

Sequence comparaison. II.

						motif 2						
			193	195	198		217	ATP	ATP	231	233	235
YstASP	293	KGKAY.LAQS	PQFNKQQLIV	ADFERVYEIG	PVFRAENSNT	HRHMTEFTGL	DMEMAFEEHY	HEVL	355		
ThtASP	190	PGLFYALPQS	PQLFKQMLMV	AGLDRYFQIA	RCFRDEDLRA	.DRQPDEFTQL	DLEMSF.VEV	EDVL	251		
EcoASP	184	KGKFYALPQS	PQLFKQLLMM	SGFDRIYQIV	KCFRDEDLRA	.DRQPEFTQII	DVETSF.MTA	PQVR	245		
EcoASN	191	QGKVFDFKDF	FGKESFLTVS	GQLNGEYAC	.ALSKIYTFG	PTFRAENSNT	SRHLAEFWML	EPEVAFAN.L	NDIA	262		
YstASN	198	.NTSPTASSY	FGKPTYLTVS	TQLHLEILAL	.SLSRCWTLS	PCFRAEKSDT	PRHLSEFWML	EVEMCFVNSV	NELT	269		
ThtLYS	205	.PFKTYHNAL	DHEFY.LRIS	LELYLKRLLV	GGYEKVFIEG	RNFRNEGIDH	.NHNPEFTML	EAYWAYAD.Y	QDMA	274		
EcoLYS	221	.PFITHHNAL	DLDMY.LRIA	PELYLKRLVV	GGFERVFEIN	RNFRNEGISV	.RHNPEFTMM	ELYMAYAD.Y	KDLI	290		
YstLYS	284	.PFITHHNDL	DMDMY.MRIA	PELFLKQLVV	GGLDRVYEIG	RQFRNEGIDM	.THNPEFTTC	EFYQAYAD.V	YDLM	353		



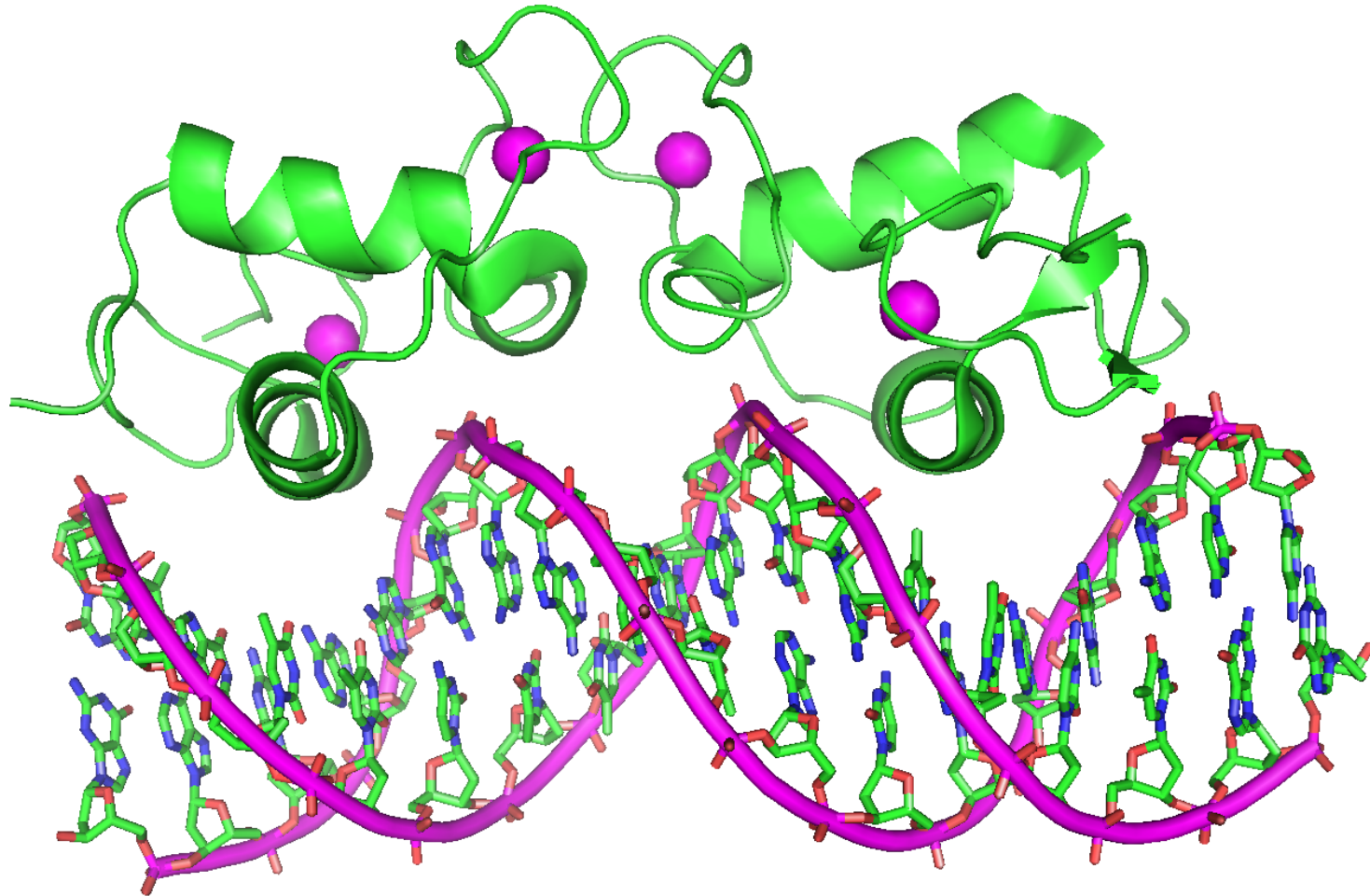
Multiple sequence alignments: method

Applications

Functional sites and their identification

See book chapters 2 and 4

Androgen receptor



Homologues of the androgen receptor identified using BLAST

#	ID Swissprot	Name	Description	Score	E	% Identity	Length aligned
1	P15207	ANDR_RAT	Androgen receptor.	162	1e-40	100	73
6	P19091	ANDR_MOUSE	Androgen receptor.	162	1e-40	100	73
14	Q63449	PRGR_RAT	Progesterone receptor	136	1e-32	80	72
17	P06401	PRGR_HUMAN	Progesterone receptor	136	1e-32	80	72
21	P08235	MCR_HUMAN	Mineralocorticoid receptor	136	1e-32	79	72
33	P04150	GCR_HUMAN	Glucocorticoid receptor	131	3e-31	77	72
41	Q9YH32	ESR2_ORENI	Estrogen receptor beta	99	3e-21	58	72
42	Q9YH33	ESR1_ORENI	Estrogen receptor alpha	98	4e-21	55	72
:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:
343	Q9N4Q7NH13_CAEEL		Nuclear hormone receptor nhr-13	54	8e-08	39	66
344	Q23294 NH11_CAEEL		Nuclear hormone receptor nhr-11	54	8e-08	42	66
345	O45460 NH54_CAEEL		Nuclear hormone receptor nhr-54	54	1e-07	37	67
346	Q09565 NH20_CAEEL		Nuclear hormone receptor nhr-20	51	7e-07	34	66
347	Q09587 NH22_CAEEL		Nuclear hormone receptor nhr-22	45	5e-05	32	66
349	P17672 E75B_DROME		Ecdysone-induced protein 75B	40	0.001	37	47

How can we combine these into a single alignment?

Dynamic programming algorithm for two sequences (seen above)

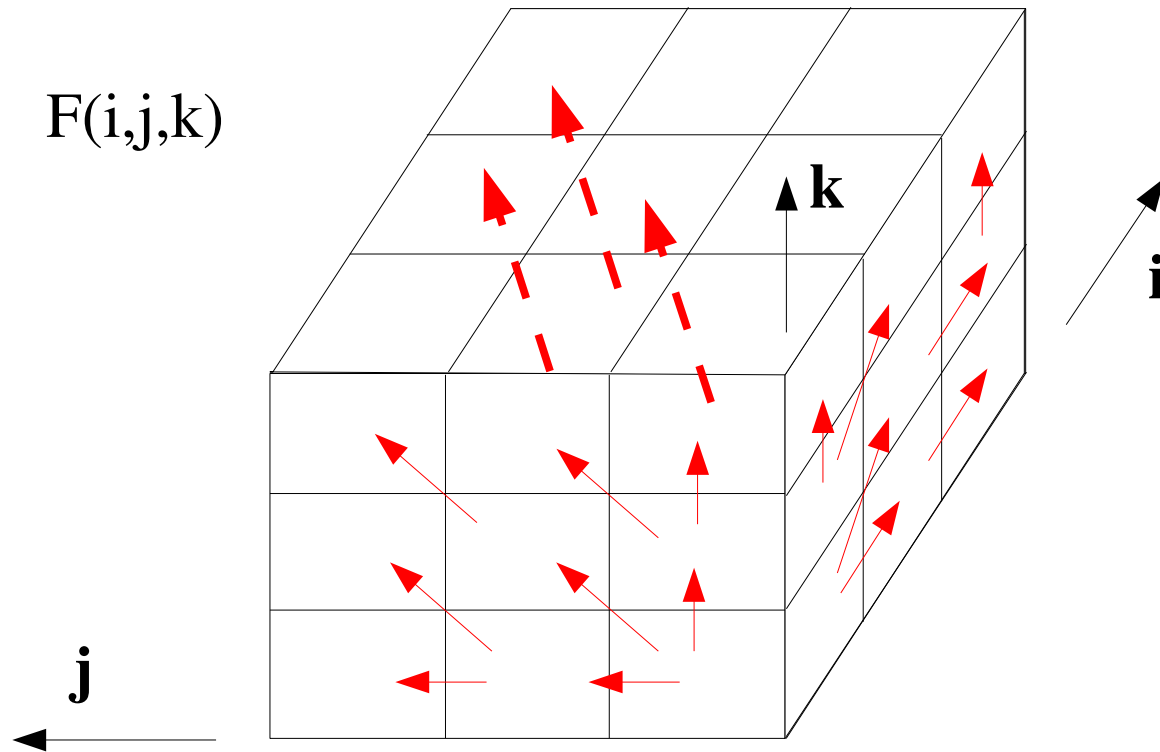
$F(i,j)$ is computed recursively from smaller alignments

$$F(i,j) = \max \begin{cases} F(i-1,j-1) + s(x_i, y_j) \\ F(i-1,j) - d \\ F(i,j-1) - d \end{cases}$$

		S	T	A	R
P	0 ← -8 -16				
I	-8 ↖ -1 ↘ -9				
T					

$s(P, S) = -1 = s(P, T)$

Dynamic programming algorithm for 3 sequences



Too expensive for more than ~10 sequences

$$50^{10} = 10^{17}$$

We need a more efficient, ‘heuristic’ method

#	ID Swissprot	Name	Description	Score	E	% Identity	Length aligned
1	P15207	ANDR_RAT	Androgen receptor.	162	1e-40	100	73
6	P19091	ANDR_MOUSE	Androgen receptor.	162	1e-40	100	73
14	Q63449	PRGR_RAT	Progesterone receptor	136	1e-32	80	72
17	P06401	PRGR_HUMAN	Progesterone receptor	136	1e-32	80	72
21	P08235	MCR_HUMAN	Mineralocorticoid receptor	136	1e-32	79	72
33	P04150	GCR_HUMAN	Glucocorticoid receptor	131	3e-31	77	72
41	Q9YH32	ESR2_ORENI	Estrogen receptor beta	99	3e-21	58	72
42	Q9YH33	ESR1_ORENI	Estrogen receptor alpha	98	4e-21	55	72
:	:	:	:	:	:	:	
:	:	:	:	:	:	:	
343	Q9N4Q7	NH13_CAEEL	Nuclear hormone receptor nhr-13	54	8e-08	39	66
344	Q23294	NH11_CAEEL	Nuclear hormone receptor nhr-11	54	8e-08	42	66
345	O45460	NH54_CAEEL	Nuclear hormone receptor nhr-54	54	1e-07	37	67
346	Q09565	NH20_CAEEL	Nuclear hormone receptor nhr-20	51	7e-07	34	66
347	Q09587	NH22_CAEEL	Nuclear hormone receptor nhr-22	45	5e-05	32	66
349	P17672	E75B_DROME	Ecdysone-induced protein 75B	40	0.001	37	47

How to merge these homologues into a single alignment?

#	ID Swissprot	Name	Description	Score	E	% Identity	Length aligned
1	P15207	ANDR_RAT	Androgen receptor.	162	1e-40	100	73
6	P19091	ANDR_MOUSE	Androgen receptor.	162	1e-40	100	73
14	Q63449	PRGR_RAT	Progesterone receptor	136	1e-32	80	72
17	P06401	PRGR_HUMAN	Progesterone receptor	136	1e-32	80	72
21	P08235	MCR_HUMAN	Mineralocorticoid receptor	136	1e-32	79	72
33	P04150	GCR_HUMAN	Glucocorticoid receptor	131	3e-31	77	72
41	Q9YH32	ESR2_ORENI	Estrogen receptor beta	99	3e-21	58	72
42	Q9YH33	ESR1_ORENI	Estrogen receptor alpha	98	4e-21	55	72
:	:	:	:	:	:	:	:
:	:	:	:	:	:	:	:
343	Q9N4Q7	NH13_CAEEL	Nuclear hormone receptor nhr-13	54	8e-08	39	66
344	Q23294	NH11_CAEEL	Nuclear hormone receptor nhr-11	54	8e-08	42	66
345	O45460	NH54_CAEEL	Nuclear hormone receptor nhr-54	54	1e-07	37	67
346	Q09565	NH20_CAEEL	Nuclear hormone receptor nhr-20	51	7e-07	34	66
347	Q09587	NH22_CAEEL	Nuclear hormone receptor nhr-22	45	5e-05	32	66
349	P17672	E75B_DROME	Ecdysone-induced protein 75B	40	0.001	37	47

One method is to infer the alignment from the Blast results: how?

A heuristic method for multiple sequence alignment in three stages

- 1) Align each pair of sequences
- 2) *Classify* the sequences by similarity: “guide” tree
- 3) Incorporate the sequences progressively into an alignment, in a sensible order (determined by the tree)

A pairwise alignment corresponds to simple evolutionary model, including

- An assumption of “minimal” divergence from a *common ancestor*
- An empirical estimation of mutation probabilities
- An assumption of equivalent and independent positions
- A well-defined reference, or “null” model

C E C H A T E S T G R A S
L E R - A T E S T P R I S

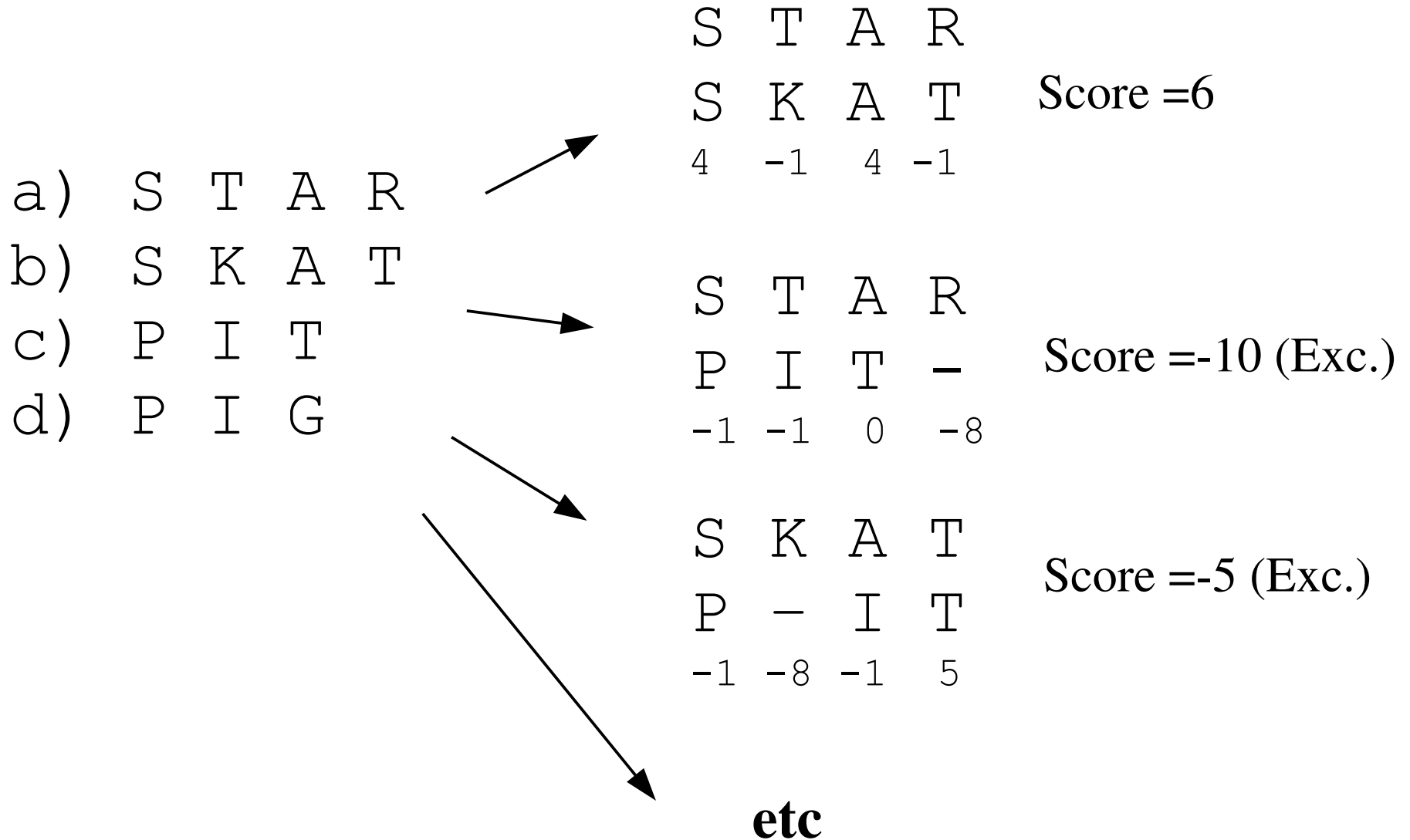
An empirical scoring matrix

BLOSUM62

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
C	9	-1	-1	-3	0	-3	-3	-3	-4	-3	-3	-3	-3	-1	-1	-1	-1	-2	-2	-2
S		4	1	-1	1	0	1	0	0	0	-1	-1	0	-1	-2	-2	-2	-2	-2	-3
T			5	-1	0	-2	0	-1	-1	-1	-2	-1	-1	-1	-1	-1	0	-2	-2	-2
P				7	-1	-2	-2	-1	-1	-1	-2	-2	-1	-2	-3	-3	-2	-4	-3	-4
A					4	0	-2	-2	-1	-1	-2	-1	-1	-1	-1	-1	0	-2	-2	-3
G						6	0	-1	-2	-2	-2	-2	-2	-3	-4	-4	-3	-3	-3	-2
N							6	1	0	0	1	0	0	-2	-3	-3	-3	-3	-2	-4
D								6	2	0	-1	-2	-1	-3	-3	-4	-3	-3	-3	-4
E									5	2	0	0	1	-2	-3	-3	-2	-3	-2	-3
Q										5	0	1	1	0	-3	-2	-2	-3	-1	-2
H											8	0	-1	-2	-3	-3	-3	-1	2	-2
R												5	2	-1	-3	-2	-3	-3	-2	-3
K													5	-1	-3	-2	-2	-3	-2	-3
M														5	1	2	1	0	-1	-1
I															4	2	3	0	-1	-3
L																4	1	0	-1	-2
V																	4	-1	-1	-3
F																		6	3	1
Y																			7	2
W																				11

also gap penalties

Stage 1: alignment of all sequence pairs

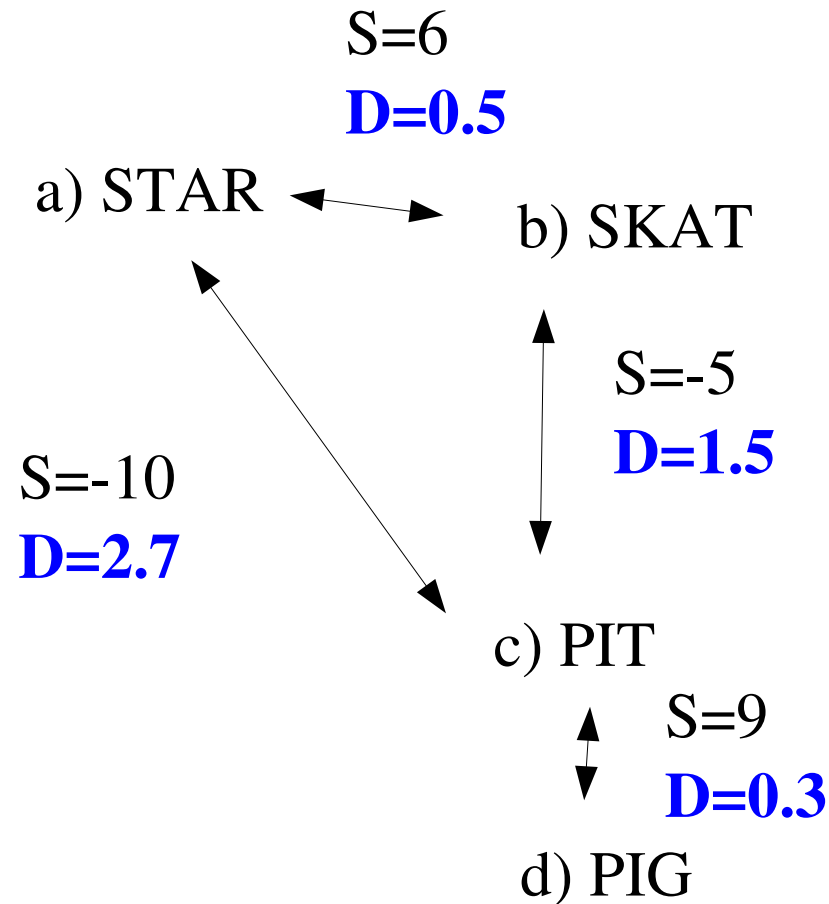


Stage 2: Classify the sequences by similarity: “guide” tree

Stage 3: Incorporate the sequences progressively into the alignment, in a sensible order

Stage 2: compute “distances” between sequences

Classification methods use distances, rather than similarity



What would be your choice of distance?

Stage 2: compute “distances” between sequences

$$D(a,b) = -\log \frac{S(a,b) - \bar{S}_{\text{rand}}}{S_{\text{max}} - \bar{S}_{\text{rand}}}$$

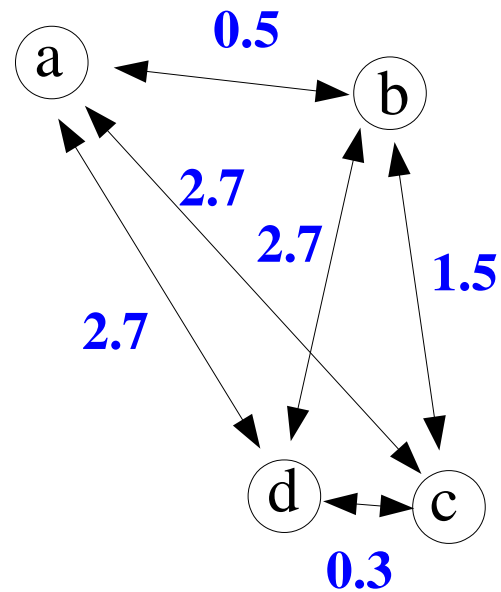
$$S_{\text{max}} = \max(S(a,a), S(b,b))$$

S_{rand} = score to align two random sequences of same length, composition

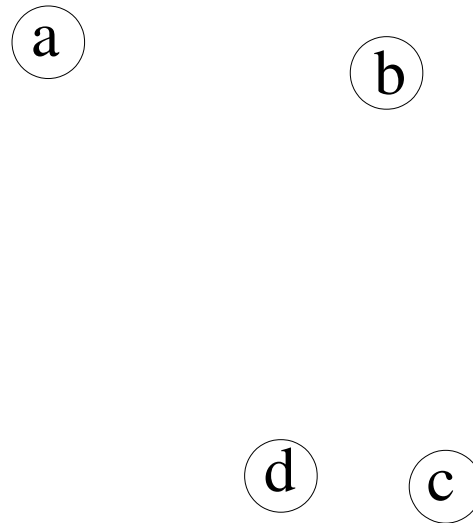
$$\begin{aligned} S &= 6 \\ S_{\text{max}} &= 18 && [\text{Blosum62}] \\ S_{\text{rand}} &= -12 \\ \mathbf{D} &= \mathbf{0.5} \end{aligned}$$

STAR \longleftrightarrow SKAT

Stage 2: compute “distances” between sequences



Stage 2: hierarchical classification

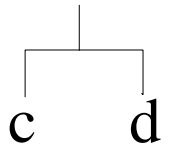
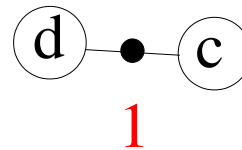


“Unweighted Pair Group Joining with Arithmetic Mean”
or UPGMA

Stage 2: hierarchical classification

a

b



“Unweighted Pair Group Joining with Arithmetic Mean”
or UPGMA

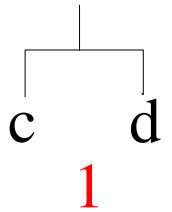
Stage 2: hierarchical classification

a

b

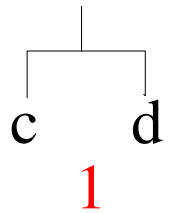
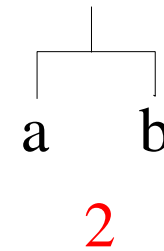
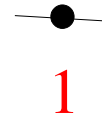
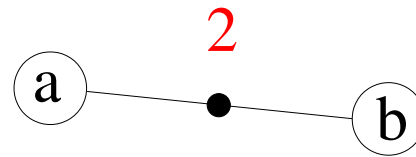


1



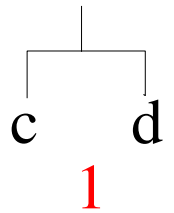
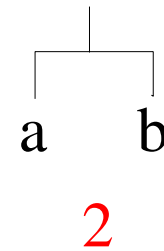
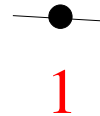
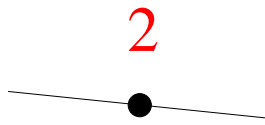
“Unweighted Pair Group Joining with Arithmetic Mean”
or UPGMA

Stage 2: hierarchical classification



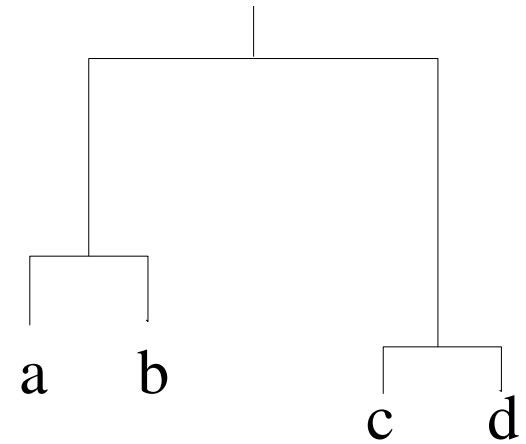
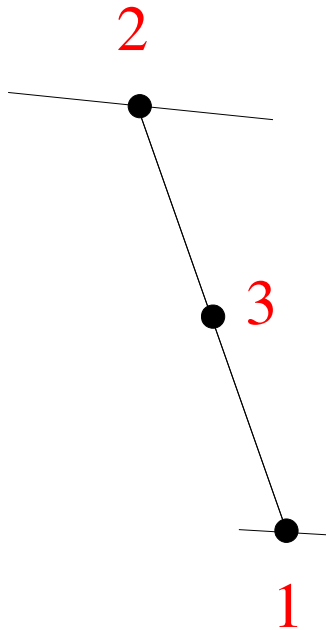
“Unweighted Pair Group Joining with Arithmetic Mean”
or UPGMA

Stage 2: hierarchical classification



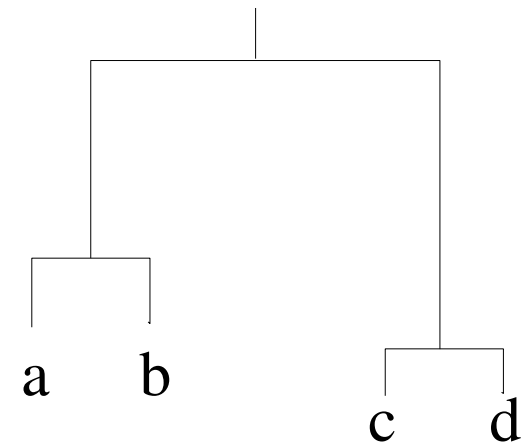
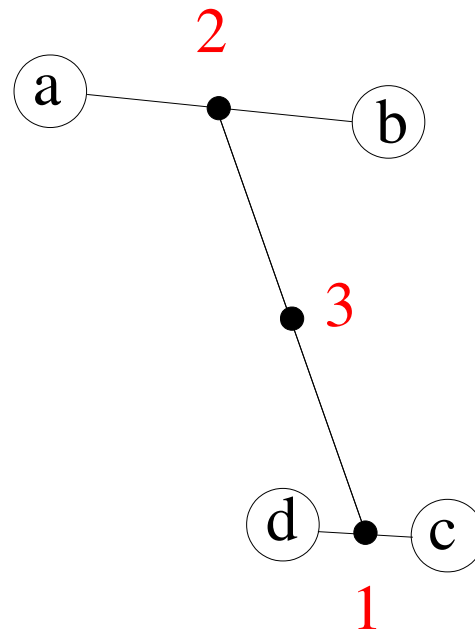
“Unweighted Pair Group Joining with Arithmetic Mean”
or UPGMA

Stage 2: hierarchical classification, or guide tree



“Unweighted Pair Group Joining with Arithmetic Mean”
or UPGMA

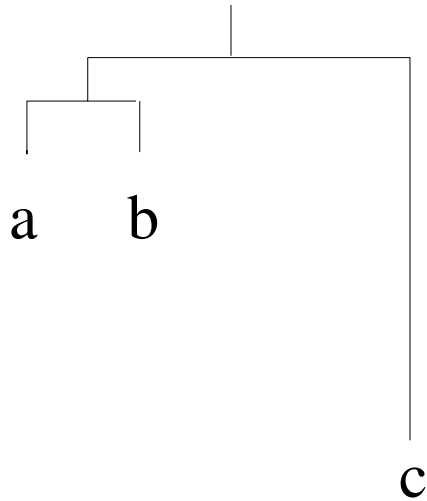
Stage 2: hierarchical classification, or guide tree



“Unweighted Pair Group Joining with Arithmetic Mean”
or UPGMA

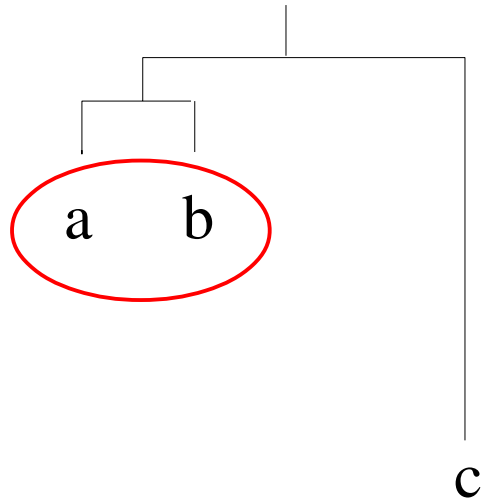
Stage 3: sequence alignment, moving up the tree

Example with 3 sequences



Stage 3: sequence alignment, moving up the tree

Example with 3 sequences

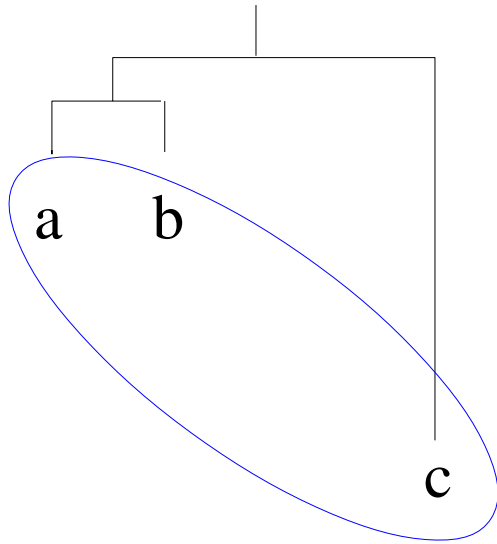


a with b

a) S T A R
b) S K A T

Stage 3: sequence alignment, moving up the tree

Example with 3 sequences



a with b

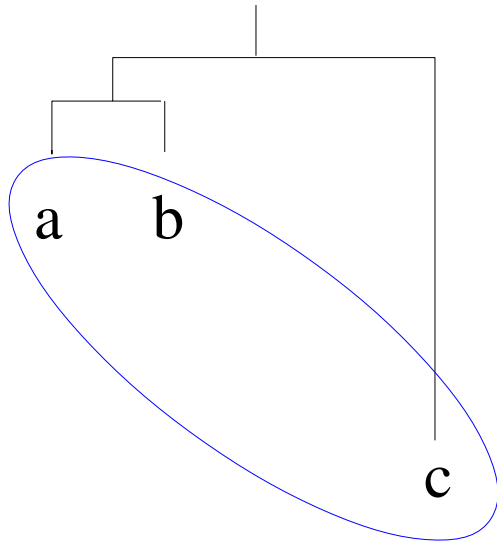
a) S T A R

b) S K A T

c with {a, b}

Stage 3: sequence alignment, moving up the tree

Example with 3 sequences



a with b

a) S T A R
b) S K A T

c with {a, b}

**How to align a sequence with
an alignment?**

“sequence-profile” alignment

We need to generalize slightly our pairwise alignment method....

Aligning a sequence with an alignment

An alignment of 3 sequences:

S	T	A	R
S	T	I	R
S	K	A	T

4th sequence y, to be aligned:

P I T

Aligning a sequence with an alignment

An alignment of 3 sequences:

S	T	A	R
S	T	I	R
S	K	A	T

4th sequence y , to be aligned:

P I T

View the alignment (above) as *a sequence of columns* X_i

Aligning a sequence with an alignment

An alignment of 3 sequences = a *sequence* of columns X_i

S	T	A	R
S	T	I	R
S	K	A	T

X_2

4th sequence y , to be aligned:

P	I	T
---	---	---

y_1

Aligning a sequence with an alignment: computing a “mean” score

An alignment of 3 sequences:

S	T	A	R			P
S	T	I	R		T	-1
S	K	A	T		T	-1
	X_2			\rightarrow	K	-1

4th sequence y, to be aligned:

P	I	T
---	---	---

y_1

Total	-3	=	$S(y_1, X_2)$
-------	----	---	---------------

Aligning a sequence with an alignment: computing a “mean” score

An alignment of 3 sequences:

S	T	A	R		P	
S	T	I	R		T	-1
S	K	A	T		T	-1
				→	K	-1
					total	-3

4th sequence y, to be aligned:

P	I	T
---	---	---

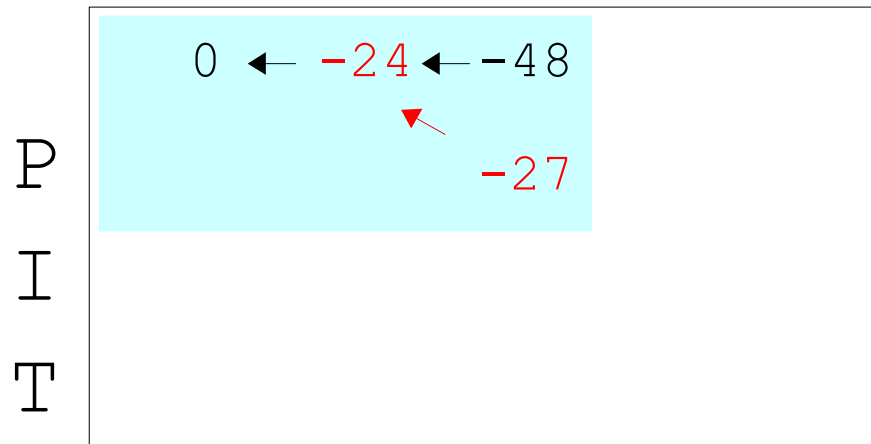
$s(a, \begin{bmatrix} b \\ b' \\ b'' \end{bmatrix}) = s(a,b) + s(a,b') + s(a,b'')$
--

Sum over pairs

Sequence-“profile” alignment : dynamic programming, as before

S T A R
S T I R
S K A T

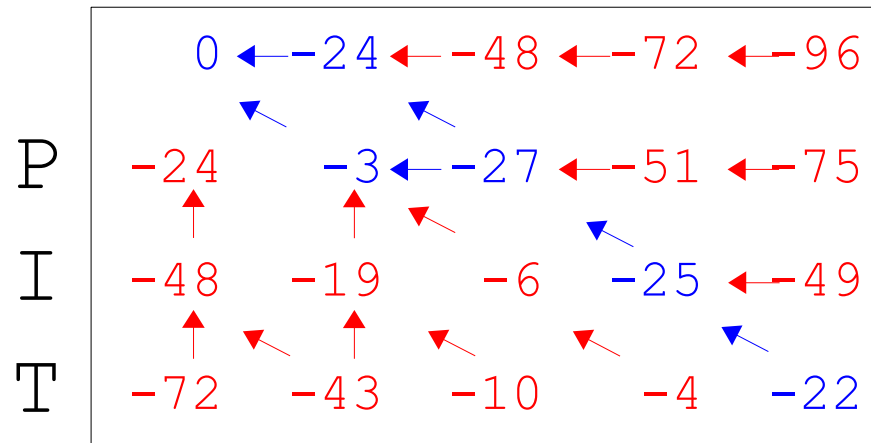
	P
T	-1
T	-1
K	-1
total	-3



gap penalties are tripled

Sequence-“profile” alignment : dynamic programming

S	T	A	R
S	T	I	R
S	K	A	T



S	T	A	R
S	T	I	R
S	K	I	T
P	-	I	T
↙	←	↖	↖

Aligning two alignments: “profile-profile” alignment

S	T	A	R
S	T	I	R
S	K	A	T

P	P	$s\left(\begin{matrix} a \\ a' \end{matrix}, \begin{matrix} b \\ b'' \end{matrix}\right) = s(a,b) + s(a,b') \\ + s(a,b'') + s(a',b) \\ + s(a',b') + s(a',b'')$
I	I	
G	T	

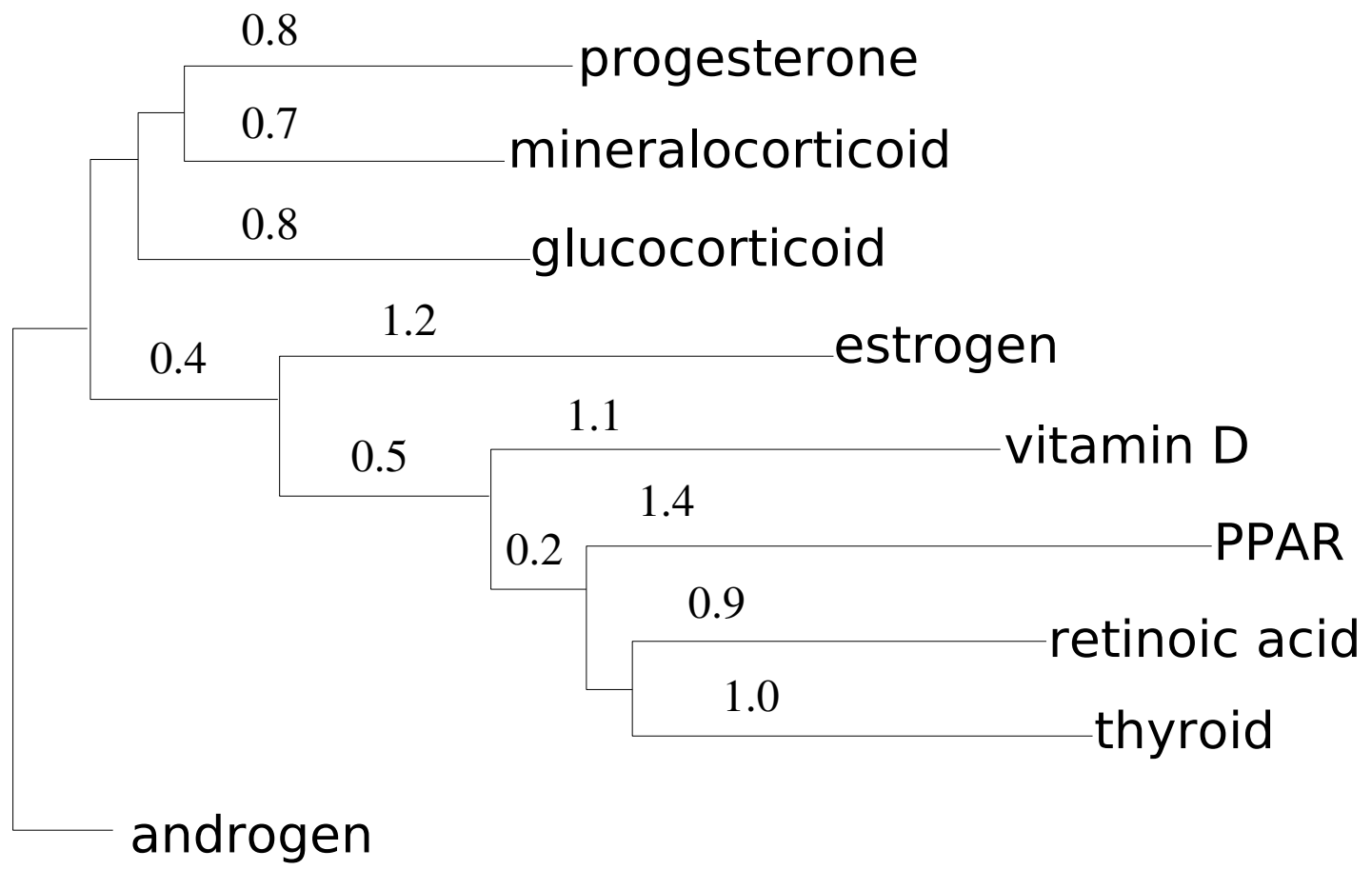


S	T	A	R
S	T	I	R
S	K	A	T
P	-	I	T
P	-	I	G

Sum over pairs

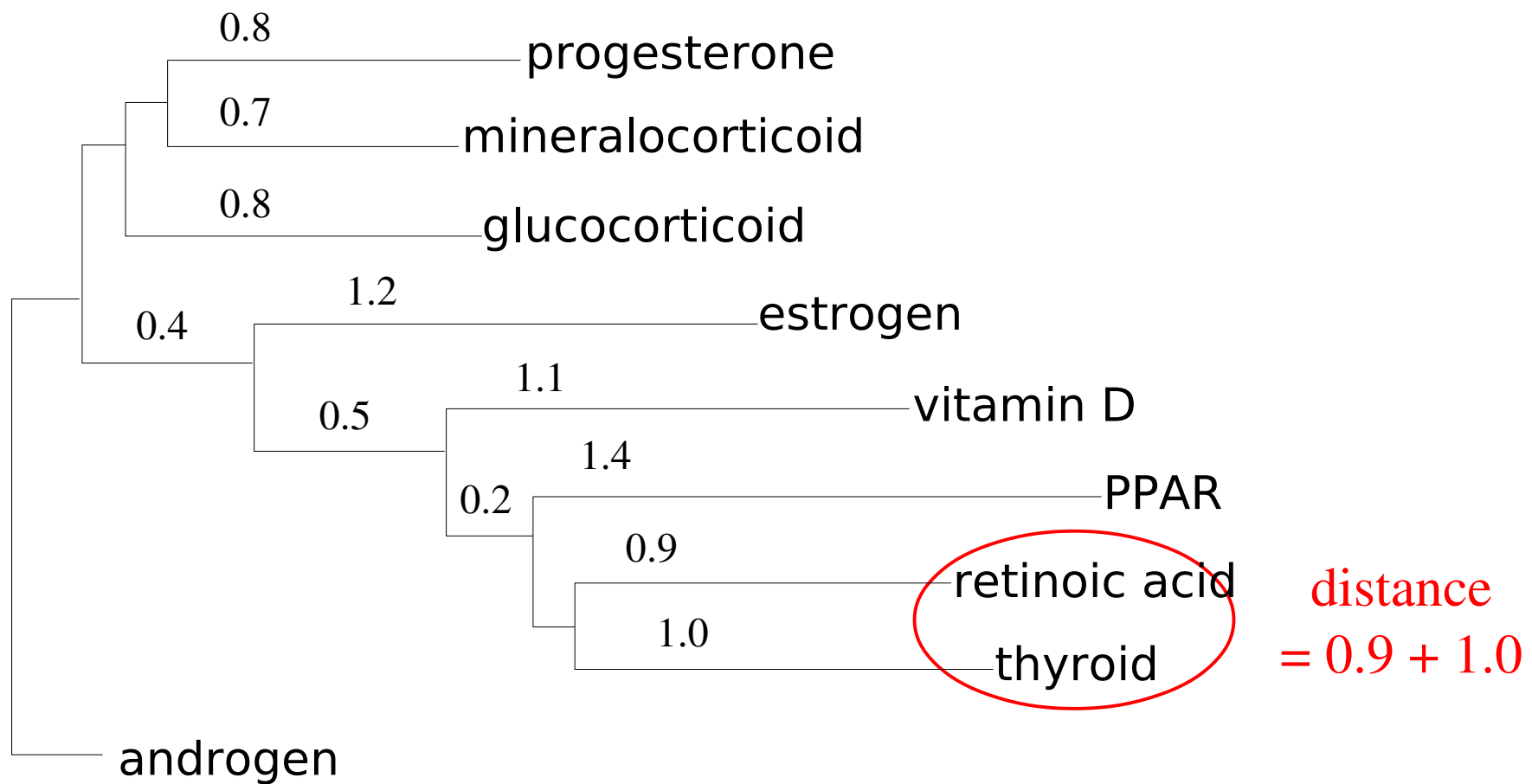
slight generalization of sequence alignment

Stage 3: align the sequences in the order of the guide tree



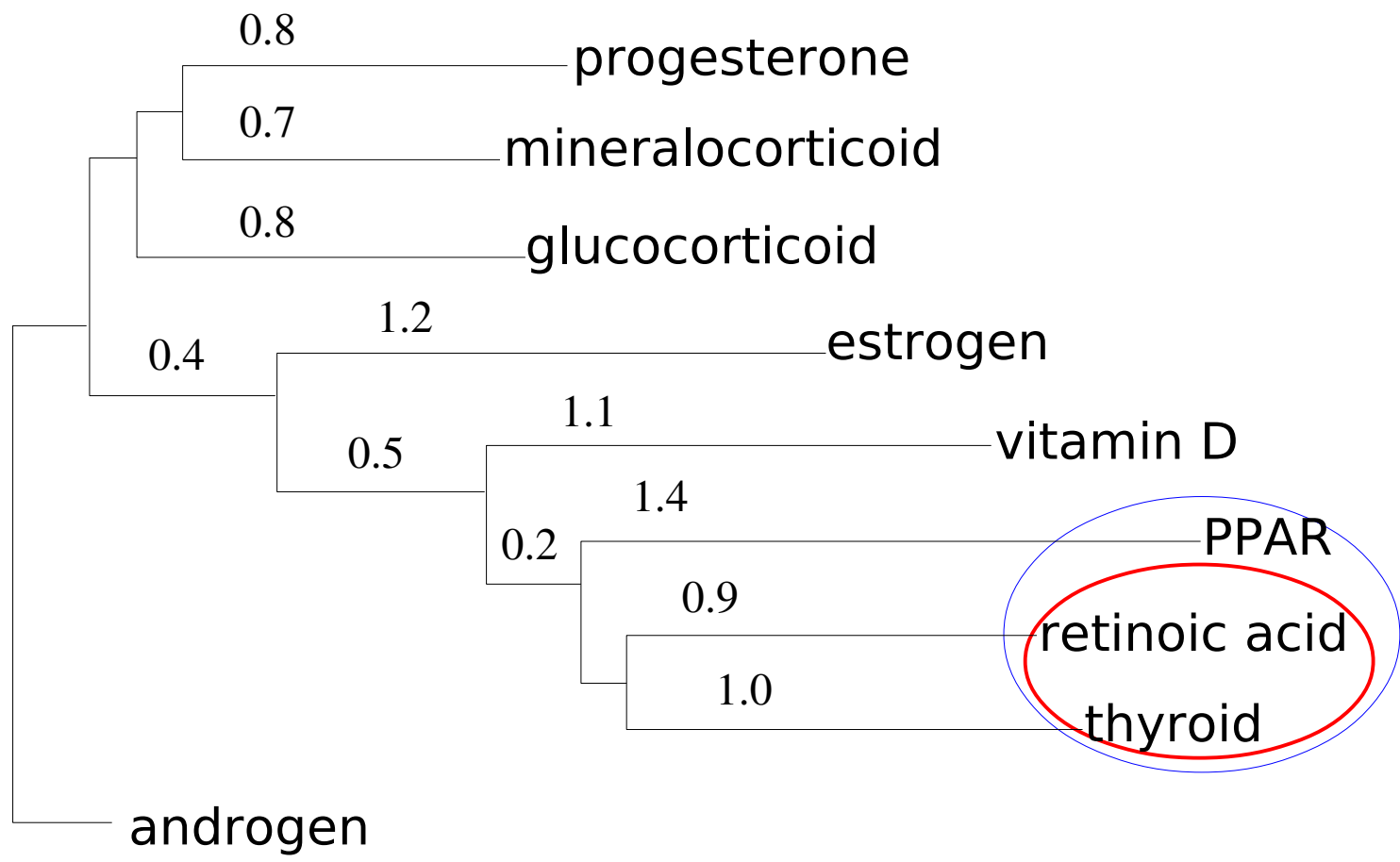
Branch lengths measure distances between sequences

Stage 3: align the sequences in the order of the guide tree

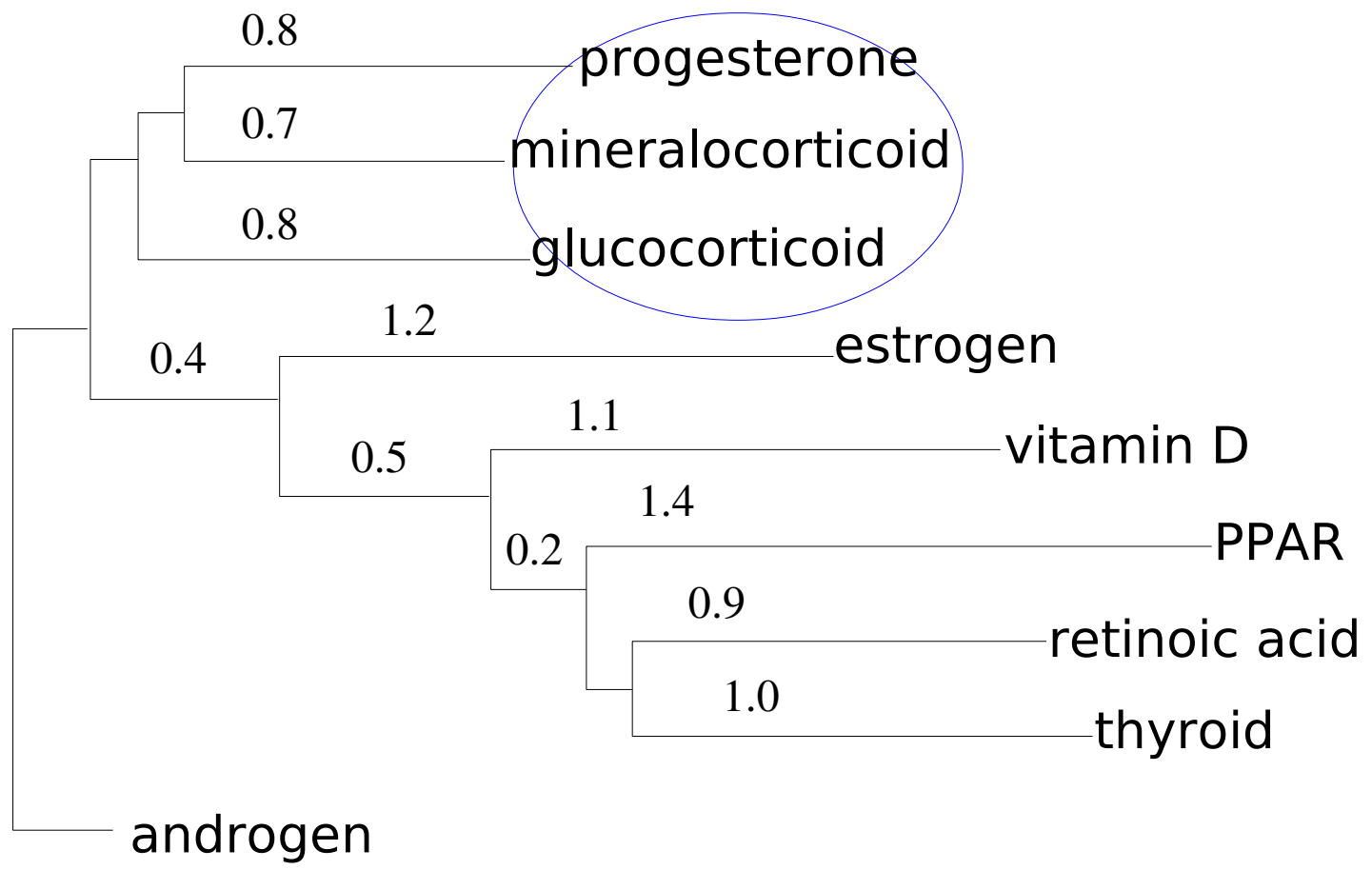


Branch lengths measure distances between sequences

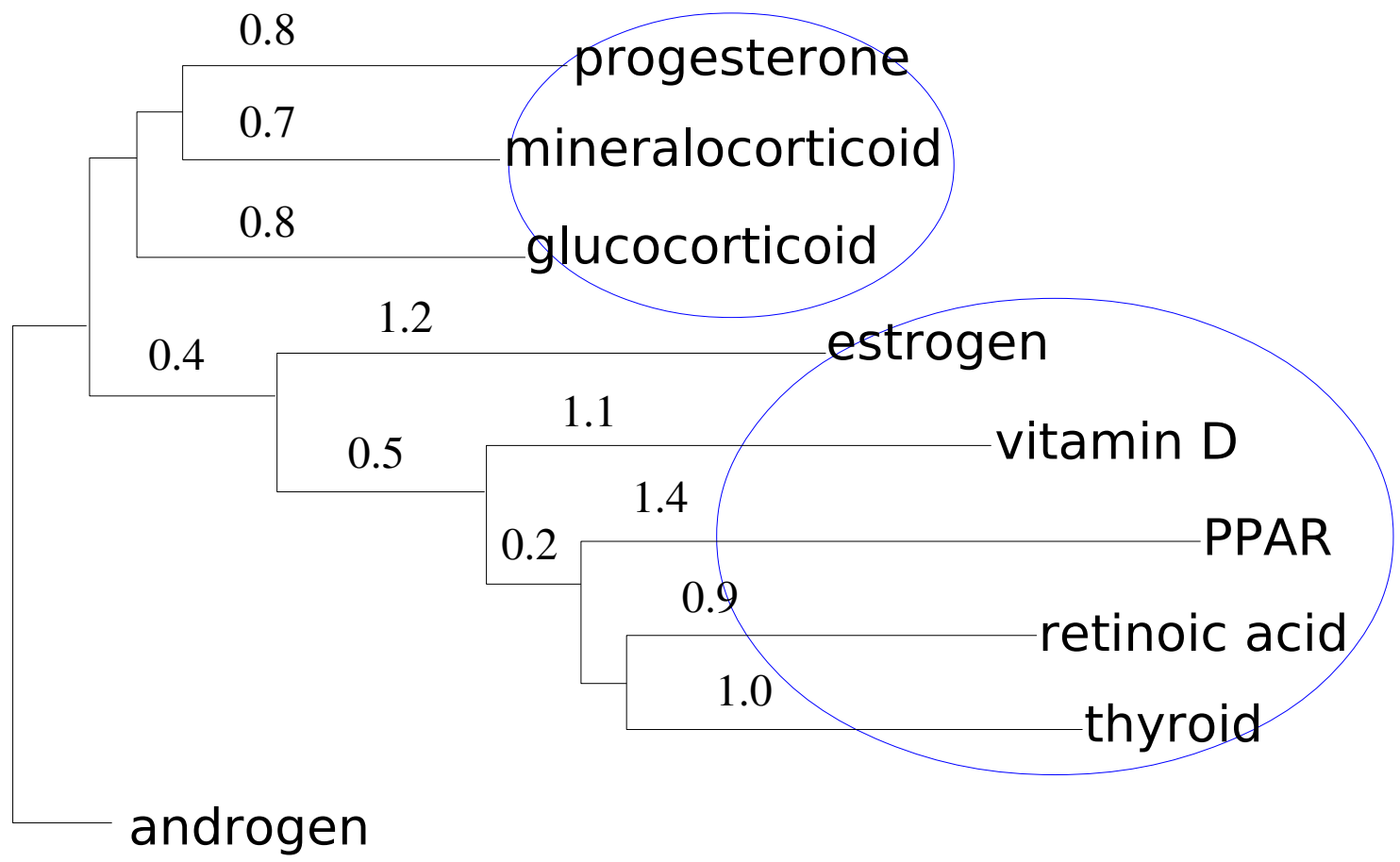
Stage 3: align the sequences in the order of the guide tree



Stage 3: align the sequences in the order of the guide tree

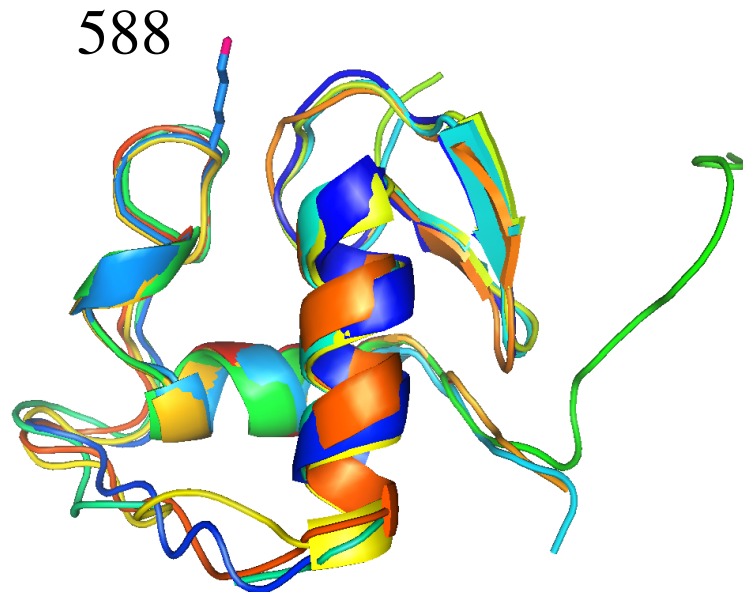


Stage 3: align the sequences in the order of the guide tree

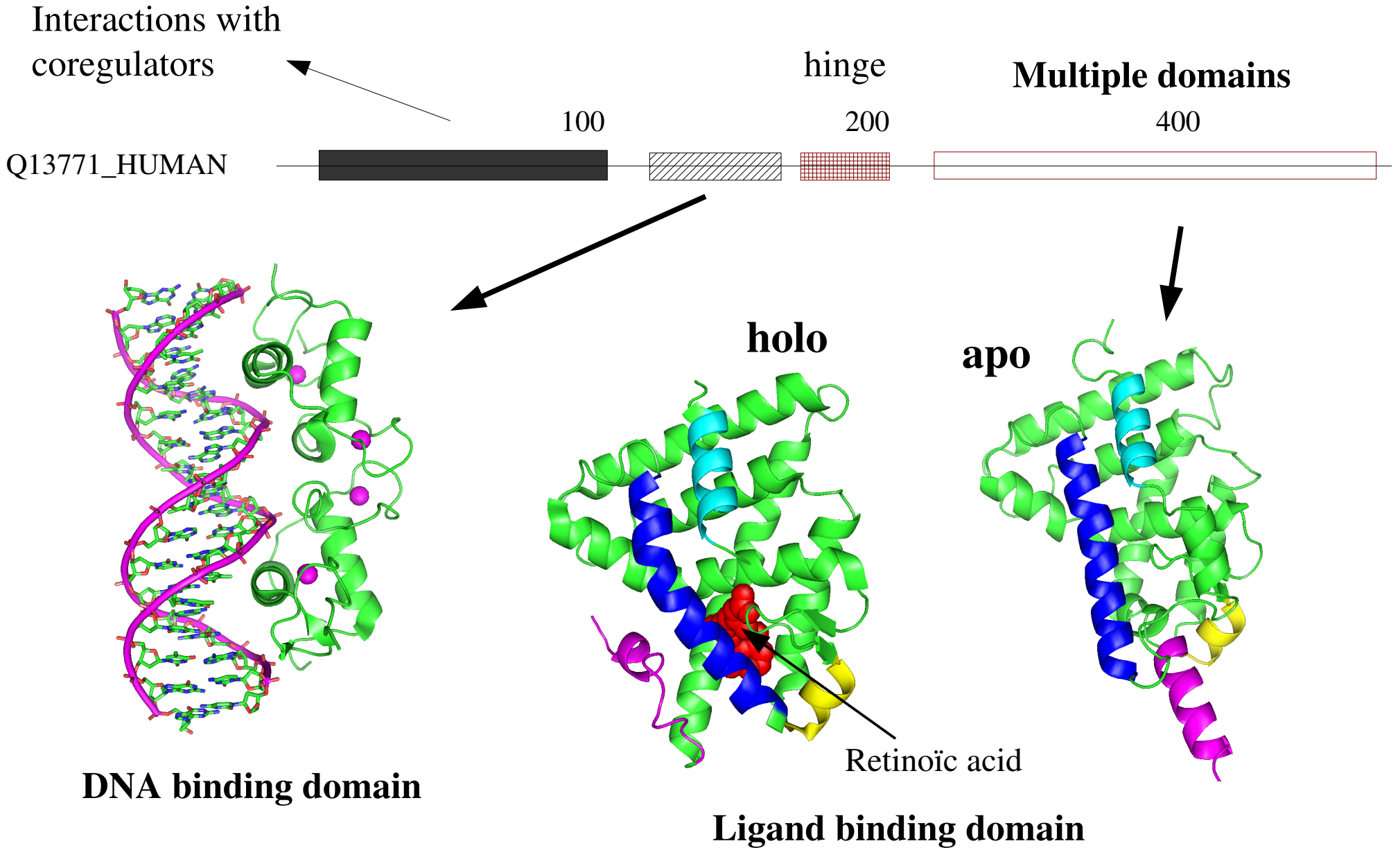


The final result

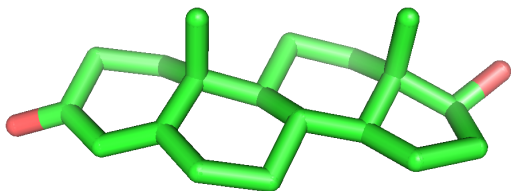
		584	588	
androgen	VFFKRAAEG--KQKYL	CASRNDCTIDK	FRRKNCPSCLR	LRKCY
progesterone	VFFKRAVEG--HHNYL	CAGRND CIVDK	IRRKNC PACRL	LRKCY
Mineralocorticoid	VFFKRAVEG--QHNYL	CAGRND CIIDK	IRRKNC PACRL	LQKCL
Glucocorticoid	VFFKRAVEG--QHNYL	CAGRND CIIDK	IRRKNC PACRY	RKCL
Estrogen	AFFKRSIQG--HNDYM	CPATNQCTIDK	NRRKSCQACRL	LRKCY
Retinoic acid	GFFRRSIQK--NMVYT	CHRDKNCIINK	VTRNRCQYCRL	LQKCF
Vitamin D3	GFFRRSMKR--KALFT	CPFNGDCRITK	DNRRHCQACRL	LKRCV
Thyroid	GFFRRTIQKNLHPTYS	CKYDSCCVI	DKITRNQCQLC	RFFKKCL
	**:* : . : :	* : *	* . *	** : : *



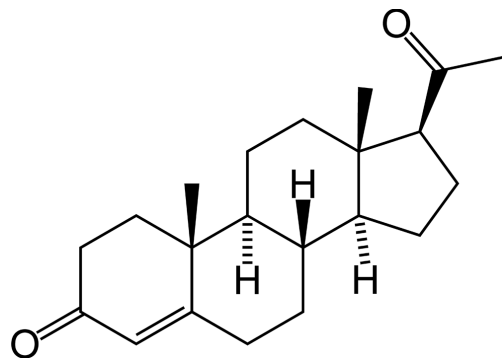
Nuclear receptors



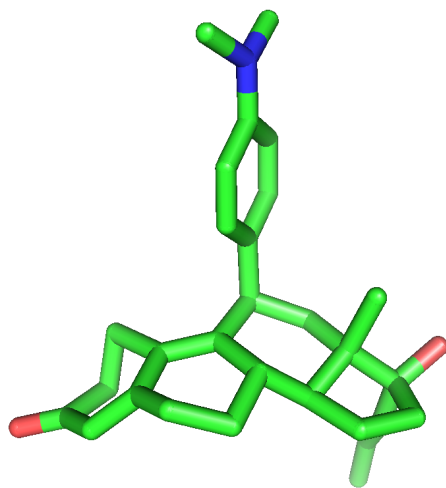
testosterone



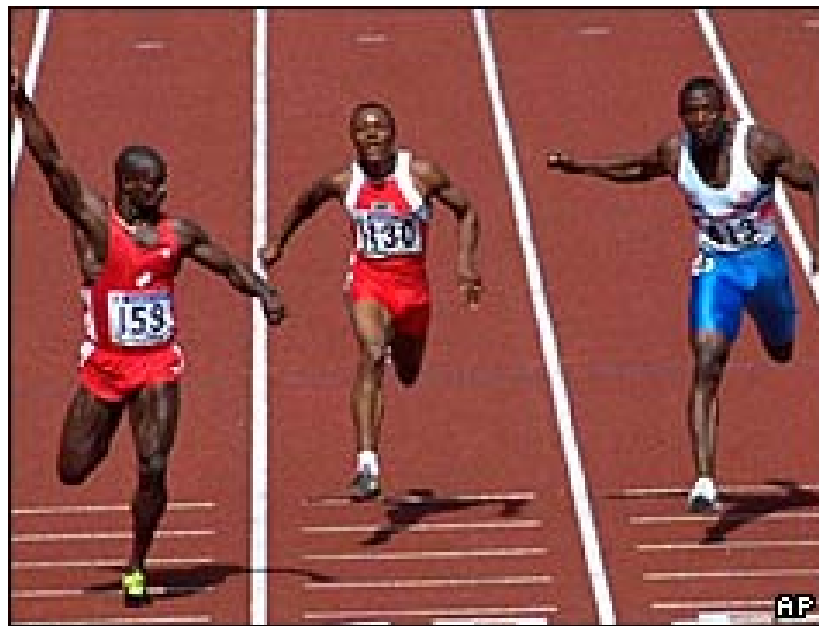
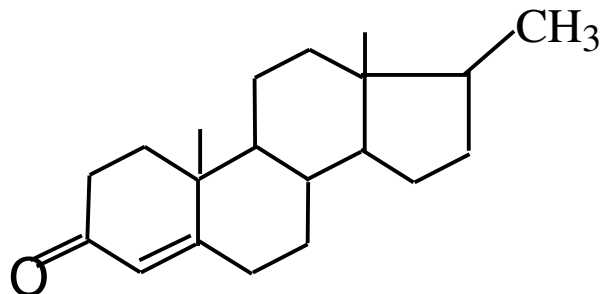
progesterone



mifepristone
or RU486
or Mifegyne™



DMT



NEW PRODUCTS convert into the exact same esters of testosterone found in the body after oral administration of the following anabolic steroids:

- *Methandrostenolone (Dianabol®)
- *Boldenone (Equi-gan®)
- *Stanozolol (Winstrol V®)
- *Drostanolone (Masteron®)

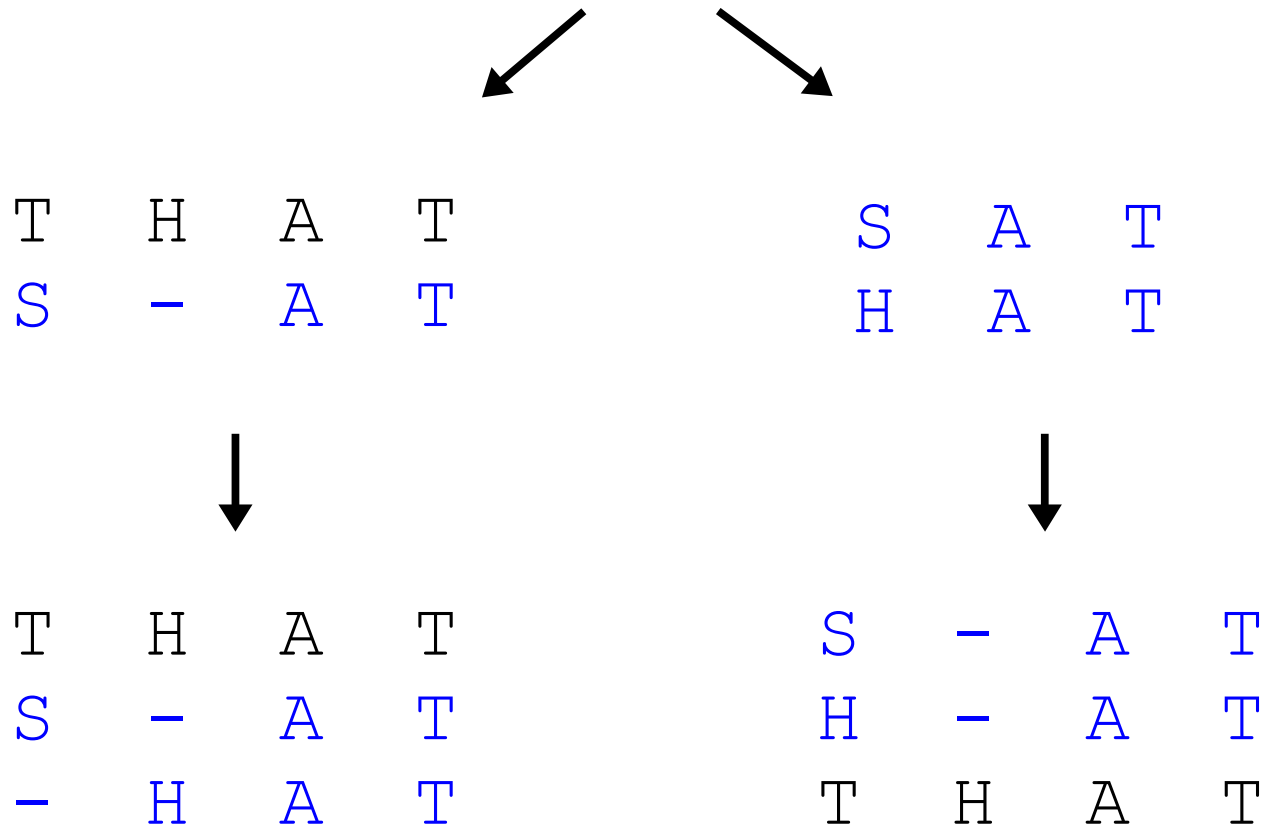


- *Methandrostenolone is the active ingredient in the anabolic steroid DIANABOL®
- *Boldenone is the active ingredient in the anabolic steroid EQUI-GAN®
- *Stanozolol is the active ingredient in the anabolic steroid WINSTROL V®
- *Drostanolone is the active ingredient in the anabolic steroid MASTERON®

GAIN INCREDIBLE MUSCLE SIZE AND STRENGTH!

The final alignment depends on its “history”

Example 1: Fixation of gaps in the alignment



The pairwise sum approximation for the score has no simple probabilistic interpretation

Probabilistic score, two sequences: $\log [P(x_i, y_j) / q_{x_i} q_{y_j}]$

Probabilistic score, three sequences: $\log [P(x_i, y_j, z_k) / q_{x_i} q_{y_j} q_{z_k}]$

$\neq \log [P(x_i, y_j) / q_{x_i} q_{y_j}] + \log [P(x_i, z_k) / q_{x_i} q_{z_k}] + \log [P(y_j, z_k) / q_{y_j} q_{z_k}]$



As if x descended from y AND z

There is no single parsimonious ancestor shared by all 3 pairs...

$((x_i \cup y_j) \cap (x_i \cup z_k) \cap (y_j \cup z_k) = \emptyset \text{ in general })$

The pairwise sum approximation: tree-based weighting

Sequence to align:

H A D

Existing alignment:

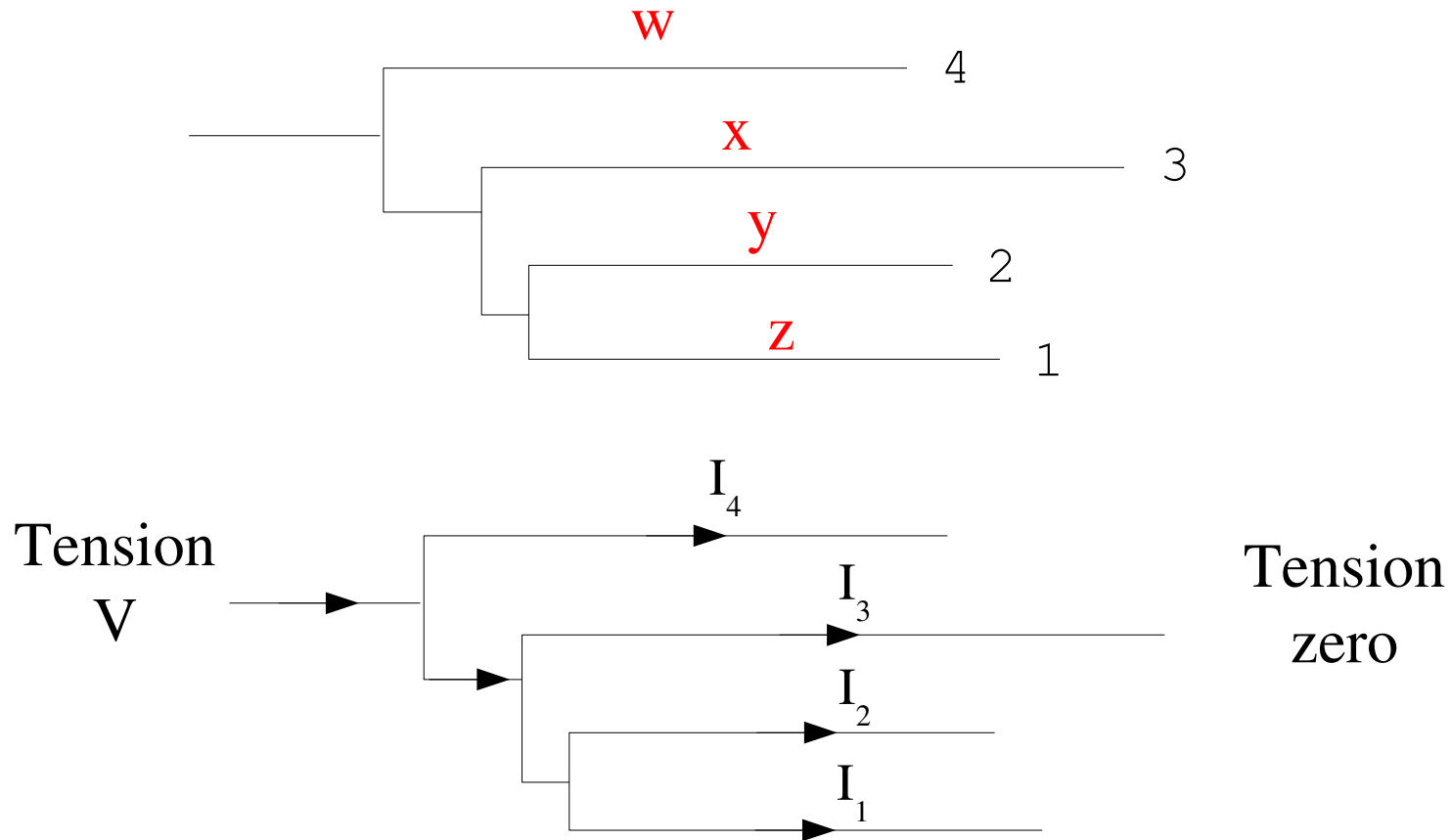
S T A R
S T I R
P - A R



				H			
S	-1	x	0.5				
S	-1	x	0.5				
P	-2	x	1				
				<hr/>			
total	-3						

The two upper sequences are fairly redundant: **reduced weight**

Tree-based weighting: an... electric method



Resistances = branch lengths (w, x, y, z)

Weight of a leaf = electric current

Where should the list of homologues stop?

#	ID Swissprot	Name	Description	Score	E	% Identity	Length aligned
1	P15207	ANDR_RAT	Androgen receptor.	162	1e-40	100	73
6	P19091	ANDR_MOUSE	Androgen receptor.	162	1e-40	100	73
14	Q63449	PRGR_RAT	Progesterone receptor	136	1e-32	80	72
17	P06401	PRGR_HUMAN	Progesterone receptor	136	1e-32	80	72
21	P08235	MCR_HUMAN	Mineralocorticoid receptor	136	1e-32	79	72
33	P04150	GCR_HUMAN	Glucocorticoid receptor	131	3e-31	77	72
41	Q9YH32	ESR2_ORENI	Estrogen receptor beta	99	3e-21	58	72
42	Q9YH33	ESR1_ORENI	Estrogen receptor alpha	98	4e-21	55	72
:	:	:	:	:	:	:	
:	:	:	:	:	:	:	
343	Q9N4Q7	NH13_CAEEL	Nuclear hormone receptor nhr-13	54	8e-08	39	66
344	Q23294	NH11_CAEEL	Nuclear hormone receptor nhr-11	54	8e-08	42	66
345	O45460	NH54_CAEEL	Nuclear hormone receptor nhr-54	54	1e-07	37	67
346	Q09565	NH20_CAEEL	Nuclear hormone receptor nhr-20	51	7e-07	34	66
347	Q09587	NH22_CAEEL	Nuclear hormone receptor nhr-22	45	5e-05	32	66
349	P17672	E75B_DROME	Ecdysone-induced protein 75B	40	0.001	37	47

The E-value: see exercise



The procedure above is commonly used, and implemented in a standard program: ClustalW

Main refinements used in ClustalW and elsewhere:

Ideas?

Multiple sequence alignment is still a research topic

The procedure above is commonly used, and implemented in a standard program: ClustalW

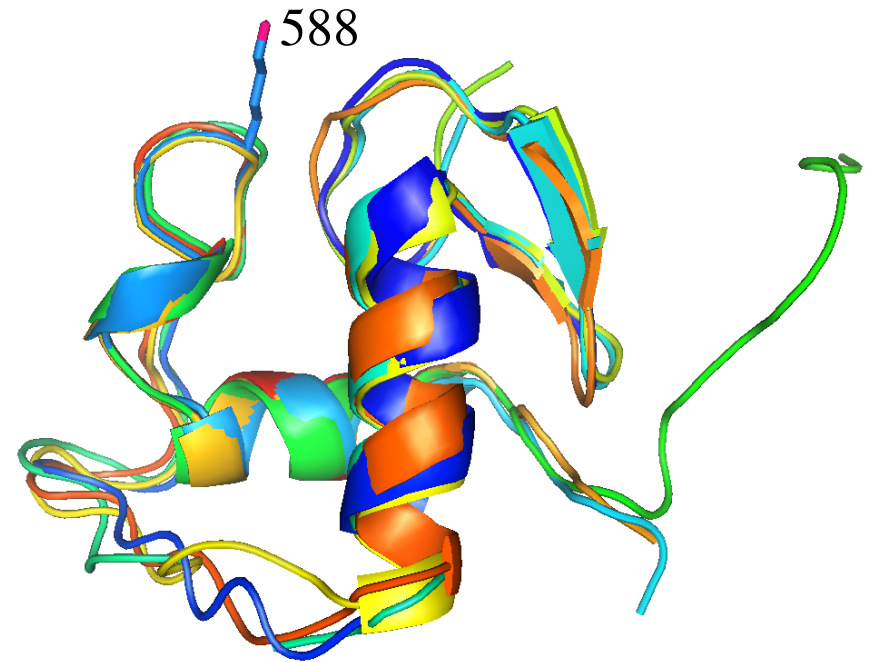
Main refinements used in ClustalW and elsewhere:

- position-dependent gap scores**
- use of different scoring matrices in different parts of the tree**
- schemes to iteratively reshuffle and refine the alignment**

Multiple sequence alignment is still a research topic

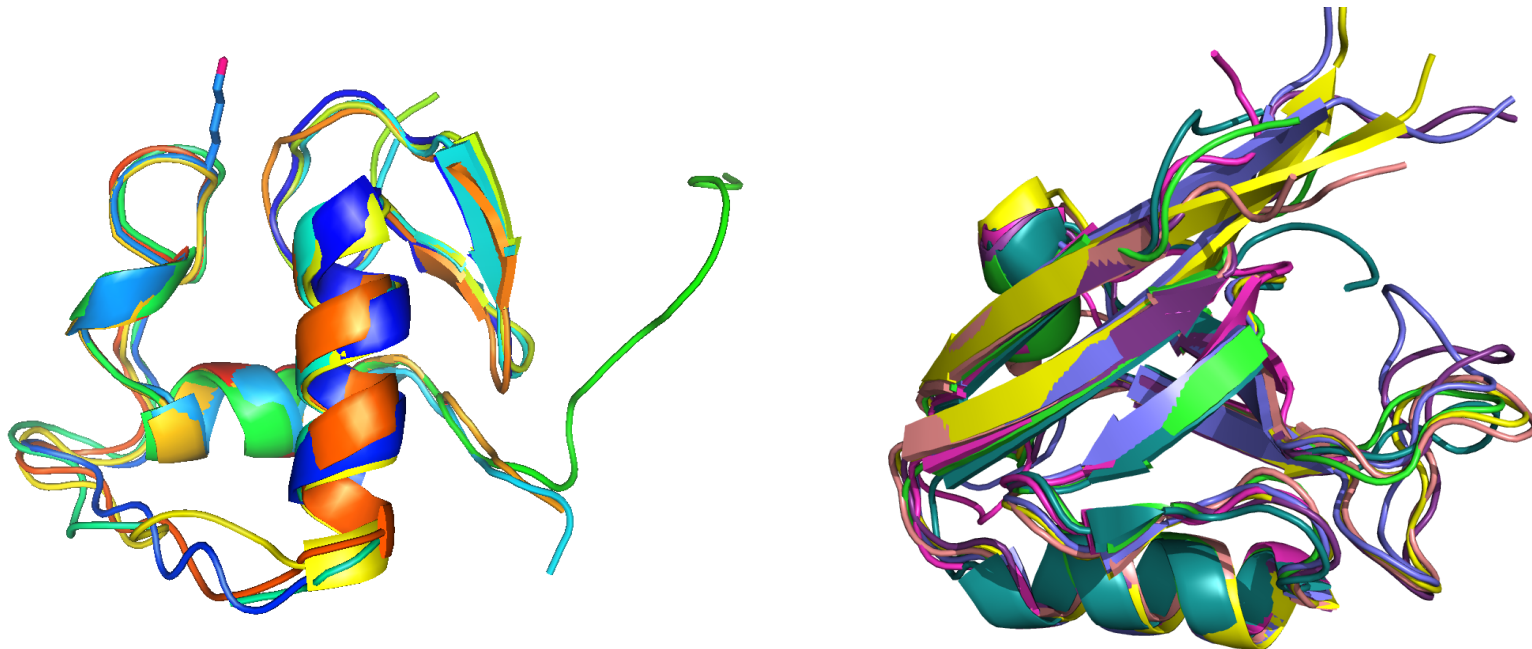
Experimental testing is needed

- Site-directed mutagenesis of conserved residues
- Detecting an interaction with a substrate or inhibitor
- Determination of the 3D-structure!

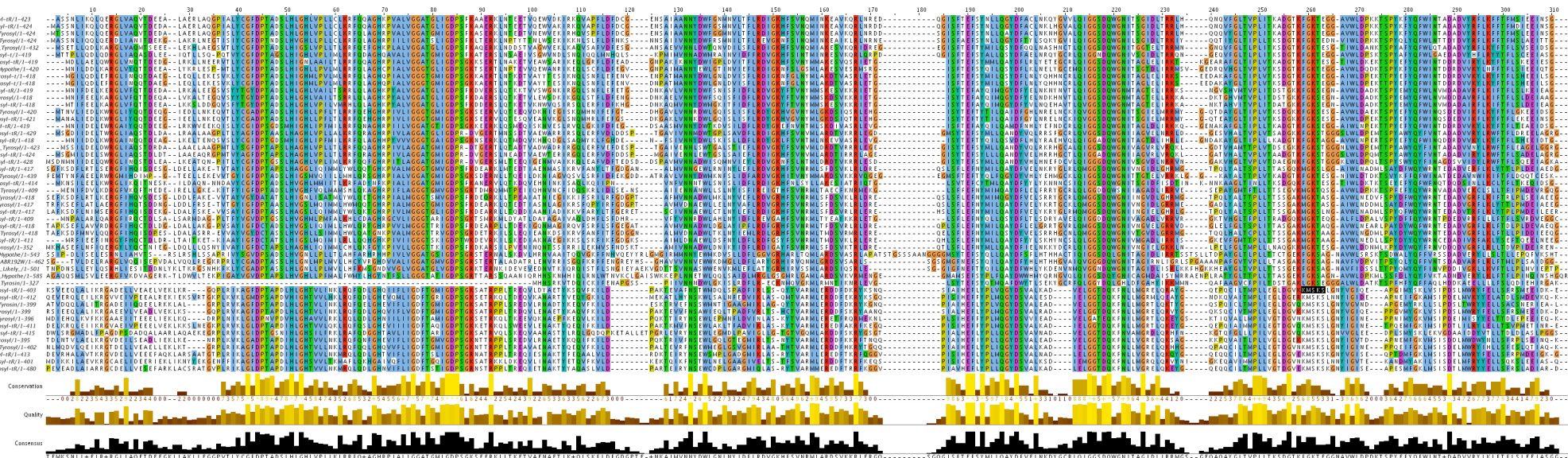


Alignment as a structural model

- Proteins have a specific 3D structure, stabilized by the physical properties of their sequences
- An alignment can be seen as a physical model without explicit evolutionary content: the “structure alignment” problem
- We can adopt this view even if the structures are unknown



Alignment of 30 tyrosyl-tRNA synthetase catalytic domains



Database of protein alignments: one per structural family

<http://pfam.sanger.ac.uk>

**Multiple alignments:
some applications**

Identifying conserved and thus functional residues

Alignment of 340 aminoacyl-ARNt synthetases,
showing conserved residues only

```

Consensus .....KQ.....R.Y.I...FR.E.....EF..ID...AF.....
1  SYD_SHIF .....KQ.....R.Y.I...FR.E.....-EF..ID...sF....-.....
21 SYD_BUC   .....KQ.....k.Y.I...FR.E.....-EF..ID...sF....-.....
41 SYD_BAC   .....KQ.....R.Y.v...FR.E.....-EF..ID...sF....-.....
61 SYD_SYN   .....KQ.....R.Y.I...FR.E.....-EF..lD...sF....-.....
81 SYD_HEL   .....KQ.....k.f.I...FR.E.....-EF..ID...sF....-.....
101 SYD_BR   .....KQ.....R.f.I...FR.E.....EF..lD...sF....-.....
121 SYK1_S    .....Kr....-...R.f.I...FR.E...-...EF..me...Ay....-.....
141 SYD_HA   .....KQ.....R.f.v...FR.E.....E...D...y.....
161 SYK_LI    .....Kr....-...k.Y.I...FR.E...-...EF..le...Ay....-.....
181 SYK_ST    .....Kr....-...k.Y.I...FR.E...-...EF..Ie...Ay....-.....
201 SYK_RA    .....Kr....-...R.Y.I...FR.E...-...EF..me...Ay....-.....
221 SYK_NE    .....Kr....-...R.f.I...FR.E...-...EF..Ie...AF....-.....
241 SYK_CH    .....Kk....-...R.Y.I...FR.E...-...EF..Ie...Ay....-.....
261 SYK3_H    .....Kr....-...f.l...FR.E...-...EF..le-----
281 SYK3_B    .....Kr....-...Y.I...FR.k...-...EF..le-----
301 SYH_ME    .....e.....R.Y...FR.E...-...EF..m...-----
321 SYN_CL    .....e.....-...Y...FR.E.....EF..Ie...AF....-.....
Consensus .....KQ.....R.Y.I...FR.E.....EF..ID...AF.....

```

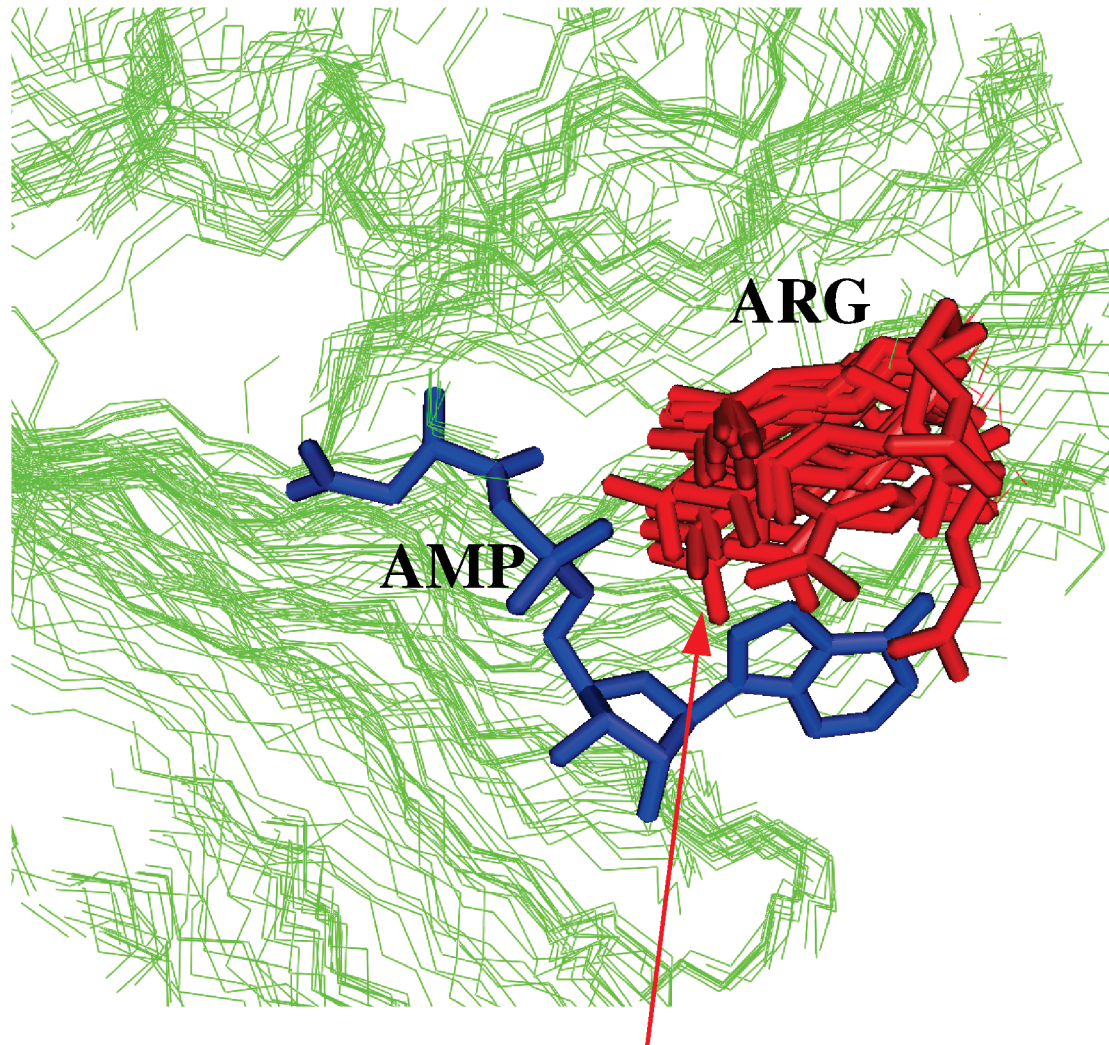
Identifying functional residues

Alignment of 340 aminoacyl-ARNt synthetases,
showing conserved residues only

Consensus	KQ	R.Y.I	FR.E	EF	ID	AF
1 SYD_SHIF	KQ	R.Y.I	FR.E	-EF	ID	sF
21 SYD_BUC	KQ	k.Y.I	FR.E	-EF	ID	sF
41 SYD_BAC	KQ	R.Y.v	FR.E	-EF	ID	sF
61 SYD_SYN	KQ	R.Y.I	FR.E	-EF	lD	sF
81 SYD_HEL	KQ	k.f.I	FR.E	-EF	ID	sF
101 SYD_BR	KQ	R.f.I	FR.E	EF	lD	sF
121 SYK1_S	Kr	-R.f.I	FR.E	-EF	me	Ay
141 SYD_HA	KQ	R.f.v	FR.E	E	D	y
161 SYK_LI	Kr	-k.Y.I	FR.E	-EF	le	Ay
181 SYK_ST	Kr	-k.Y.I	FR.E	-EF	Ie	Ay
201 SYK_RA	Kr	-R.Y.I	FR.E	-EF	me	Ay
221 SYK_NE	Kr	-R.f.I	FR.E	-EF	Ie	AF
241 SYK_CH	Kk	-R.Y.I	FR.E	-EF	Ie	Ay
261 SYK3_H	Kr	-f.l	FR.E	-EF	le	-----	-----
281 SYK3_B	Kr	-Y.I	FR.k	-EF	le	-----	-----
301 SYH_ME	e	R.Y	FR.E	-EF	m	-----	-----
321 SYN_CL	e	-Y	FR.E	EF	Ie	AF
Consensus	KQ	R.Y.I	FR.E	EF	ID	AF

Identifying functional residues

Structural alignment of 30 aminoacyl-ARNt synthetases



Consensus K Q R . Y . I . . F R . E E F . . I D . . A F

Template-based 3D structure prediction

Estrogen receptor example

1) align its sequence with homologues, including at least one of known 3D structure: the template

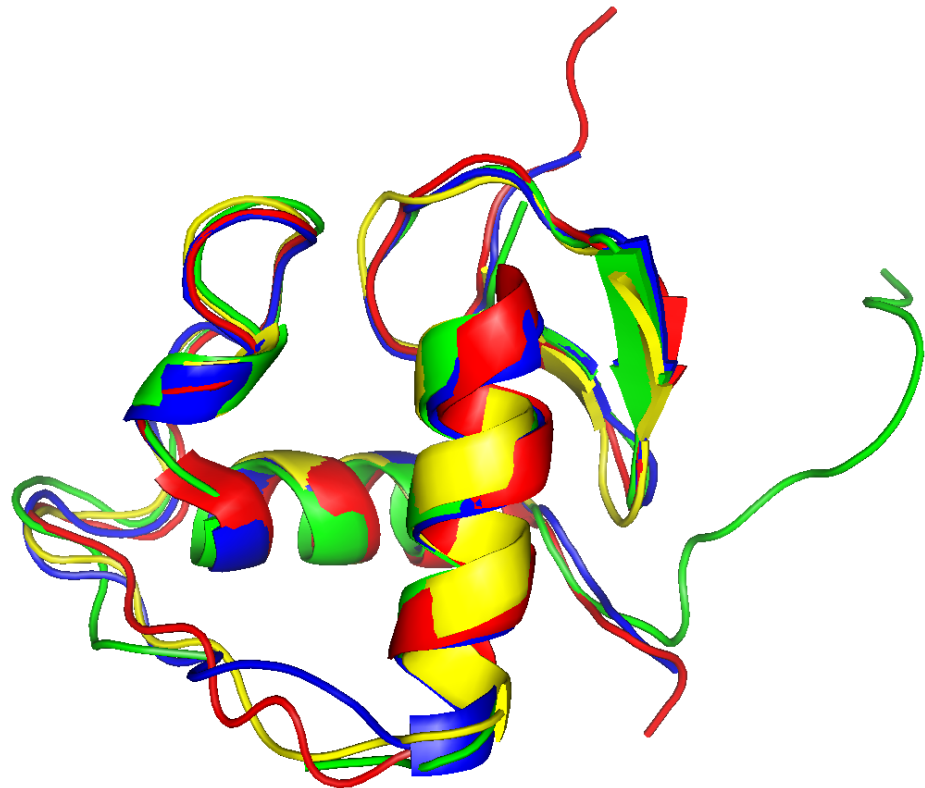
androgen	CLICGDEASGAHYGALTCGSCKVFFKRAAEGKQKYL-CASRNDCTIDKFRRKNCPSCLRLRKCYEAGMTLGA
Rev Erb	CKVCGDVASGFHYGVLACEGCKGFFRRSIQQNIQYKRCLKNENCSTIVRINRNRCQQCRFKKCLSVGMSRD-
glucocorticoid	CLVCSDEASGCHYGVLTCGCKAFFKRAVEGQHNYL-CKYEGKCIIDKIRRKNCPCRYRKCLQAGMNLEA
retinoic acid	CAICGDRSSGKHYGVYSCEGCKGFFKRTVRKDLTYT-CRDNKDCLIDKRQRNRCQYCRYQKCLAMGM---
estrogen	CAVCNDYASGYHYGVWSCEGCKAFFKRSIQGHNDYM-CPATNQCTIDKNRRKSCQACRLRKCYEVGMMKG-

Template-based 3D structure prediction

Estrogen receptor example

2) adopt a sensible or consensus main chain trace

3) build sidechains

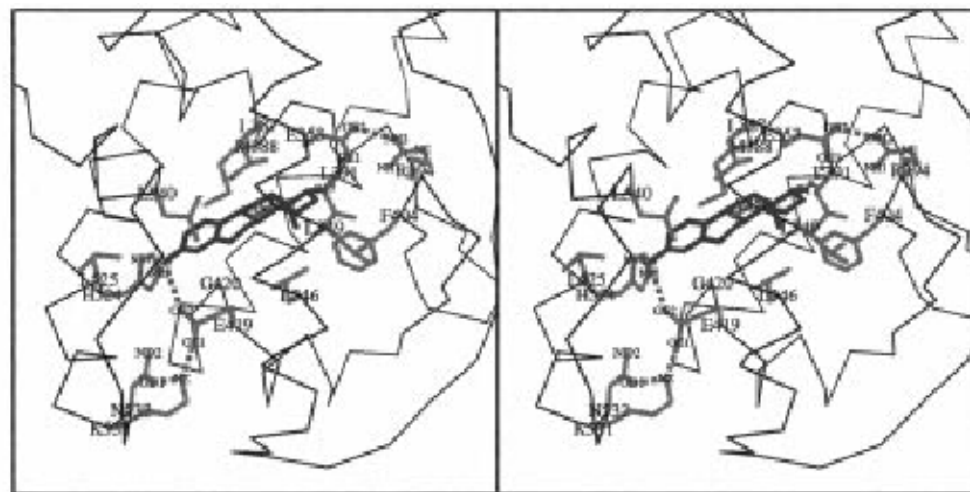


Template-based 3D structure prediction

Estrogen receptor example

2) adopt a consensus main chain trace

3) build sidechains



Wurtz et al (1998)
J Medicinal Chem

In this example, the model was good enough to explain estrogen binding

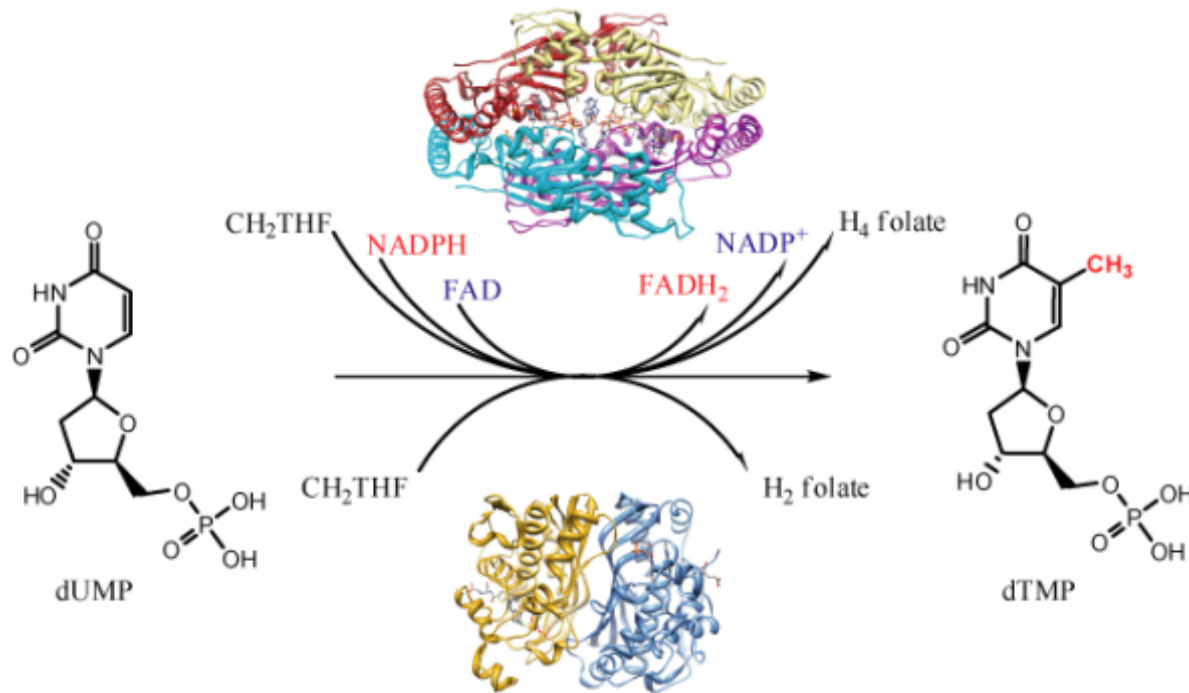
Identifying orthologues and possible drug targets

A strategy to attack a pathogen is to identify an essential protein that has no human orthologue, making it a potential drug target.



Identifying orthologues and possible drug targets

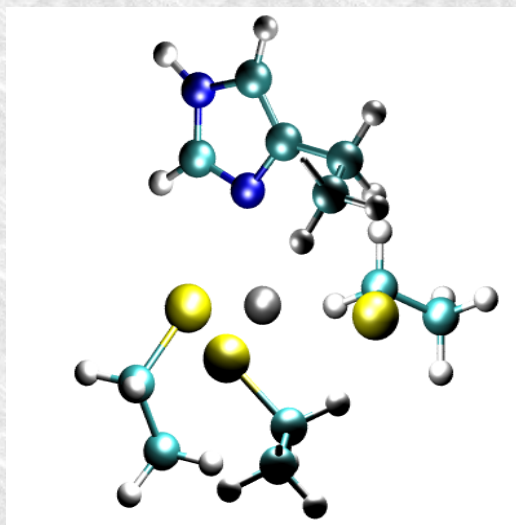
A strategy to attack a pathogen is to identify an essential protein that has no human orthologue, making it a potential drug target.



Distinct structures and chemical reactions in human thymidylate synthase ThyA (top) and the essential **analogue ThyX (bottom) present in many pathogenic bacteria, including *Helicobacter pylori***

Science (2002) 297:105

Identifying patterns and functional sites



Identifying functional residues

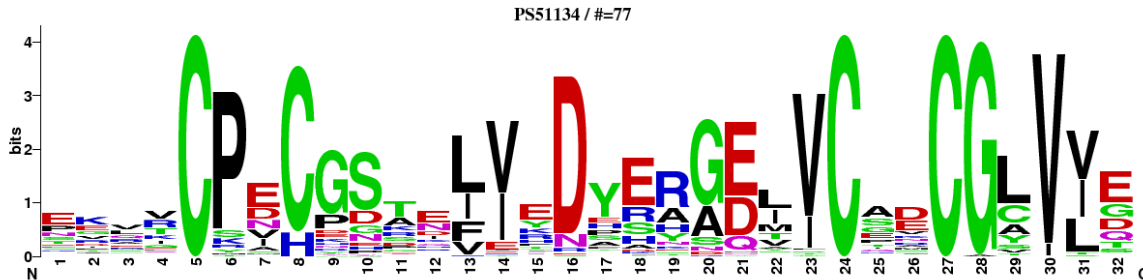
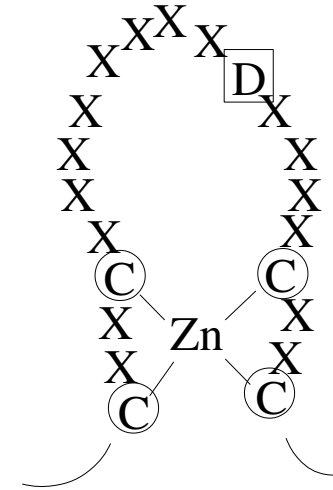
Alignment of 340 aminoacyl-ARNt synthetases,
showing conserved residues only

```
Consensus .....KQ.....R.Y.I...FR.E.....EF..ID...AF.....
1 SYD_SHIF .....KQ.....R.Y.I...FR.E.....-EF..ID...sF.....-.....
21 SYD_BUC .....KQ.....k.Y.I...FR.E.....-EF..ID...sF.....-.....
41 SYD_BAC .....KQ.....R.Y.v...FR.E.....-EF..ID...sF.....-.....
61 SYD_SYN .....KQ.....R.Y.I...FR.E.....-EF..lD...sF.....-.....
81 SYD_HEL .....KQ.....k.f.I...FR.E.....-EF..ID...sF.....-.....
101 SYD_BR .....KQ.....R.f.I...FR.E.....EF..lD...sF.....-.....
121 SYK1_S .....Kr....-...R.f.I...FR.E...-...EF..me...Ay.....-.....
141 SYD_HA .....KQ.....R.f.v...FR.E.....E...D...y.....
161 SYK_LI .....Kr....-...k.Y.I...FR.E...-...EF..le...Ay.....-.....
181 SYK_ST .....Kr....-...k.Y.I...FR.E...-...EF..Ie...Ay.....-.....
201 SYK_RA .....Kr....-...R.Y.I...FR.E...-...EF..me...Ay.....-.....
221 SYK_NE .....Kr....-...R.f.I...FR.E...-...EF..Ie...AF.....-.....
241 SYK_CH .....Kk....-...R.Y.I...FR.E...-...EF..Ie...Ay.....-.....
261 SYK3_H .....Kr....-...f.l...FR.E...-...EF..le-----
281 SYK3_B .....Kr....-...Y.I...FR.k...-...EF..le-----
301 SYH_ME .....e.....R.Y....FR.E...-...EF..m...-----
321 SYN_CL .....e.....-...Y....FR.E.....EF..Ie...AF.....-.....
Consensus .....KQ.....R.Y.I...FR.E.....EF..ID...AF.....
```

Divergence within the family may be too great to use a multiple alignment
→ structural, biochemical studies needed to infer functional pattern

Another protein example:
zinc fingers (TFIIB class)

Database of conserved patterns
in proteins: **ProSite** 
<http://prosite.expasy.org>

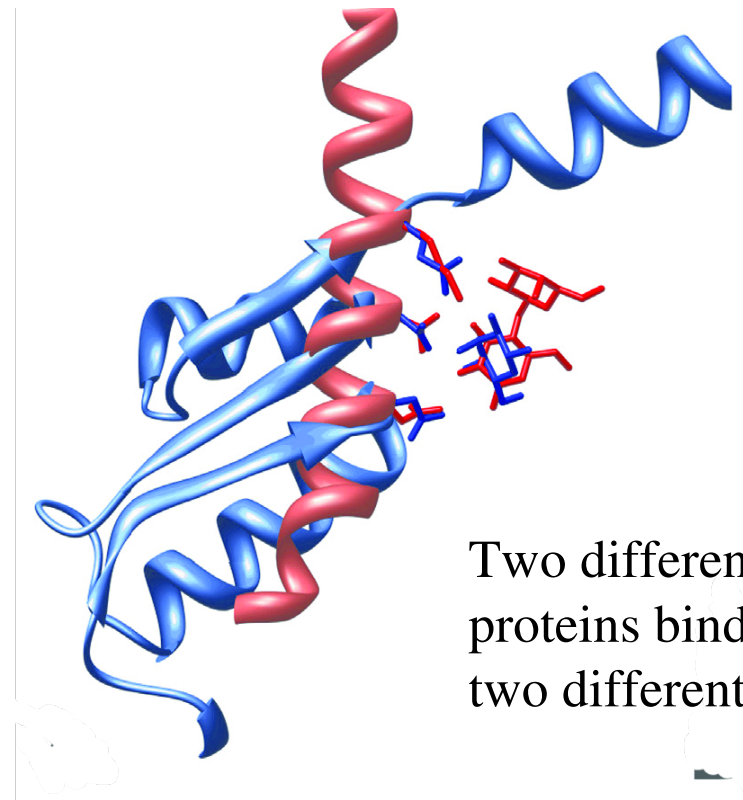


Database of conserved patterns
in proteins: **ProSite**
<http://prosite.expasy.org>



**Functional patterns sometimes
arise from convergent, not
divergent evolution... so that a
3D pattern is relevant, not a 1D
pattern**

BMC Bioinformatics 2009, 10:182



Two different
proteins binding
two different sugars

Coding, signals and patterns in DNA

- Transcription signals
- Restriction sites
- Splicing sites
- Regulation sites
- Methylation sites

To detect short patterns, special techniques are needed

Restriction sites

Restriction enzymes cut DNA at specific sites

5' GAATTC 3' EcoRI
3' CTTAAG 5'

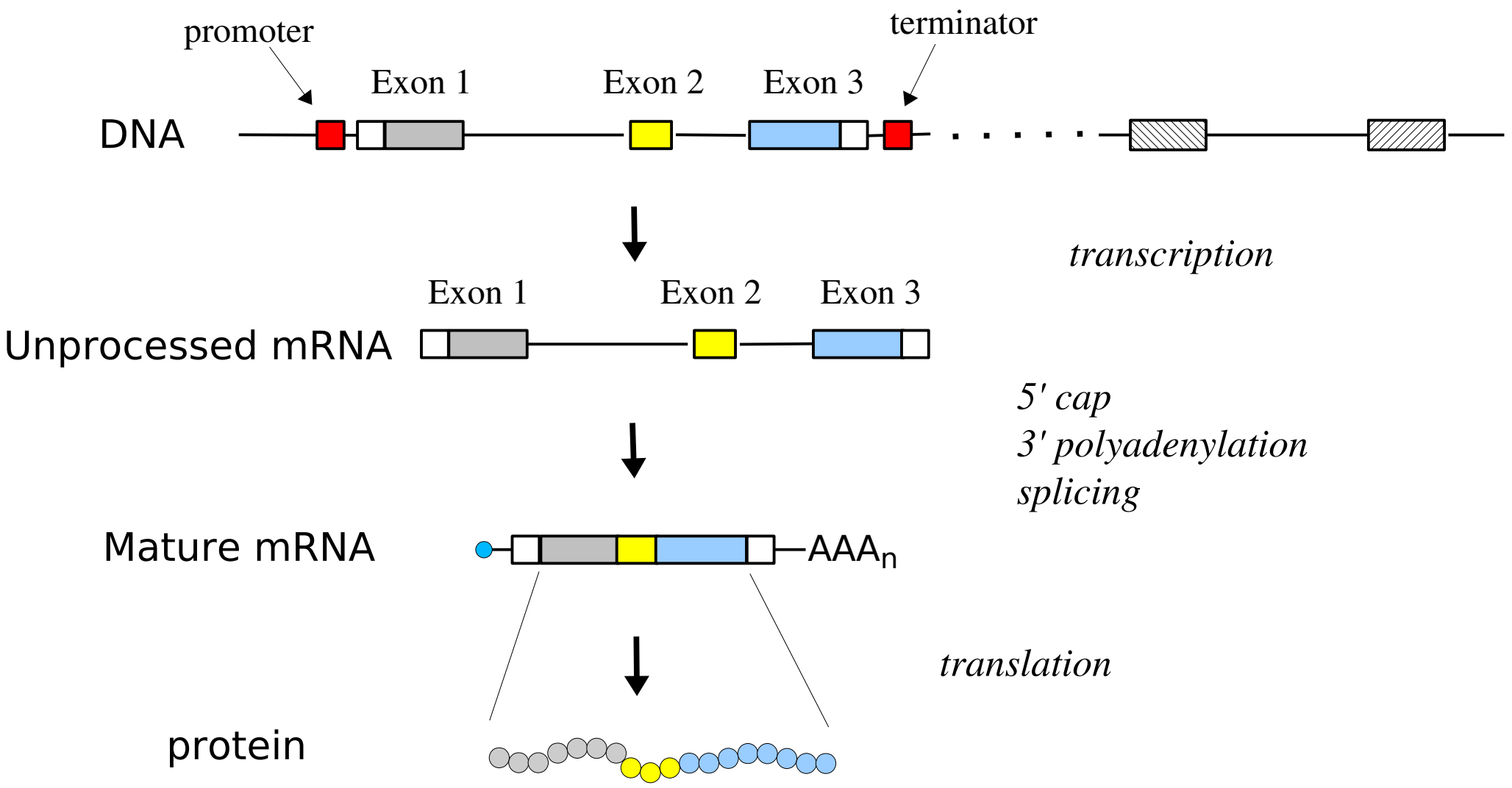
5' GTTAAC 3' HpaI
3' CAATTG 5'

5' GGATCC 3' BamHI
3' CCTAGG 5'

Antiviral defence; a tool for molecular biotechnology

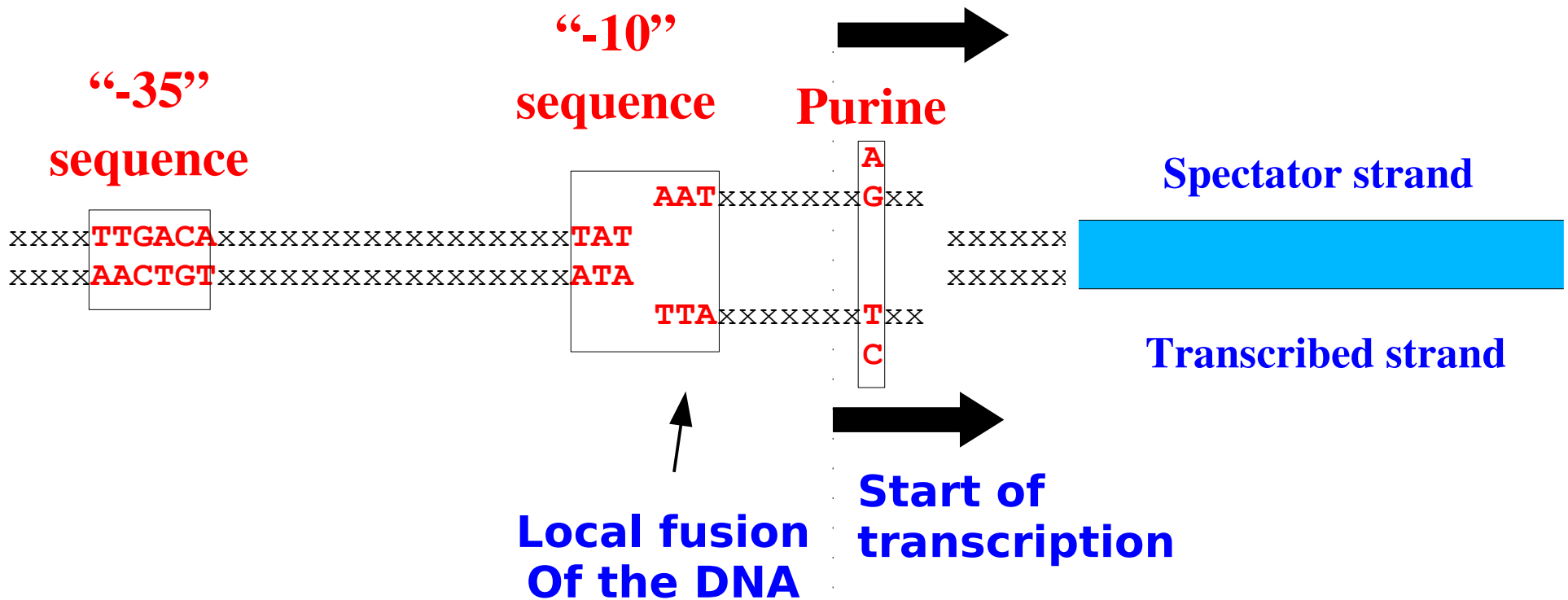


Transcription signals: promoters



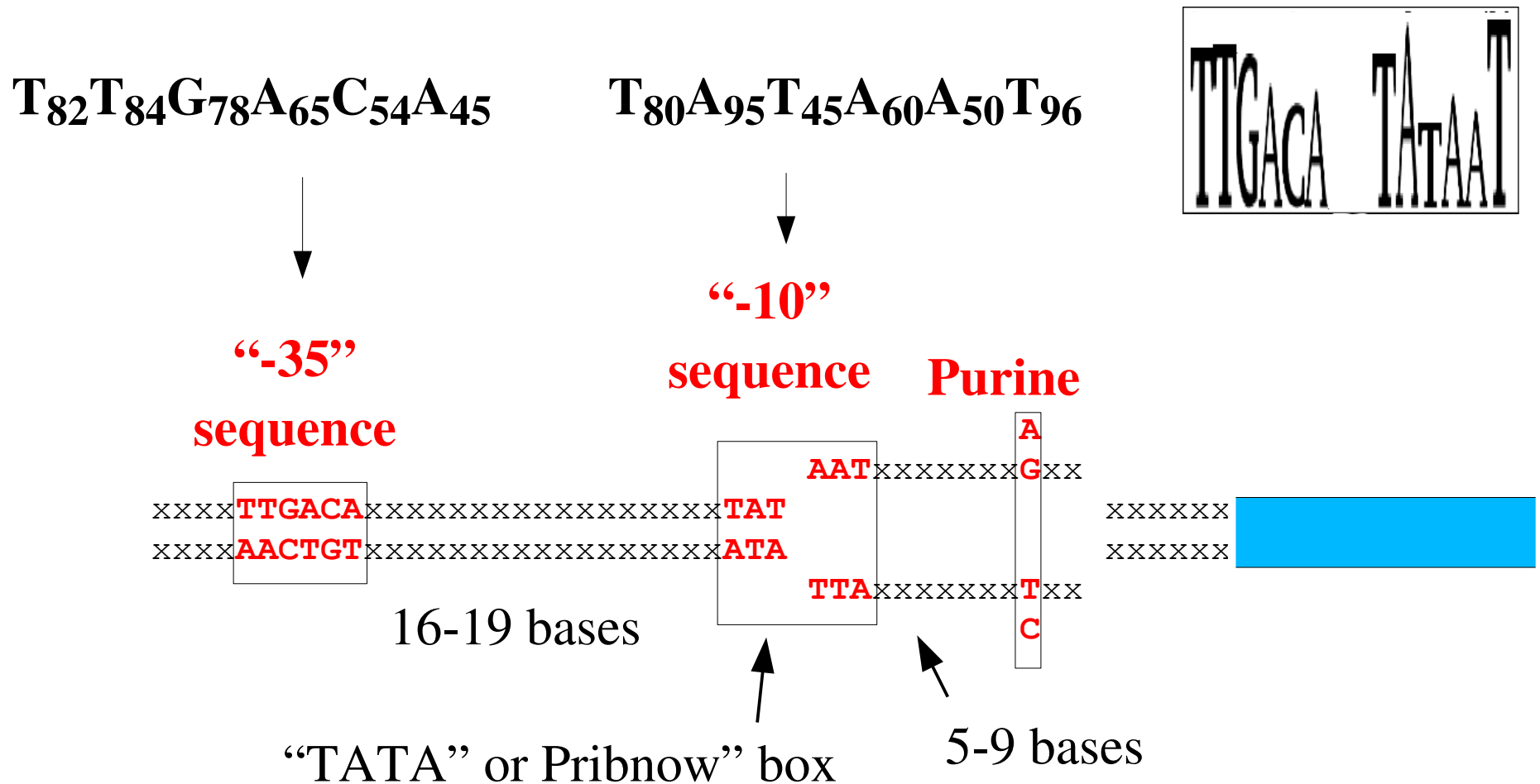
Transcription signals: promoters

3 regions upstream of the START allow RNA polymerase binding and the initiation of transcription of an ORF (open reading frame): they form the **promoter**.

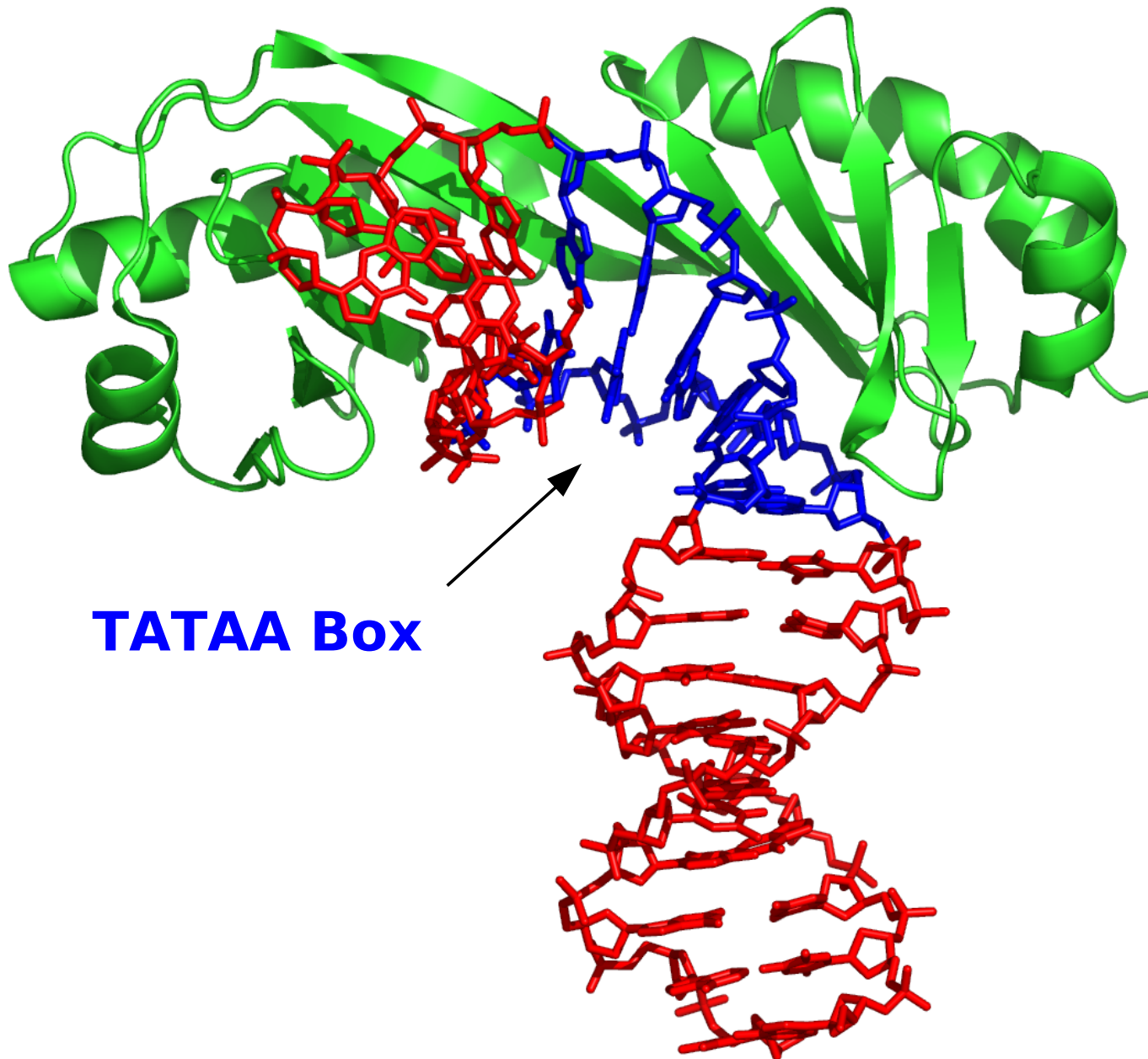


Transcription signals: promoters

The conservation (%) of the -35 and -10 consensus sequences is:



TATAA Binding Protein (TBP)



TATAA Box

Splicing sites

intron

Before splicing ...GAGGAGGG|gtgagtgtg.....cctctccccag|CTGCTC...

After splicing

...GAGGAGGG|CTGCTC...

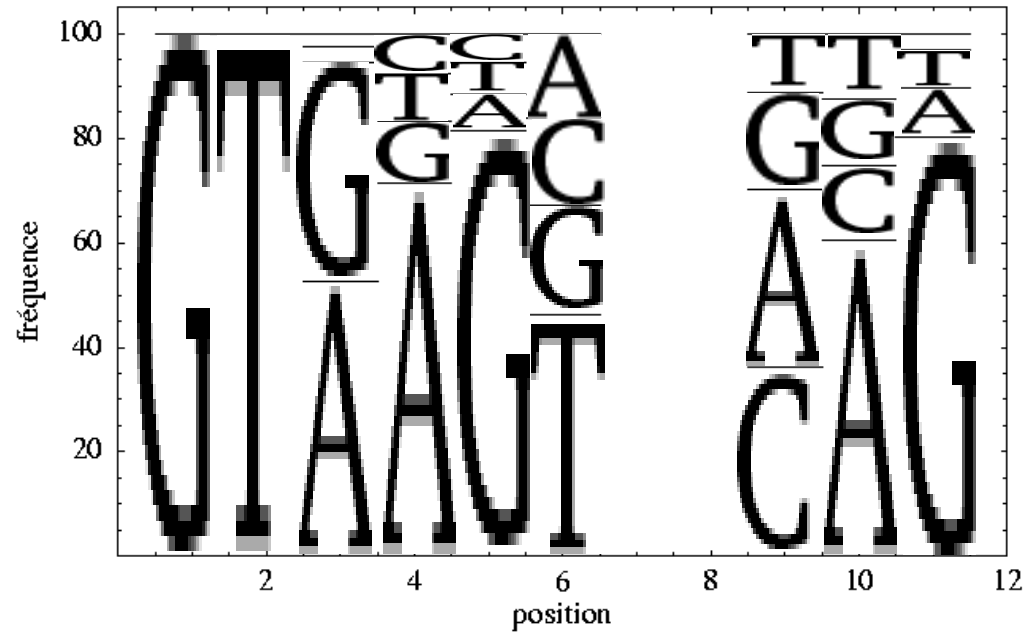
“GT-AG” splicing junction

Splicing sites

intron

...GAGGAGGG|gtgagtgtg.....cctctccccag|CTGCTC...

“GT-AG” type
consensus
sequence



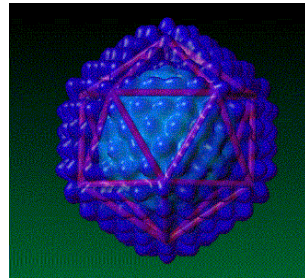
98.7% of known splicing sites

The discovery of splicing



Richard Roberts
New England Biolabs

1977
Adenovirus



Chicken
albumin
1977



Pierre Chambon
Institut de Génétique
et Biologie Moléculaire
et Cellulaire, Strasbourg



Philip Sharp, MIT

Ovalbumin gene is split in chicken DNA; BREATHNACH, R;
MANDEL, JL; CHAMBON, P (1977) NATURE 270:314-319

Prosite database of functional patterns in proteins:

<http://www.expasy.org/prosite>

Zinc coordination sites in nuclear hormone receptors

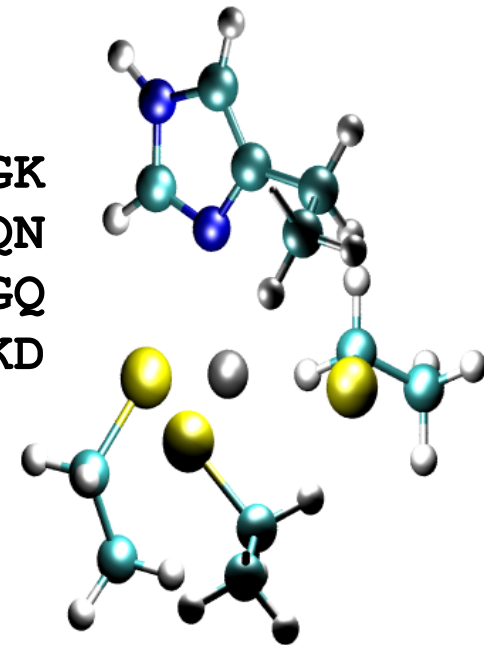
androgen
Rev Erb
glucocorticoid
retinoic acid

* * * *

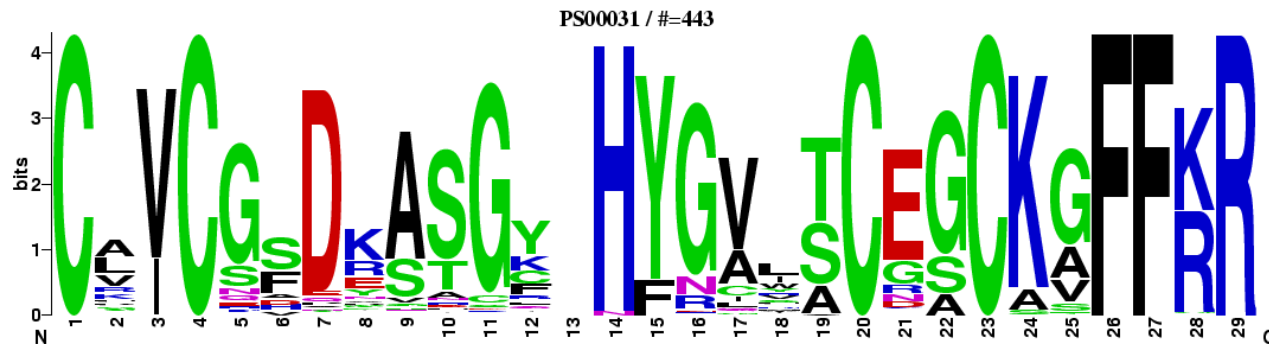
CLICGDEASGAHYGALTCGSCKVFFKRAAEGK
CKVCGDVASGFHYGVLACEGCKGFFRRSIQQN
CLVCSDEASGCHYGVLTCGCKAFFKRAVEGQ
CAICGDRSSGKHYGVYSCEGCKGFFKRTVRKD

consensus
"PROSITE"
pattern

CxxCxDxxxxxHFxxxxCxxCxxFFxR
E NY



Sequence
profile



Representation by a “regular expression” or search pattern

consensus
pattern

C _{xx} C _x D _{xxxxxx} H F _{xxxx} C _{xx} C _{xx} F F _x R
E NY

C-x(2)-C-x-[DE]-x(5)-[HN]-[FY]-x(4)-C-x(2)-C-x(2)-F-F-x-R

grep format:

C.{2}C.[DE].{5}[HN][FY].{4}C.{2}FF.R

perl format:

C\w{2}C\w[DE]\w{5}[HN][FY]\w{4}C\w{2}FF\wR

<http://www.expasy.org/prosite>

Pattern matching: a naïve algorithm

To recognize the polyadenylation motif **AAUAAA** in the sequence below:

AAAUAUUAAAUAUUAGACGAAUAAAAGUAUAUUU

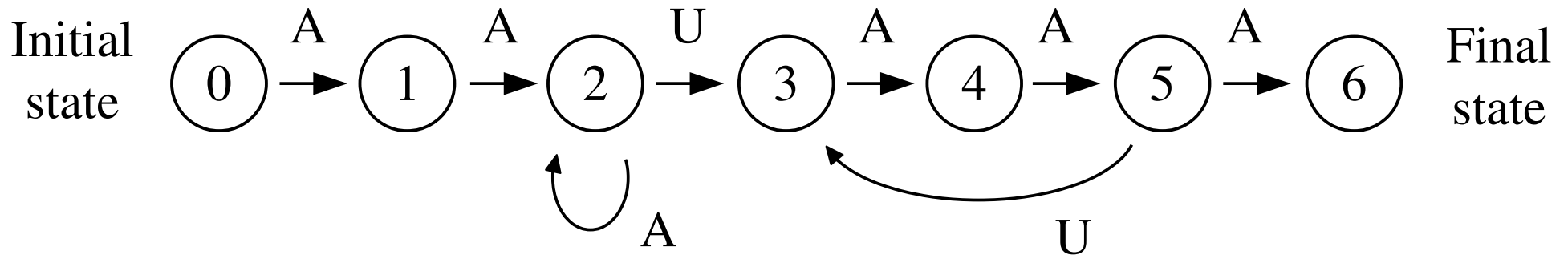
we have the naïve procedure:

- Find an A; assume it is the beginning of the pattern
- Test the following nucleotides one by one
- If the match fails, start over with the next A

The method is far from optimal ($O(nm)$)

Pattern matching with automata

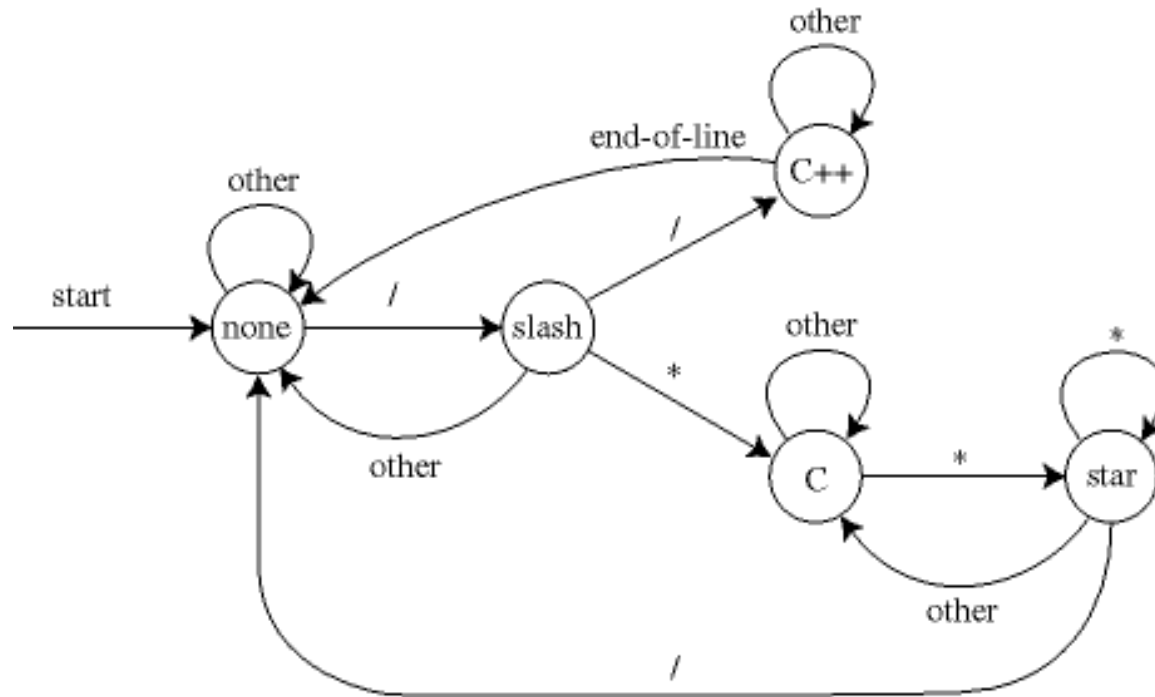
The following automaton recognizes the polyadenylation motif **AAUAAA**



All other transitions go back to 0 (not shown).

Example: target: AAAUAUUAAAUAUUAGACGAAUAAAAGUAUAUUU
state: 1223400122340010100123456

Automata: a vast area in computer science...



An automaton to remove comments from Java programs

Biological Sequence Analysis. Durbin et al, 1998, Cambridge U Press, chapter 9
Algorithms on strings, trees and sequences. D Gusfeld, 1997, Cambridge U Press

Multiple alignments: a 3 stage progressive method

**Applications: functional residues, structure prediction,
phylogeny**

Short patterns and signals: other techniques are needed