want to calculate the probability of a sequence, X, given our profile HMM model, and compare this to a background model.
- In general, multiple paths through an HMM can give us the same sequence (unlike for a PSSM)
- We need to sum over all of the possible paths through the model, weighting each by their probability
- This is done using the "forward algorithm"

### Forward algorithm for HMMs

- Even though there are exponentially many paths, they are all made out of the same pieces.
- The trick is to add up the probability of all the paths that can get you to this point in the HMM, but don't keep track of all of the paths.
- Example of "dynamic programming"



### How does HMMer score sequences?



The null model is a one-state HMM configured to generate "random" sequences of the same mean length as the target sequence, with each residue drawn from a background frequency distribution (a standard i.i.d. model: residues are treated as independent and identically distributed).

 $Log = \frac{p(X | profile HMM)}{p(X | background HMM)}$ 

## **Classic Bioinformatics topics!**

- consensus sequences
- regular expressions
- PWMs / PSSMs / matrix models
  - Motif scanning / matrix matching
  - De novo motif finding
- profile HMMs



book chapter

#### Questions for Glam2 discussion

- What is the prior model for indels in Glam2?
- How does the Glam2 model differ from a profile HMM?
  Why does this make it easier to estimate parameters?
- What's the difference between an alignment-based score and the full Forward algorithm?
- What is the optimization strategy used by Glam2? Why not use EM?