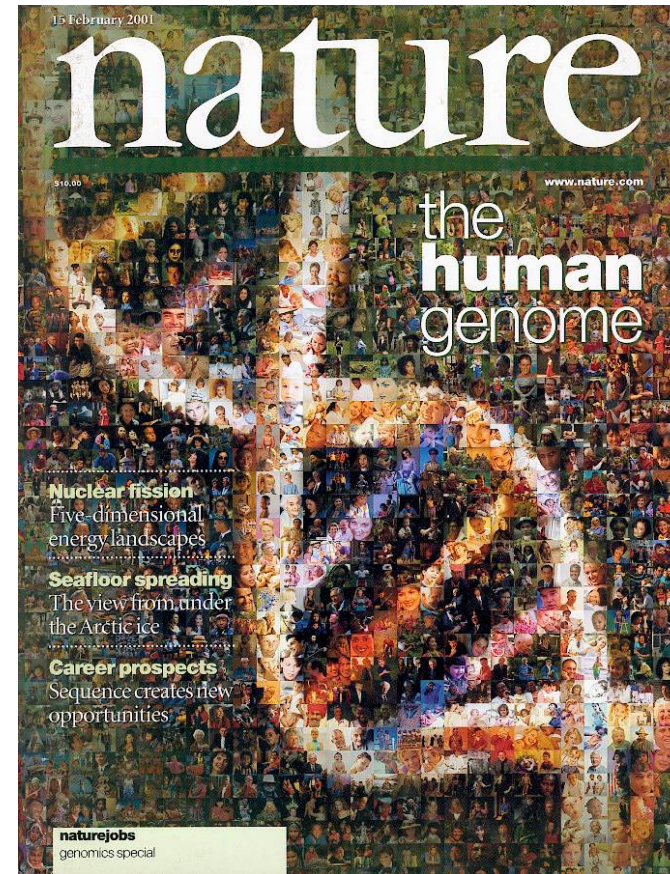

Bioinformatics

- Scientific Computing Lecture, UNICC

Marina Alexandersson

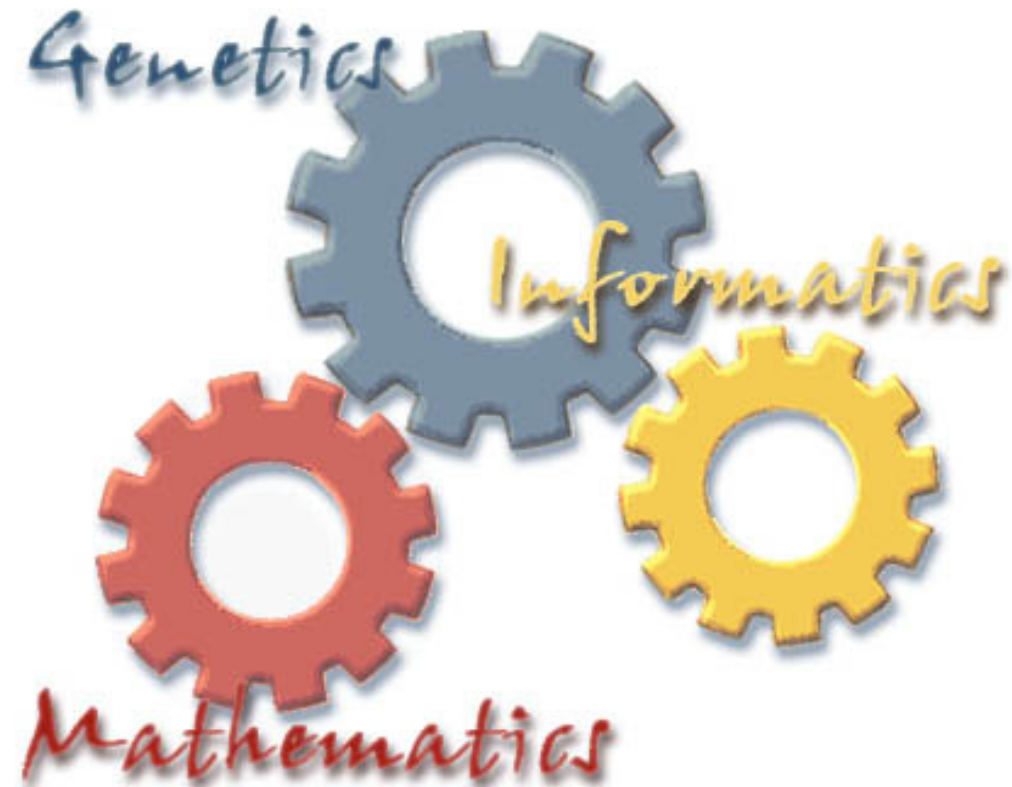
Göteborg, March 25, 2004

The human genome



Fraunhofer **CHALMERS**
Research Centre
Industrial Mathematics

What is bioinformatics?



Fraunhofer **CHALMERS**
Research Centre
Industrial Mathematics

Where is bioinformatics used?

DNA analysis



Drug

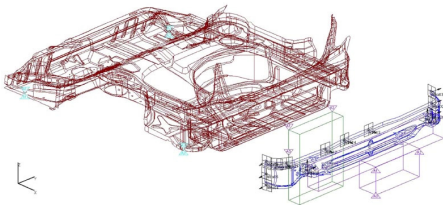
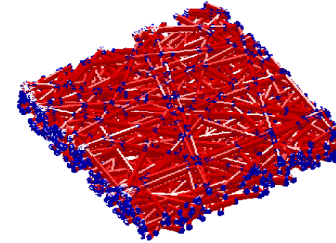
- DNA sequencing
- Sequence analysis
- Genes involved in diseases
- Protein structure and function
- Drug targets
- Lead compound
- Clinical trials



Technologies at the Fraunhofer-Chalmers centre

➤ Computational physics and fatigue

- materials, mechanics, electromagnetics
- fatigue life and load analysis



➤ Optimization and quality engineering

- design, geometry, quality management
- logistics, scheduling, production

➤ System theoretical methods

- biostatistics, bioinformatics, systems biology
- telecommunications, signal and image analysis
- finance and insurance



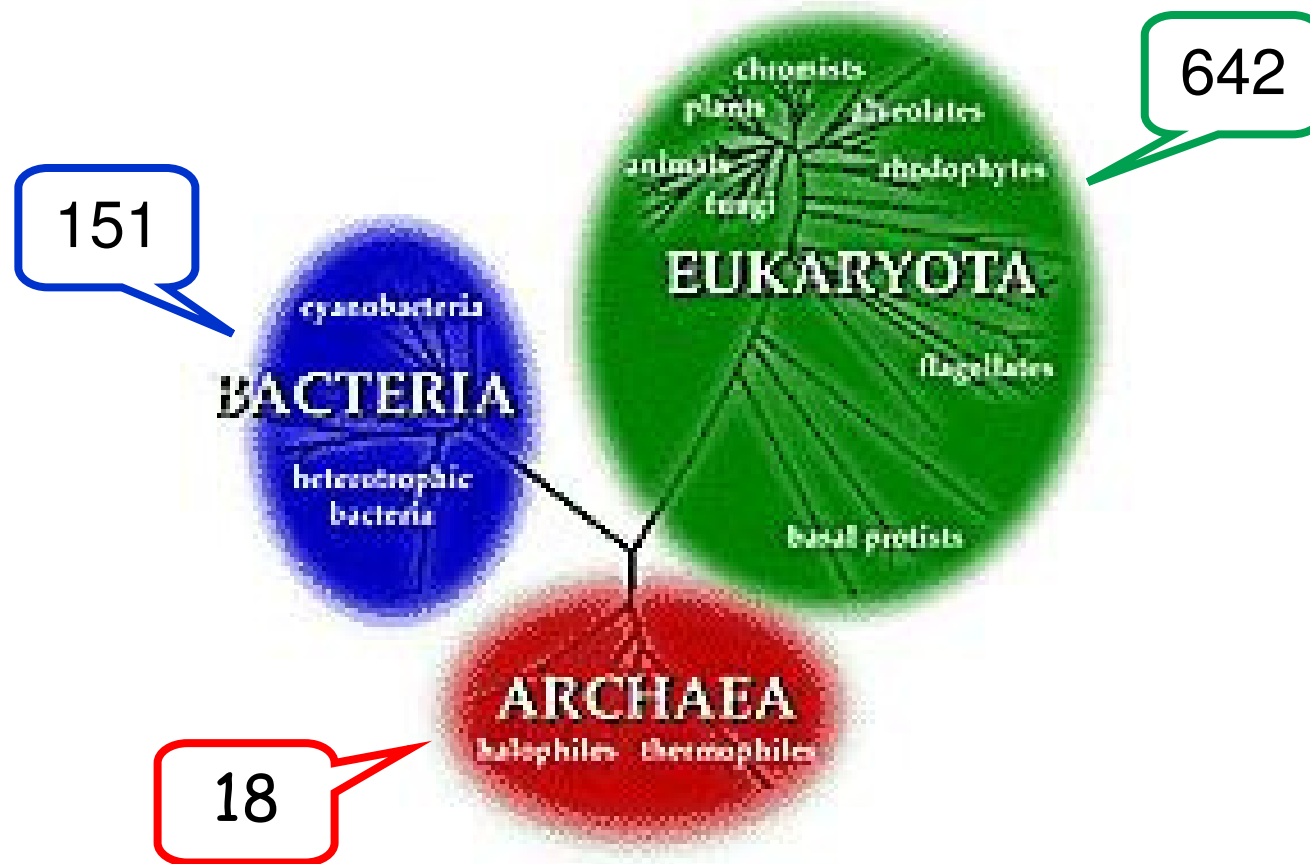
Bioinformatics

Research, development and application of computational methods and tools to expand the use of biological and medical data, including retrieving, storing, organizing, analyzing and visualizing such data.

- **Sequence analysis** to classify sequence data and identify functional elements
- **Comparative genomics** to detect similarities and differences in related organisms
- **Functional genomics** to infer biological structure and function from primary sequence
- **Proteomics** to identify, characterize and quantify proteins active in certain pathways, tissues, cells or organs
- **Biostatistics and statistical genetics** to analyze disease genotypes and phenotypes, detect environment interactions etc



Sequenced genomes



More to come...



Fraunhofer **CHALMERS**
Research Centre
Industrial Mathematics

Whole genome analysis - Why?

- Complete gene and protein sets
- Primary sequence of all genes
- Sequence relationships between genes and proteins
- Function of new proteins
- Transcriptional level of all genes
- Understanding metabolic pathways
- Trace disease genes
- ...

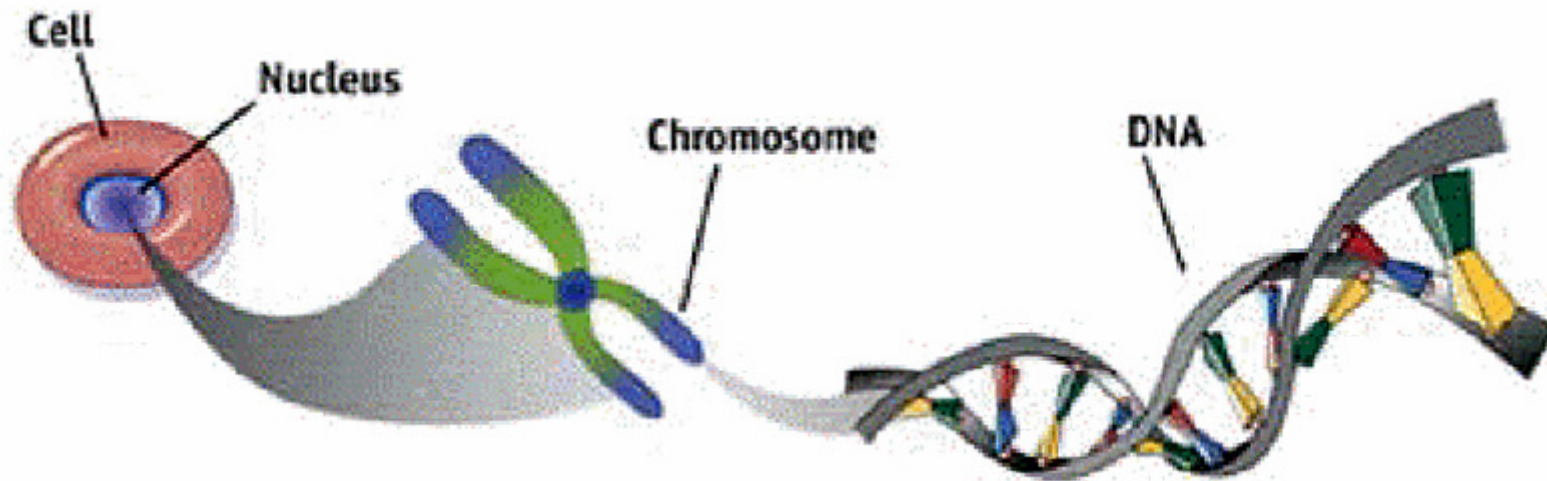
Whole genome analysis

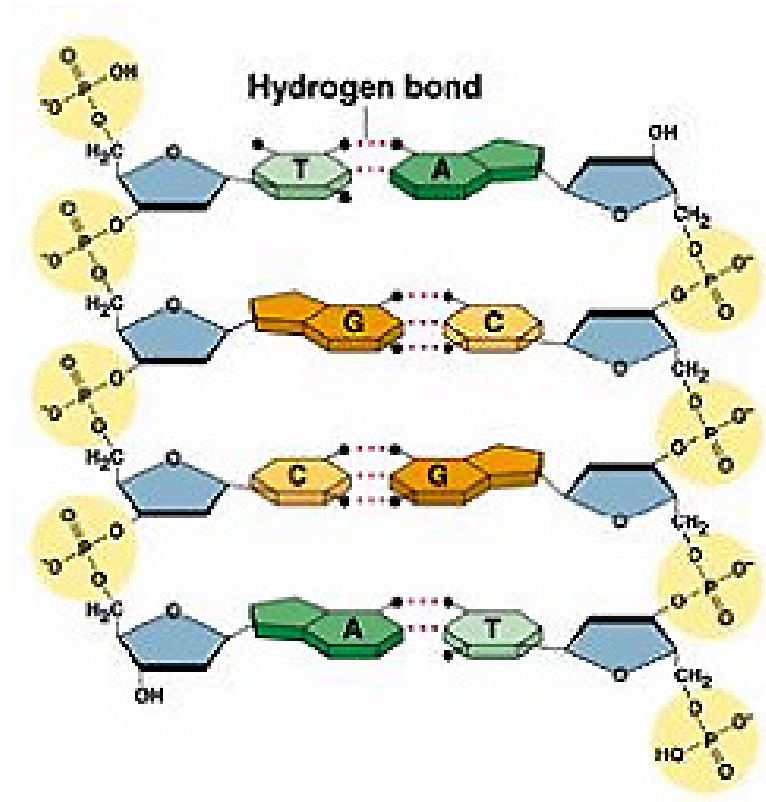
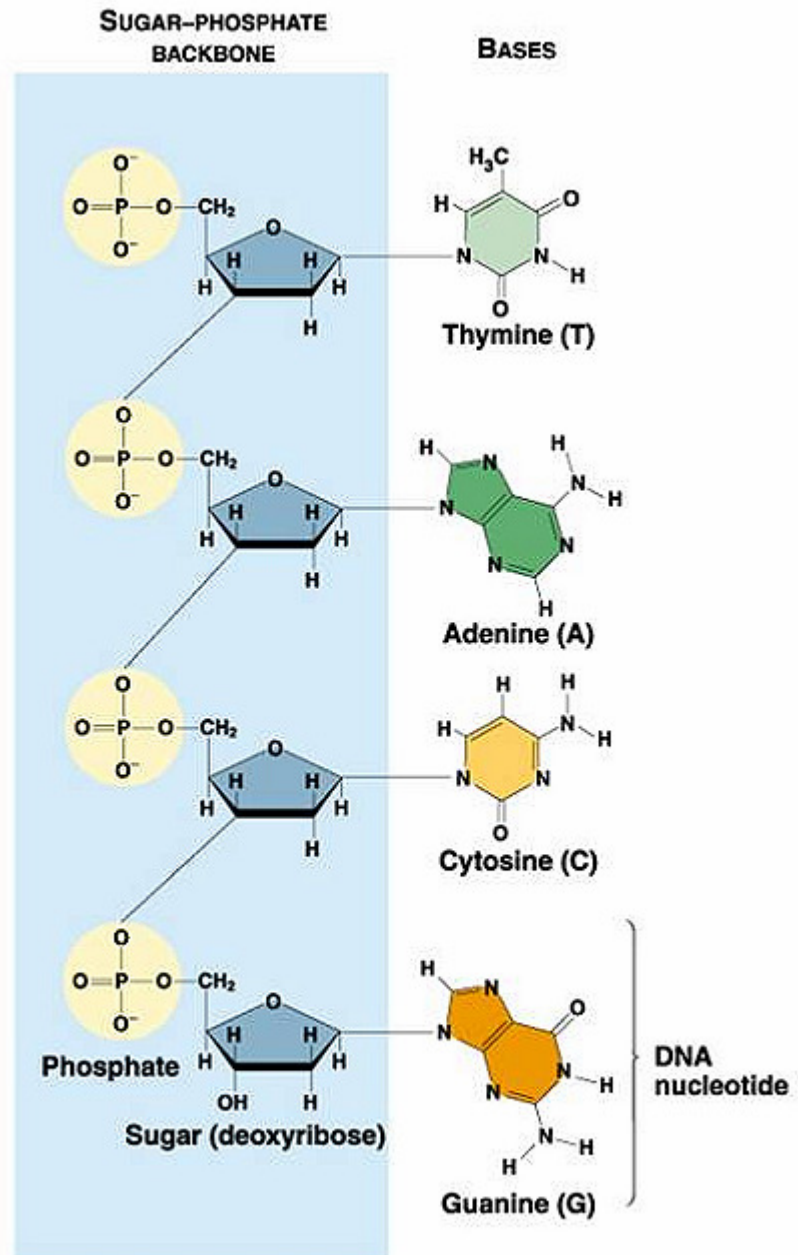
- Gene finding
- Sequence comparisons
- Regulatory region identification

Gene finding

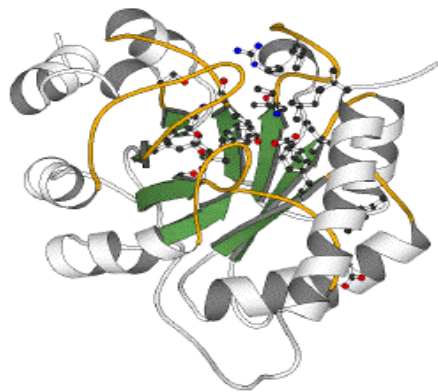
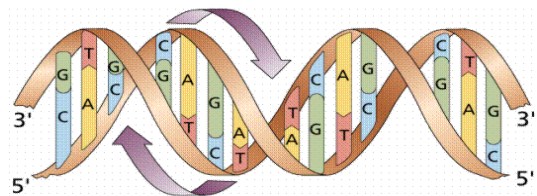
- *Generalized HMMs*

Basic genetics





Gene expression



DNA

transcription

RNA

translation

Protein

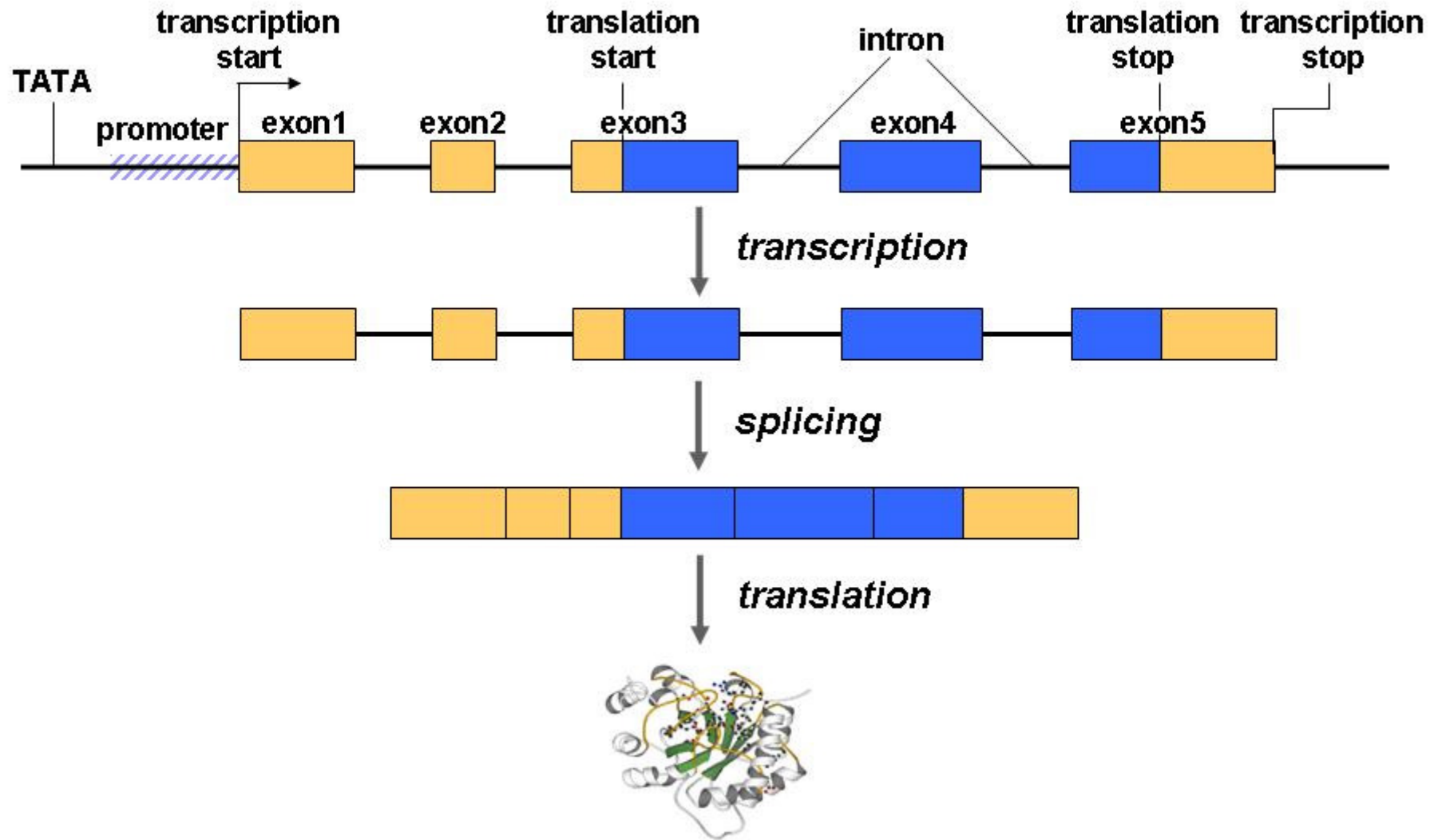
CCTGAGCCAAC TATTGATGAA

CCU**GAG**CCA**ACU**AUU**GAU**GAA

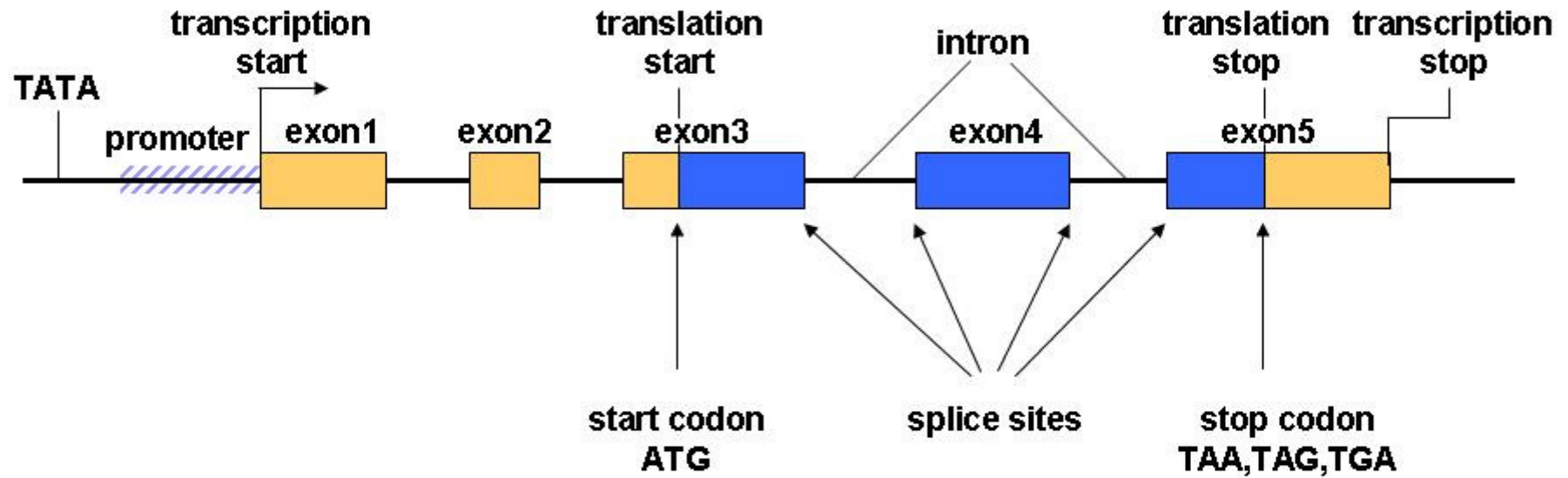
PEPTIDE



Gene structure



Finding genes



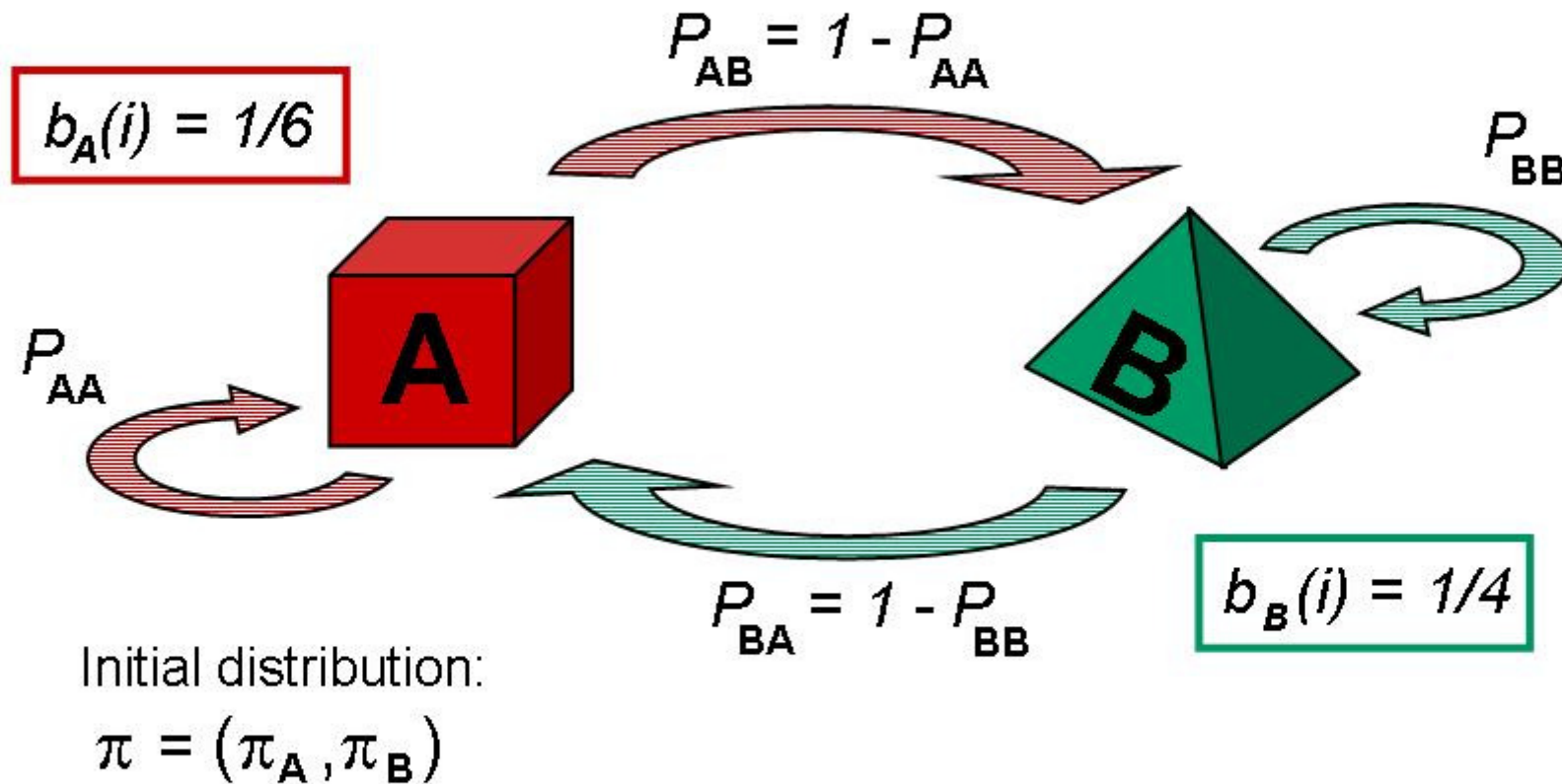
> HSKIIIBE, Human gene for casein kinase II subunit beta (EC 2.7.1.37).
ggggctgagatgtaaattagaggagctggagaggagtgcctcagagtttgggttgctttaagaaagggg
gggtccgaattctcccgtgggtggagggccgaatgtgggaggaggaggataccagaggcaggggaagga
gaacttgagcttactgacactgtcttttctagctgacgtgaagatgagcagctcagaggagggtgc
ctggatttccctgggtctgtgggctccgtggcaatgaattctctgtgaagtgagttctcttcaacctcc
ctacttgccagcttcacatatcttcccaccagacgttcttccacatatccactctacactgttctct
aaagcttttatgggagagagtgtagggtgaactagggagagacacaagtactctgctgagttgggagtg
agaaacaagcacaacagatgcagttgtgttgatgataaggcatcacttagagcattttgcccagggtcaa
agatgaggattttgatatgggttccctcttggcttccatgtcctgacaggtggatgaagactacatcca
ggacaaatttaactctactggactcaatgagcaggtccctcactatcgacaagctctagacatgatctt
ggacctggagcctgggtgaggcaccctcaggggtgttttgtgtgtgtgcgtgcactattttctcttcaa
atctctatctacttgccctgaattttgccaaatttcccttgggtctctgatctcttaacccccaaatca
tgctttatttgatccctccacctgactctgtctagttttgtgacgtatatcacttgttctcatgtttt
tgcaagggtcagaagcccagggttctgggtcccatgcccagatgttggatggggtaaggcccaaaagta
gggtgctaggcaaaactgaaatagcccgcagcccctggatattggggcagggcacctaggaaagctgaaaaaca
agtagttgcatttggccgggctgtgggttcagatgaagaactggaagacaacccccaccagagtgacctg
attgagcaggcagccgagatgctttatggatgatccacgcccgcctacatccttaccaaccgtggcatc
gcccagatgggtgaggcctctctgctcctacctgcctcctctctgagcagtaagagacacaggttccctgca
gcaagaagtcatgtttaagccctgtttaaggaagctagctgagaagaggggaagaaccccagaacttgg
ccctgccctaatttggaagaaaggcaacacagaagtttgagagcccactctagtcacagagaagggggcct
ctggacagagttggaaggagtgcgcacagagttggatagggtgggctgogaagggagttgcctctct
ttacatctacctgccaaacccttccatgtatccacctcagttggaaaagtaccagcaaggagactttg
gttactgtcctcgtgtgtactgtgagaaccagccaatgcttccattGgtgagtggtgaagaagggaaa
ggaaagcaccgtgtggcagctcttatgggaaggagtggggctcaacacatggagcctgagtcctgagg
ggaggttaggtaggaatagggggataacctggcctgctgagctctggctgtctcccaggccttccagacat
cccagggtgaagccatgggtgaagctctactgccccaaagtgcattggaatgtgtacacaccaagtcatcaag
acaccatcacacggatggcgcctacttcggcactgggttccctccacatgctctcatgggtgcattcccga
gtaccggcccaagagacctgccaaaccagtttgtgcccaggtagggagcagggagagtcattaaggggtca
aaggaaaggcccaagatccccagagagggggaggacaggggcatggcccttcttgaggctctgcttctcc
cagaatcagggcatctccctgctgagtgactgtgggaaagttatttgattatctgtgcttgagttacct
tatgtagaatgttcttgagctgagaagtgggaaaccacagggcttagctctgagcaggtccatagag
gagctcaggtggggagggtgggaaatgcaggtgactggcagggcctggatggggctcatgctgctgcctct
ctgacctctgccctggcctaggctctacggtttcaagatccatccgatggcctaccagctgcagctcca
agccgccagcaacttcaagagcccagtcagaacgatccgctgatccctccccacctgtcctgcagtc
tttgtctttcccttctttttgccacccttccaggaaccctgtatgggttttagtttaaataaagga
gtcgttatcgtgggtgggaatatgaaataaagtagaagaaaaggccatgagctagctctgctgggtgcttgc
ggaaggggggtggagcgtggccatggaaatcgggctccacggcccagggatgg

> HSKIIIBE, Human gene for casein kinase II subunit beta (EC 2.7.1.37).
ggggctgagatgtaaattagaggagctggagaggagtgcctcagagtttgggttgctttaagaaagggg
ggttcogaattctcccgtgggtggagggccogaatgtgggaggaggaggataccagaggcaggggaagga
gaacttgagcttactgacactgttcttttctagctgacgtgaagatgagcagctcagaggagggtgc
ctggatttccctgggtctgtgggctccgtggcaatgaattctctgtgaagtgagttctcttcaacctcc
ctacttgccagcttcacatatcttcccaccagacgttcttccacatatccactctacactgttctct
aaagcttttatgggagagagtgtagggtgaactagggagagacacaagtactctgctgagttgggagtg
agaaacaagcacaacagatgcagttgtgttgatgataaggcatcacttagagcattttgcccagggtcaa
agatgaggattttgatatgggttccctcttggcttccatgtcctgacaggtggatgaagactacatcca
ggacaaatttaacttactggactcaatgagcaggtccctcactatcgacaagctctagacatgatctt
ggacctggagcctgggtgaggcaccctcaggggtgttttgtgtgtgtgcgtgcactattttctcttcaa
atctctatctacttgccctgaattttgccaaatttcccttgggtctctgatctcttaacccccaaatca
tgctttatttgatcctccacctgactcttgtctagttttgtgacgtatctacttgttctcatgtttt
tgcaagggtcagaagcccagggttcttgggtcccatgcccagatgttggatggggtaaggccccaaaagta
gggtgctaggcaaaactgaaatagcccgcagcccctggataggggcagggcacctaggaaagctgaaaaaca
agtagttgcatttggccgggctgtgggttcagatgaagaactggaagacaacccccaccagagtgacctg
atgagcagggcagccgagatgctttatggattgatccacgcccgtacatccttaccacccgtggcatc
gcccagatgggtgaggcctctctgctcctacctgcctcctcttgagcagtaagagacacaggttccctgca
gcaagaagtcatgtttaagccctgtttaaggaagctagctgagaagaggggaagaaccccagaacttgg
ccctgccctaatttggaaagaaaggcaacacagaagtttgagagcccactctagtcacagagaagggggcct
ctggacagagttggaaggagtgcgcacagagttggataggggtgggctgogaagggagttgcctctct
ttacatctacctgccaaaccccttccatgtatccacctcagttggaaaagtaccagcaaggagactttg
gttactgtcctcgtgtgtactgtgagaaccagccaatgcttccattGgtgagtggtgaagaagggaaa
ggaaagcaccgtgtggcagctcttatgggaaggagtggggctcaacacatggagcctgagtcctgagg
ggaggttaggttaggaatagggggataacctggcctgctgagctctggctgtctcccaggccttccagacat
cccagggtgaagccatgggtgaagctctactgccccaaagtgcattggaatgtgtacacaccaagtcatcaag
acaccatcacacggatggcgcctacttcggcactgggttccctccacatgctctctatgggtgcattcccga
gtaccggcccaagagacctgccaaaccagtttgtgcccaggtagggagcagggagagtcattaagggta
aaggaaaggccccaaagatccccagagagggggaggacaggggcatggcccttcttgaggtctgcttctcc
cagaatcagggcatctccctgctgagtgactgtgggaaagttatttgatatactgtgcttgagttacct
tatgtagaatgttcttgagctgagaagtgggaaccacagaggcttagctctgagcaggtccatagag
gagctcaggtggggagggtgggaatgcaggtgactggcagggcctggatggggctcatgctgctgcctct
ctgacctctgccctggcctaggctctacggtttcaagatccatccgatggcctaccagctgcagctcca
agccgccagcaactcaagagcccagtcacagacgatcgcctgatccctccccacctgtcctgcagtc
tttgtcttttcccttcttttttgccaccttccaggaacctgtatgggttttagtttaaataaagga
gtcgttatcgtgggtgggaatatgaaataaagtagaagaaaaggccatgagctagctctgctgggtgcttgc
ggaaaggggggtggagcgtggccatggaaatcgggctccacggcccagggatgg

Approaches to gene finding

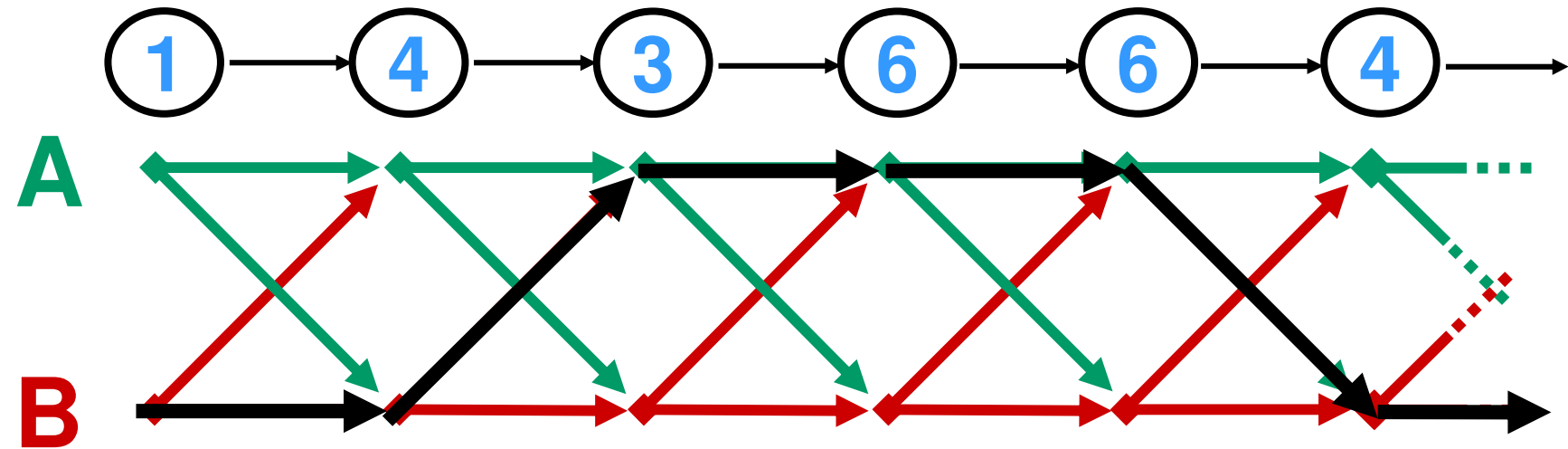
- Homology searches
 - BLAST, Procrustes
- Ab initio
 - Genscan, Genie, GeneID
- Hybrids
 - GenomeScan, GenieEST, Twinscan, SGP, ROSETTA, CEM, TBLASTX, SLAM

A simple Hidden Markov Model (HMM)

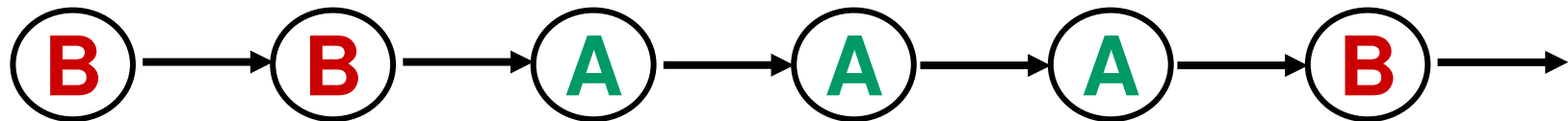


A lattice view

Observed sequence:



Hidden sequence:



Two fundamental problems

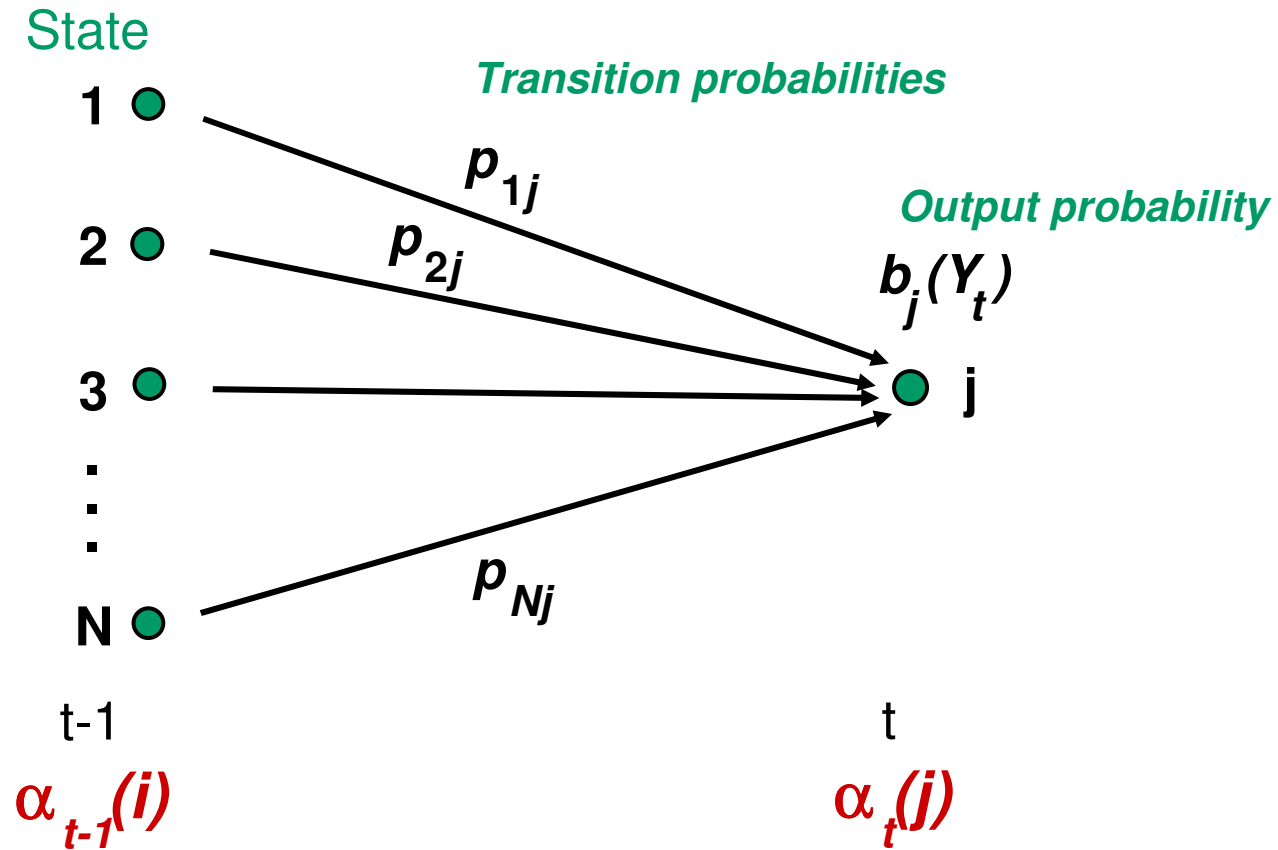
- The probability of the observed data given the model.

The forward algorithm

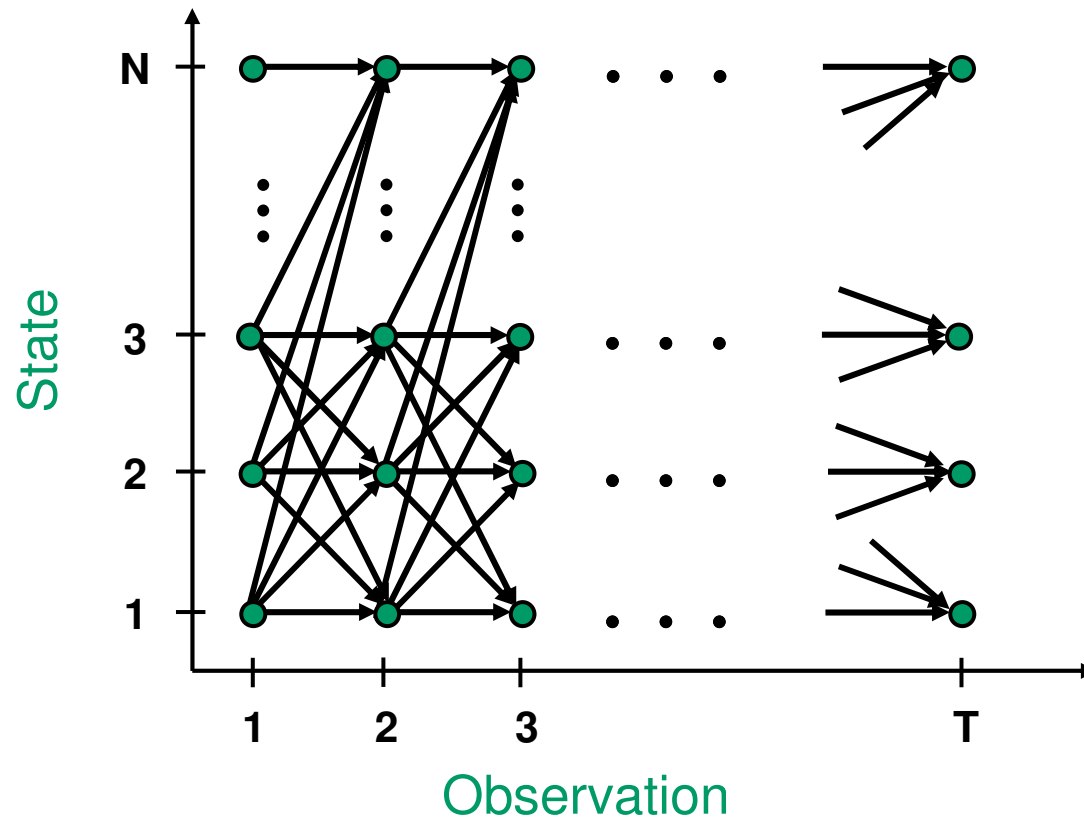
- The best hidden state sequence given the data.

The Viterbi algorithm

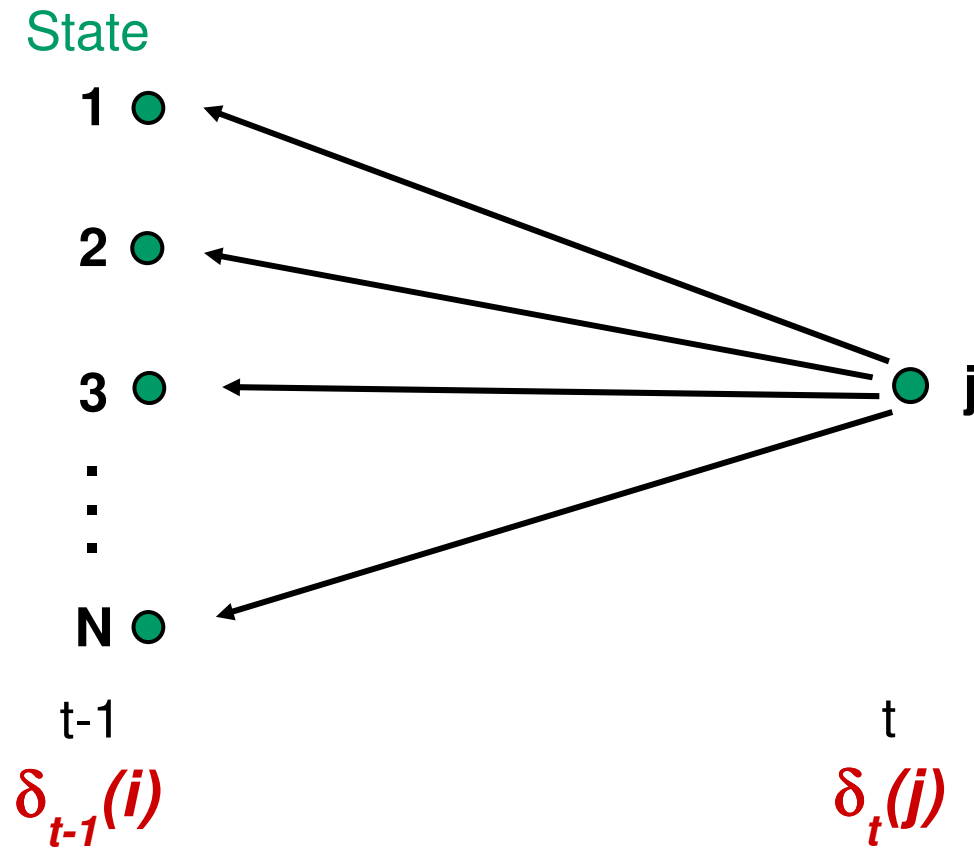
The forward algorithm - α



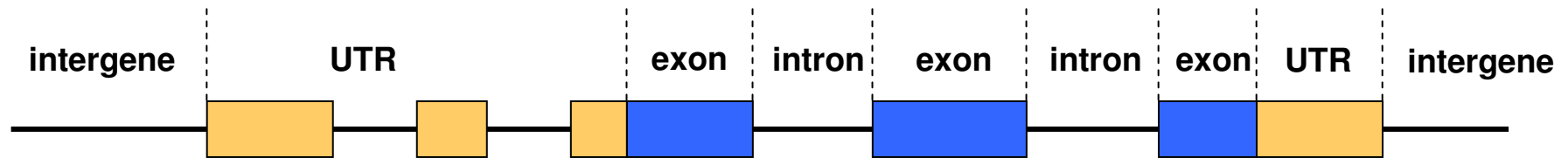
The forward algorithm



The Viterbi algorithm - δ



HMMs for gene finding



Observed:

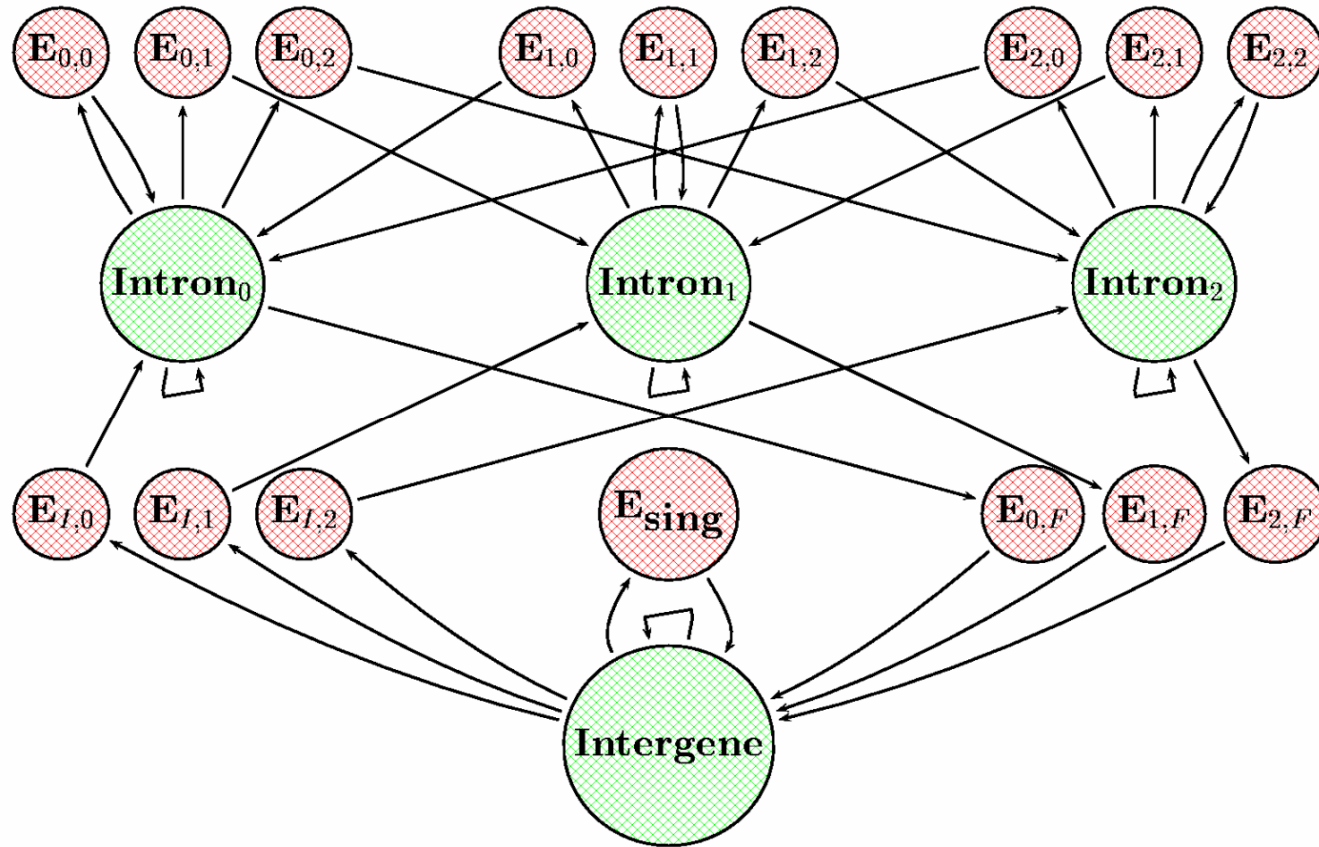
CTTGATGCTGGCACGTTCTGCTTCATCGGAGACAAATTACGGCTTTCCGGAGCA

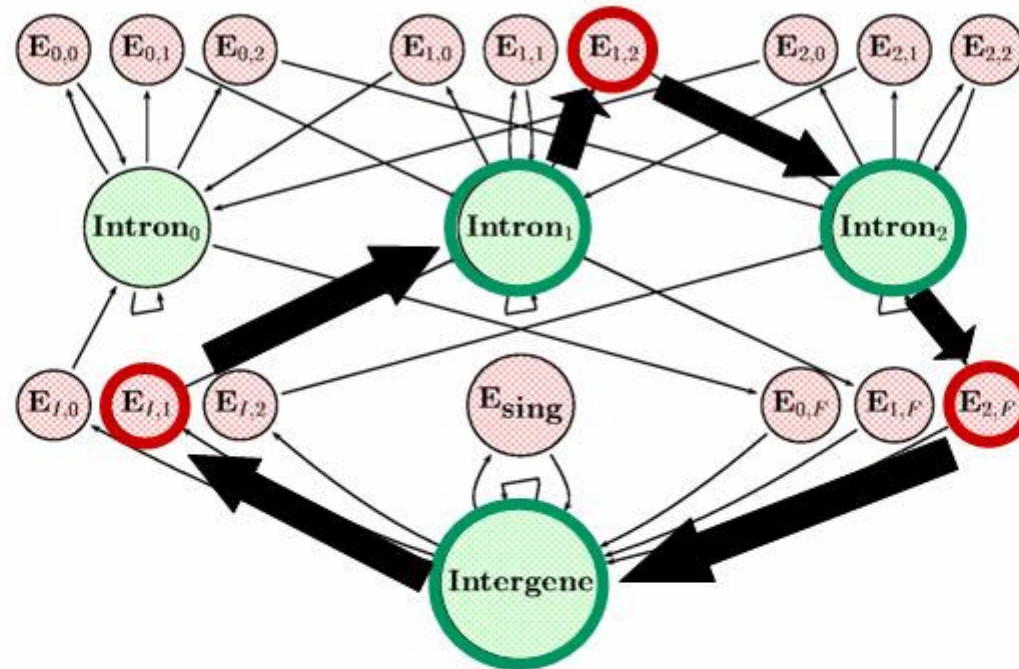
Hidden:

CTTGATGCTGGCACGTTCTGCTTCATCGGAGACAAATTACGGCTTTCCGGAGCA



Generalized HMM





TAATATGTCCACGGGTATTGAGCATTGTACACGGGGTATTGAGCATGTAATGAA

Exon1

Exon2

Exon3



Sequence alignment

- *Pair HMMs*

Sequence alignment

```
50      .      :      .      :      .      :      .      :      .      :
247 GGTGAGGTCGAGGACCCTGCA  CGGAGCTGTATGGAGGGCA  AGAGC
      |:  ||  ||||:  ||||  --:||  |||  |::|  |||---|||
368 GAGTCGGGGGAGGGGGCTGCTGTTGGCTCTGGACAGCTTGCATTGAGAGG

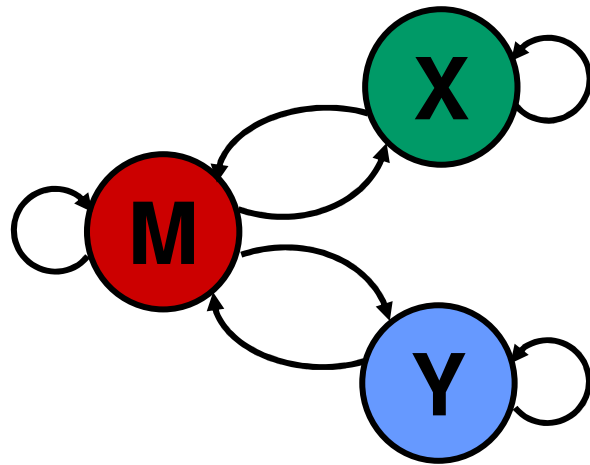
100     .      :      .      :      .      :      .      :      .      :
292 TTC                      CTACAGAAAAGTCCCAGCAAGGAGCCACACTTCACTG
      |||-----|| |  |::|  |:  |||:::|:|:-||  |:|  |
418 TTCTGGCTACGCTCTCCCTTAGGGACTGAGCAGAGGGCT  CAGGTCGCGG

150     .      :      .      :      .      :      .      :      .      :
332                      ATGTCGAGGGGAAGACATCATTCGGGATGTCAGTG
      -----||| ||||| ||||| ||||| ||||| :||| ||||| |||||
467 TGGGAGATGAGGCCAATGTCGAGGGGAAGACATCATTTGGGATGTCAGTG

200     .      :      .      :      .      :      .      :      .      :
367 TTCAACCTCAGCAATGCCATCATGGGCAGCGGCATCCTGGGACTCGCCTA
      |||||:|||||||:||||||| |||||:|  ||:|||||:|||||||
517 TTCAATCTCAGCAACGCCATCATGGGCAGTGAATTCTGGGGCTCGCCTA
```



Pair HMMs for alignment



M = (mis)match

X = insert seq1

Y = insert seq2

Output sequence:

ATCG- - G

AC- GTCA

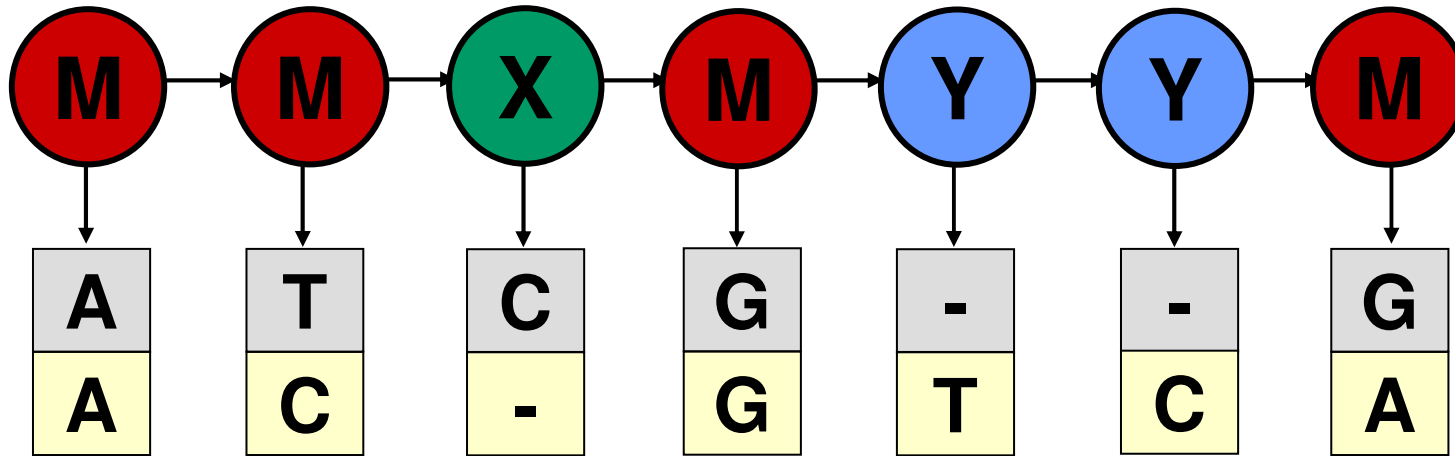
Observed sequences:

ATCGG

ACGTCA



Pair HMMs



Output sequence:

ATCG- - G
AC- GTCA

Observed sequences:

ATCGG
ACGTCA

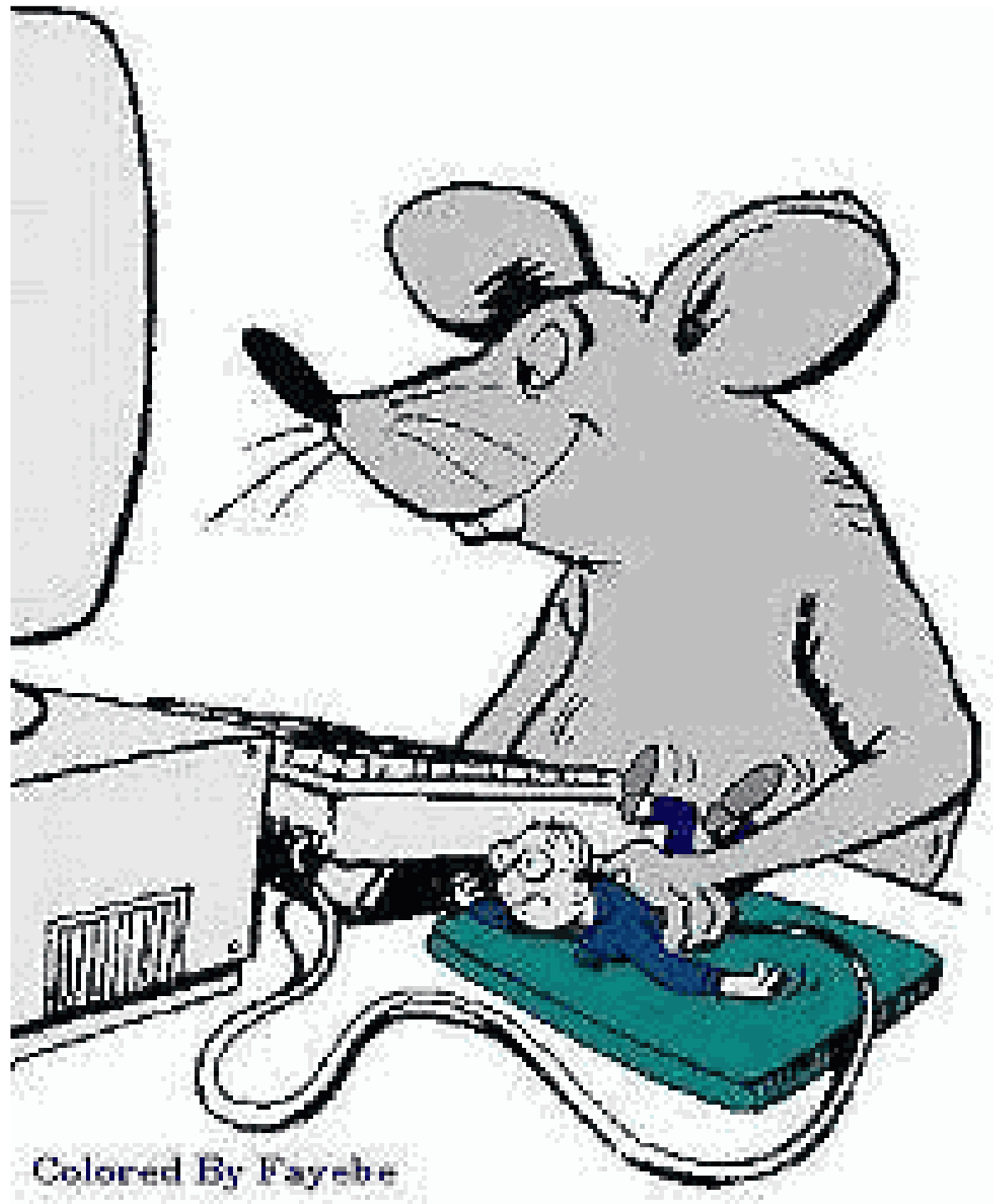
Comparative gene finding

- *Generalized Pair HMMs*

Comparing human and mouse



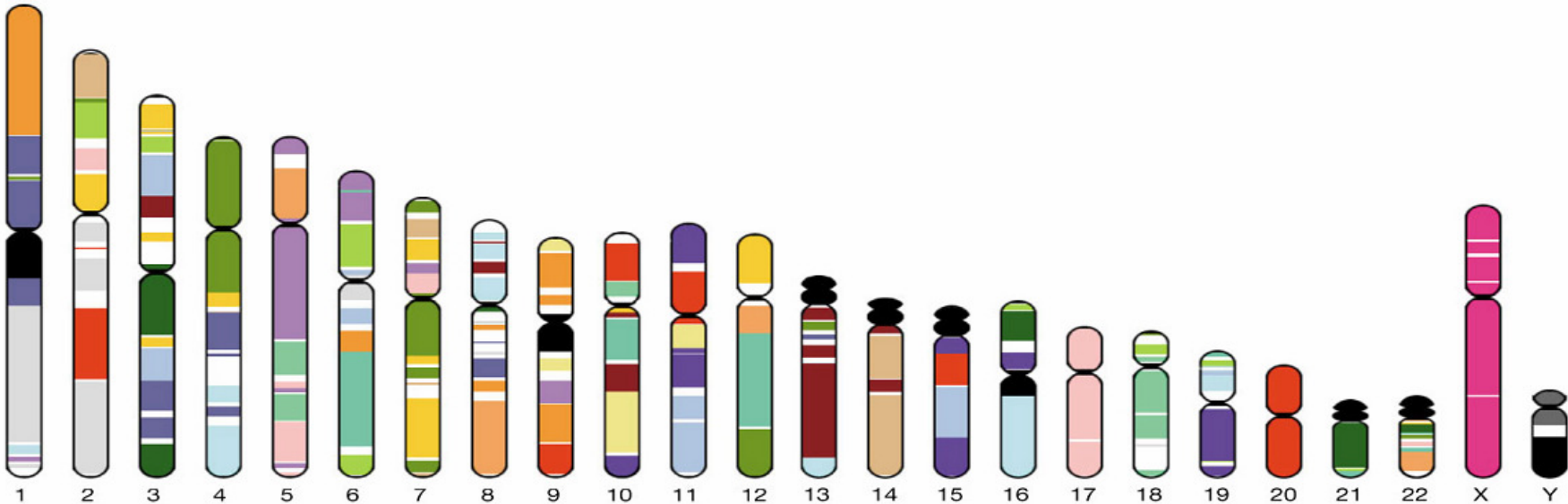
Fraunhofer **CHALMERS**
Research Centre
Industrial Mathematics



Colored By Fayebe

Why mouse?

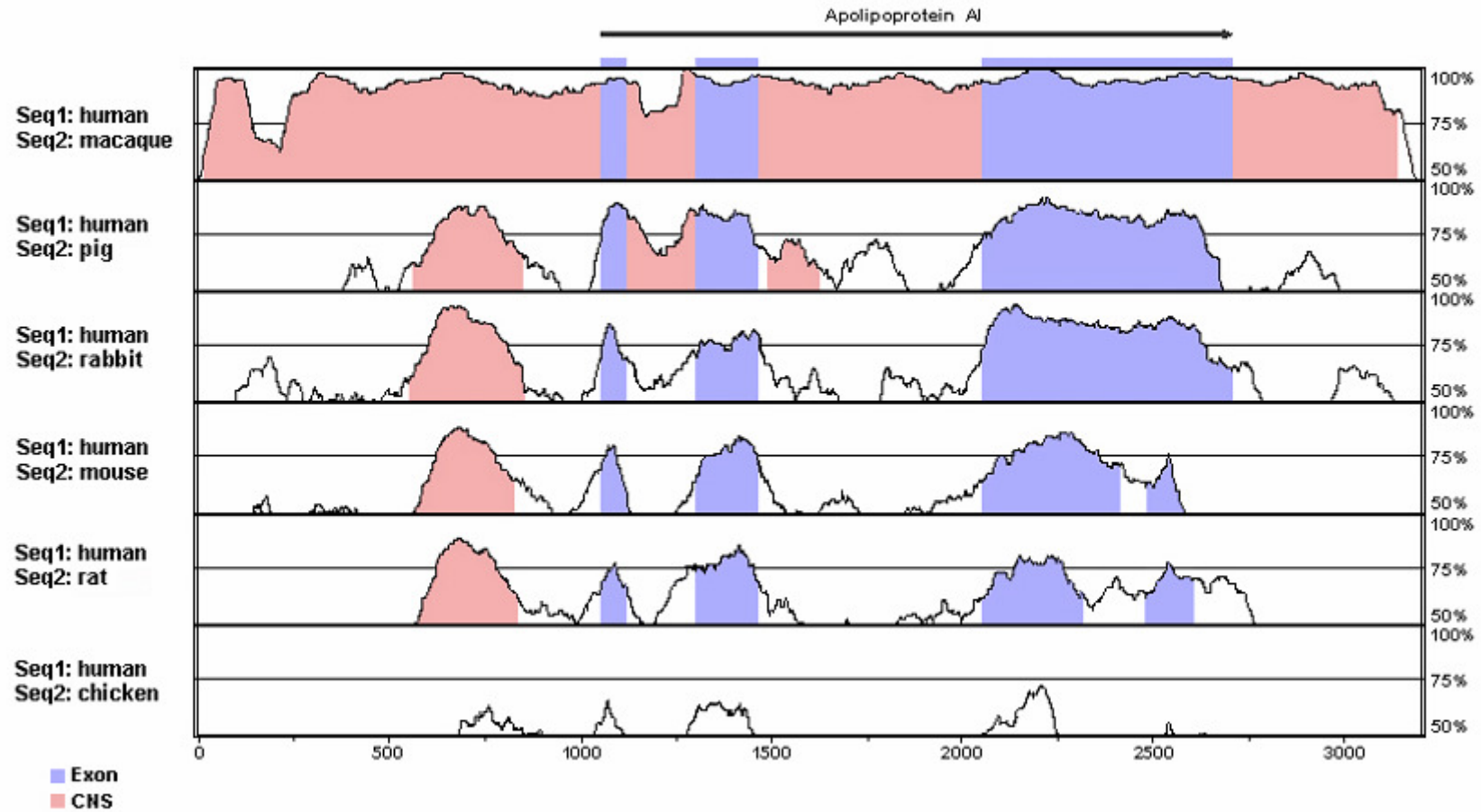
Human

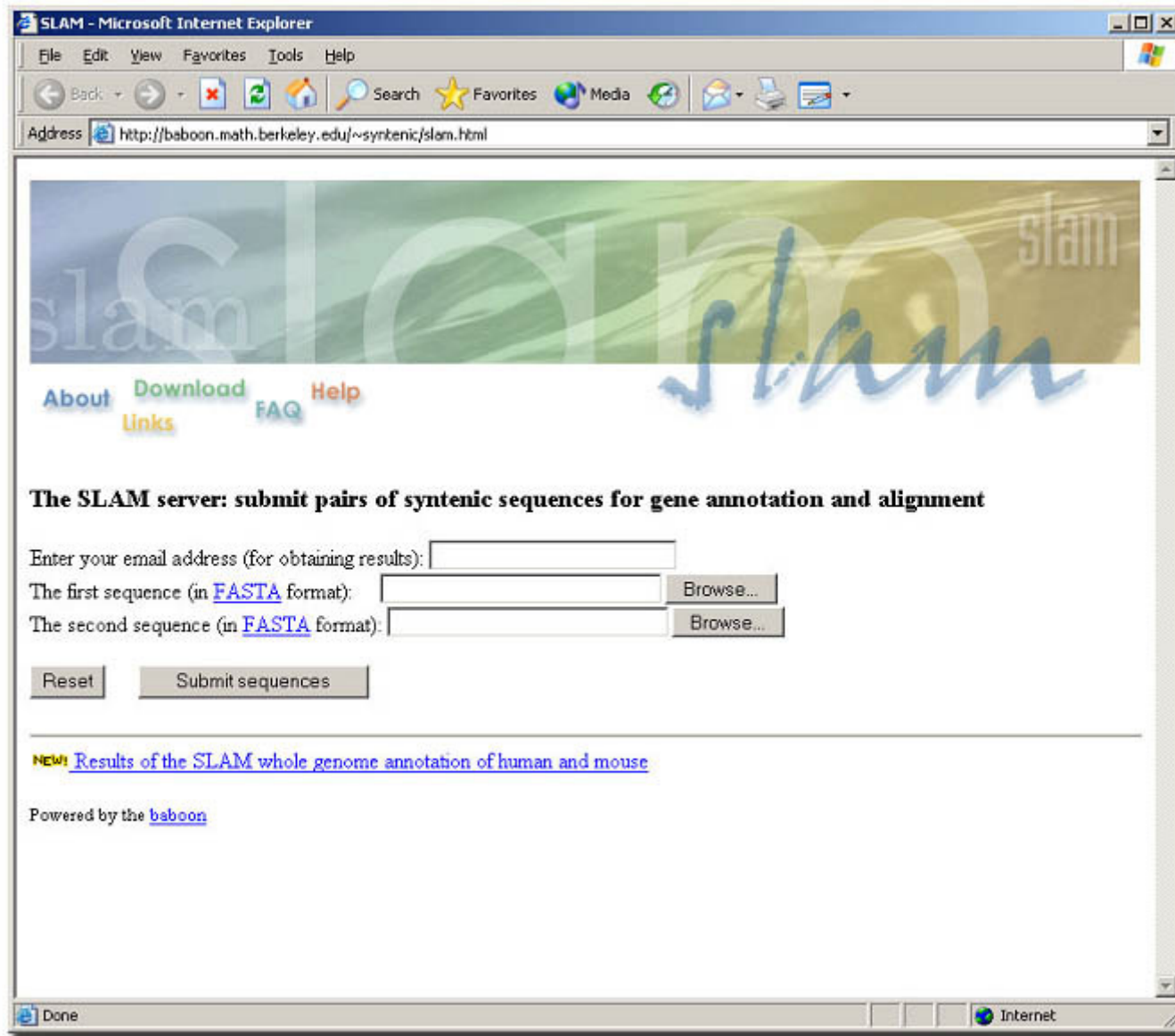


Mouse



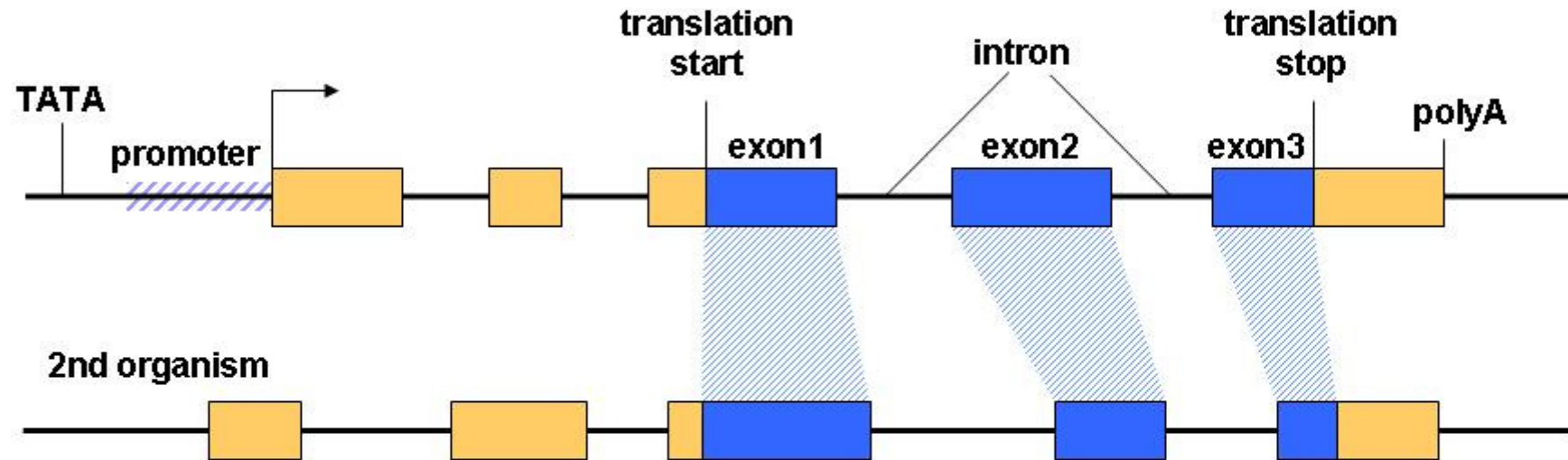
human vs macaque, pig, rabbit, mouse, rat, chicken

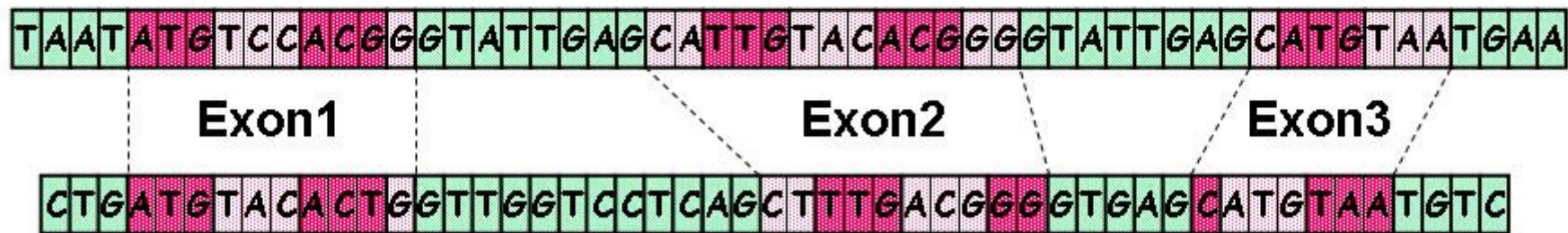
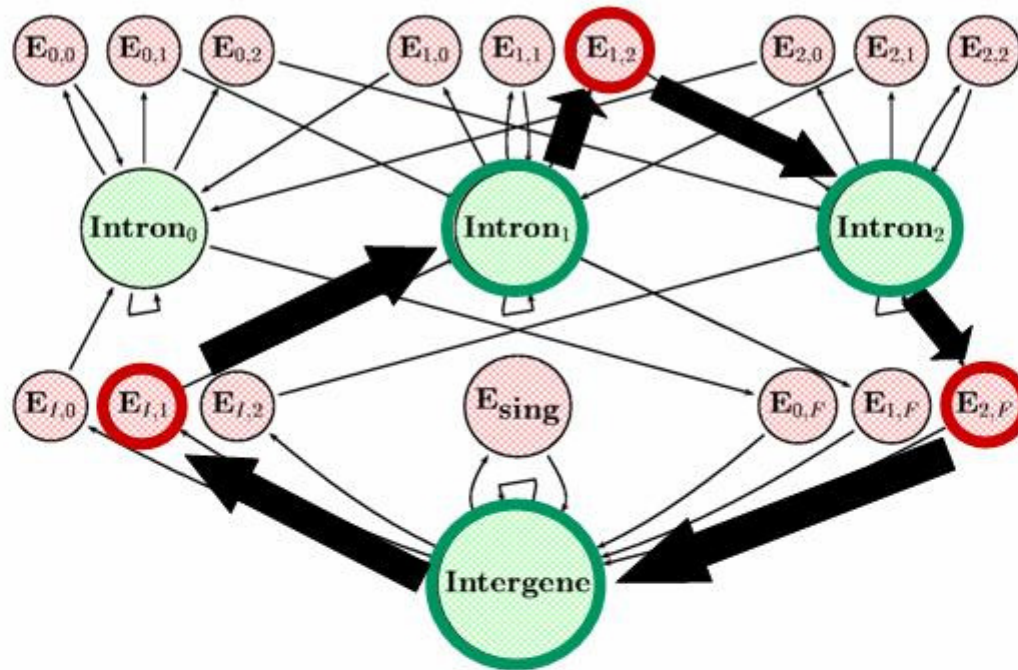




<http://bio.math.berkeley.edu/slam/>

Comparative gene finding in SLAM





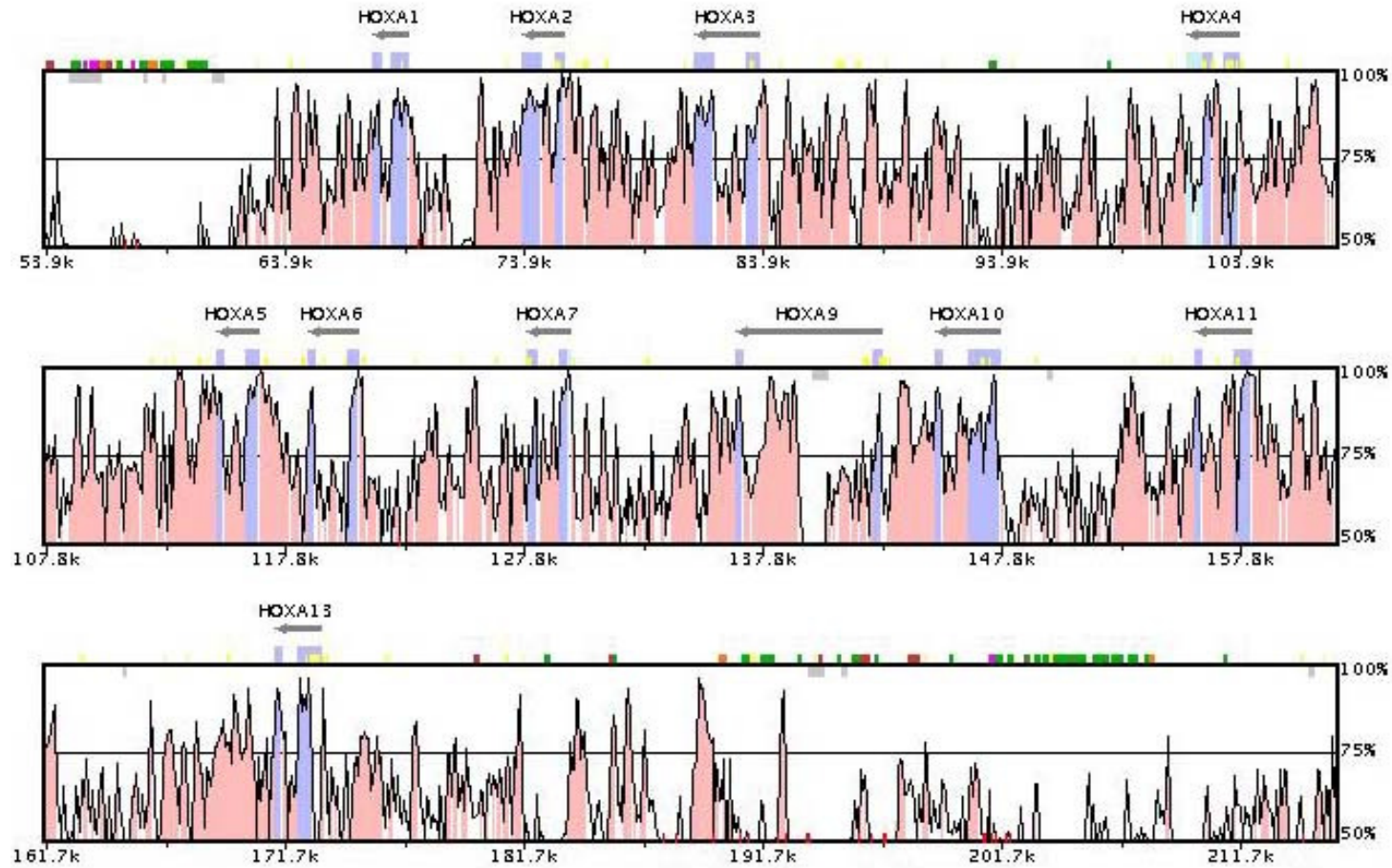
The Rosetta set

Local	SN	SP	AC
Genscan	0.96	0.91	0.93
Twinscan	0.96	0.94	0.94
SGP-1	0.94	0.96	0.94
SLAM	0.95	0.98	0.96

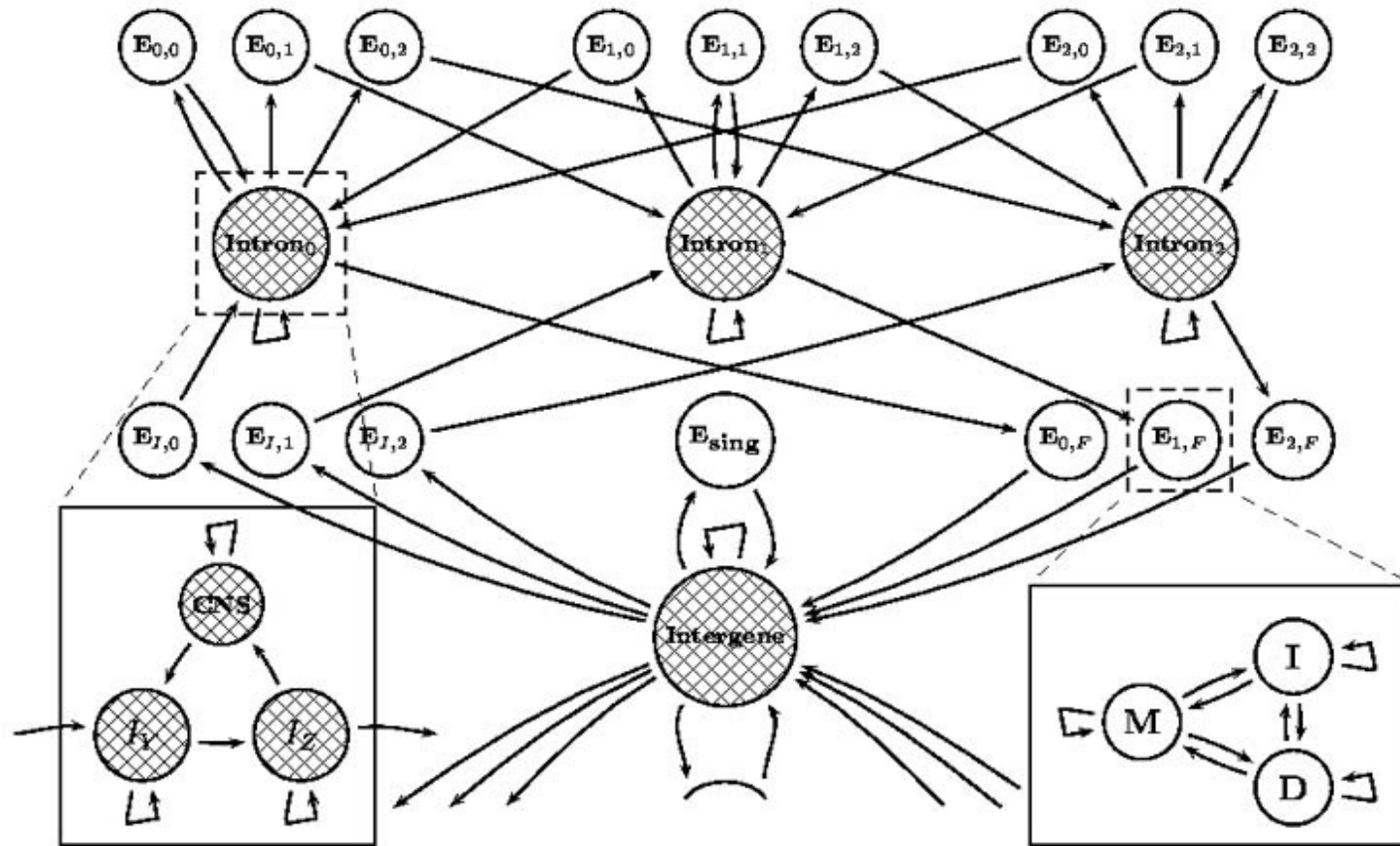
Global	SN	SP	AC
Genscan	0.87	0.87	0.84
Twinscan	0.89	0.92	0.89
SGP-2	0.88	0.94	0.89
SLAM	0.87	0.96	0.90



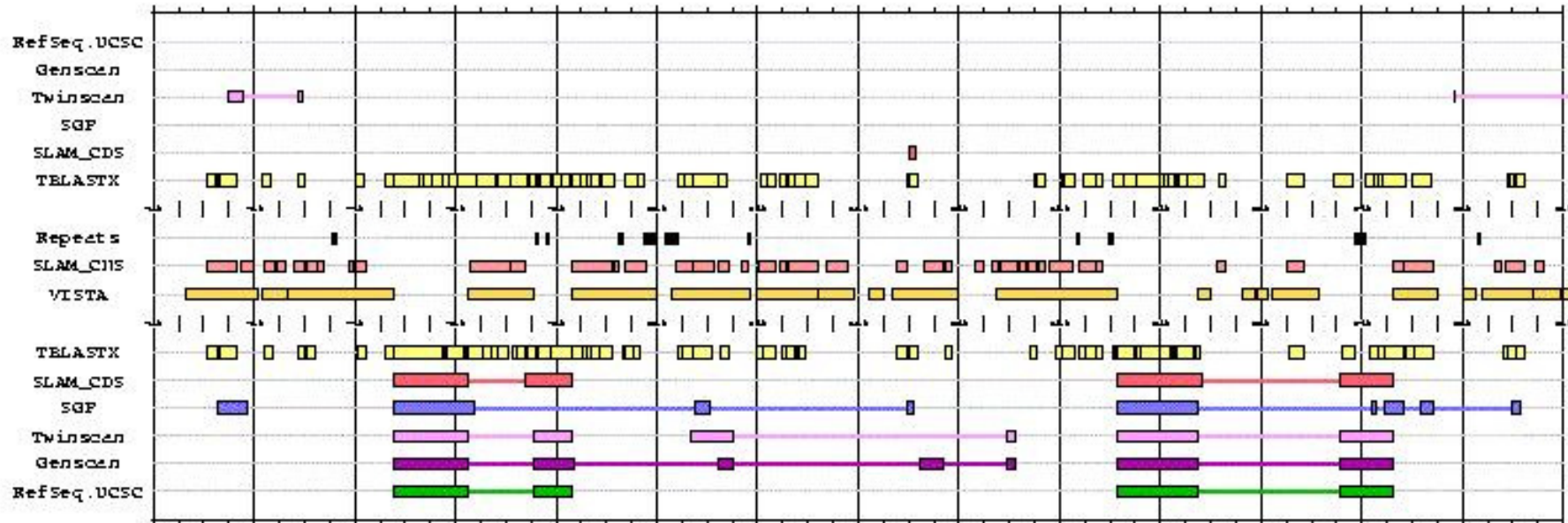
A problem: HoxA human-mouse



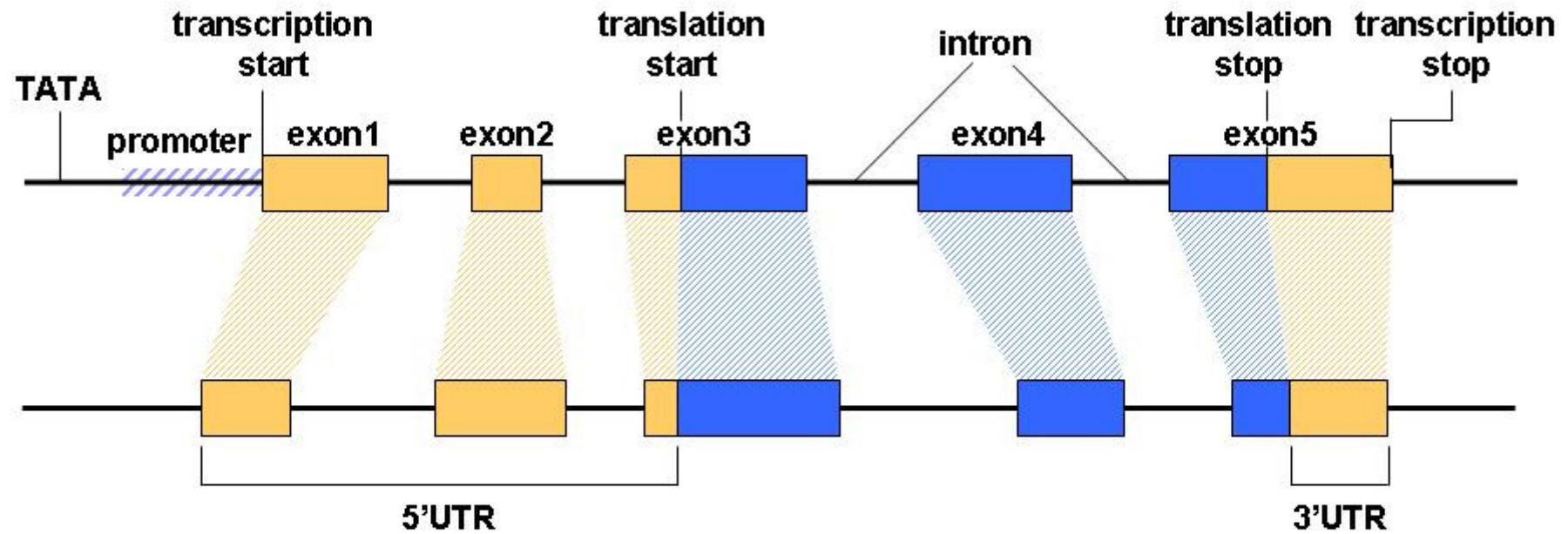
The SLAM model



Example: HoxA2 and HoxA3



Comparative UTR prediction



Reducing computational complexity

Computational complexity

Model	Time	Space
HMM	N^2T	NT
PHMM	N^2TU	NTU
GHMM	D^2N^2T	NT
GPHMM	D^4N^2TU	NTU

N = no. states

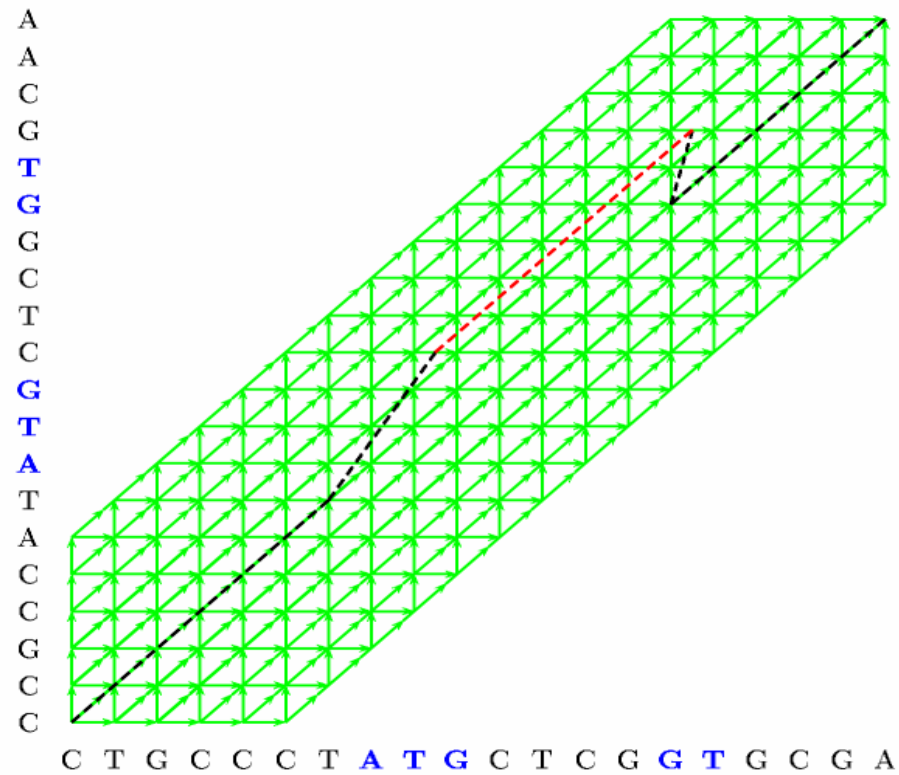
T = length seq1

D = max duration

U = length seq2

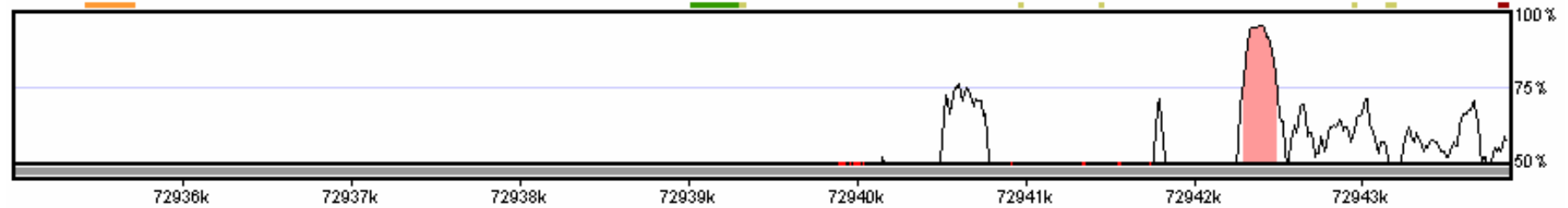
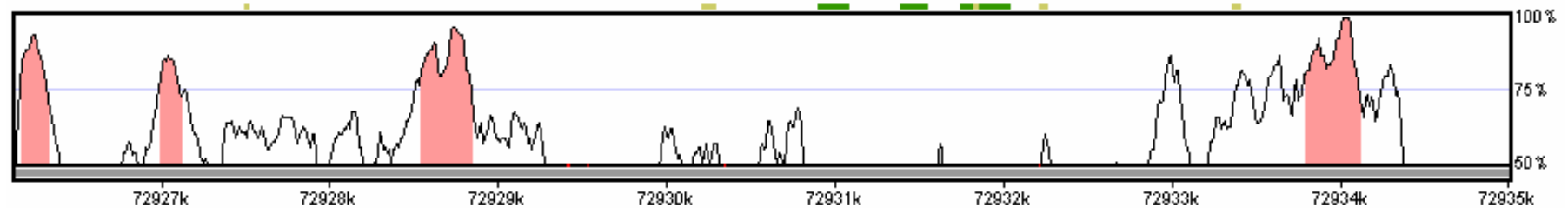
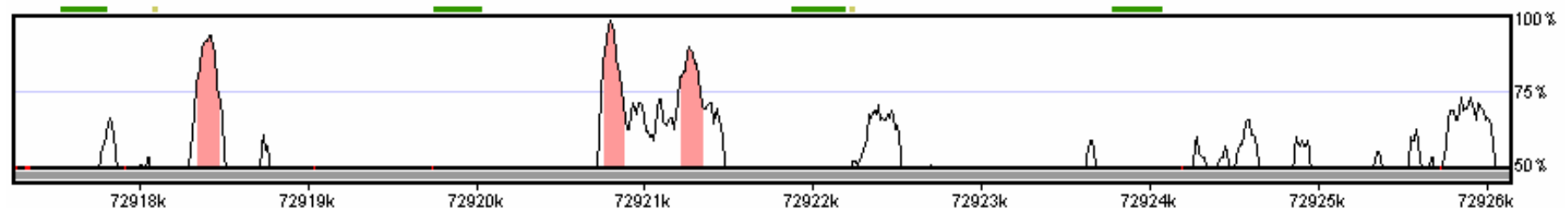
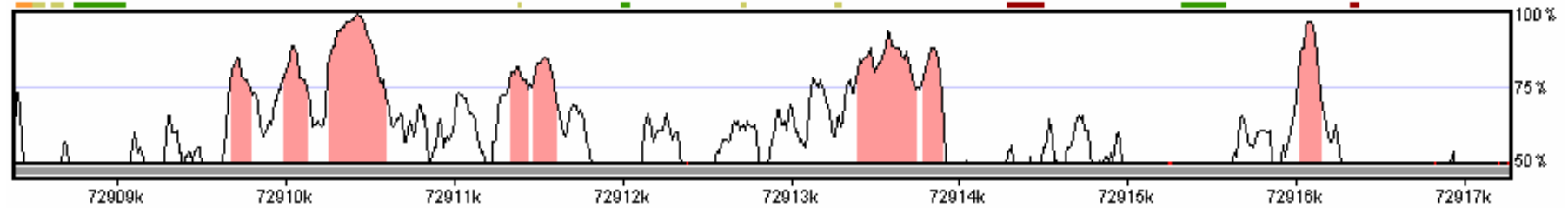


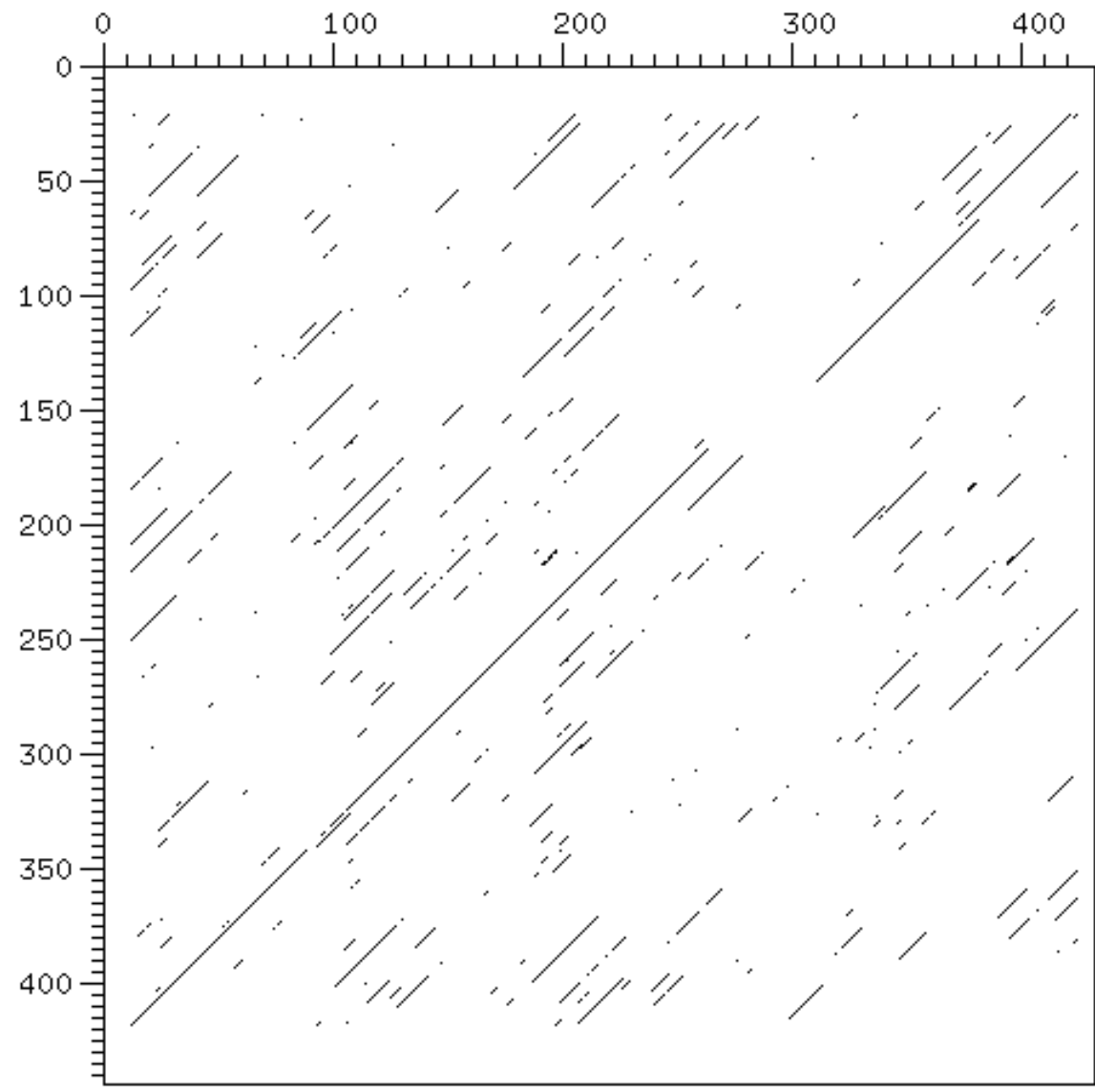
Approximate alignment

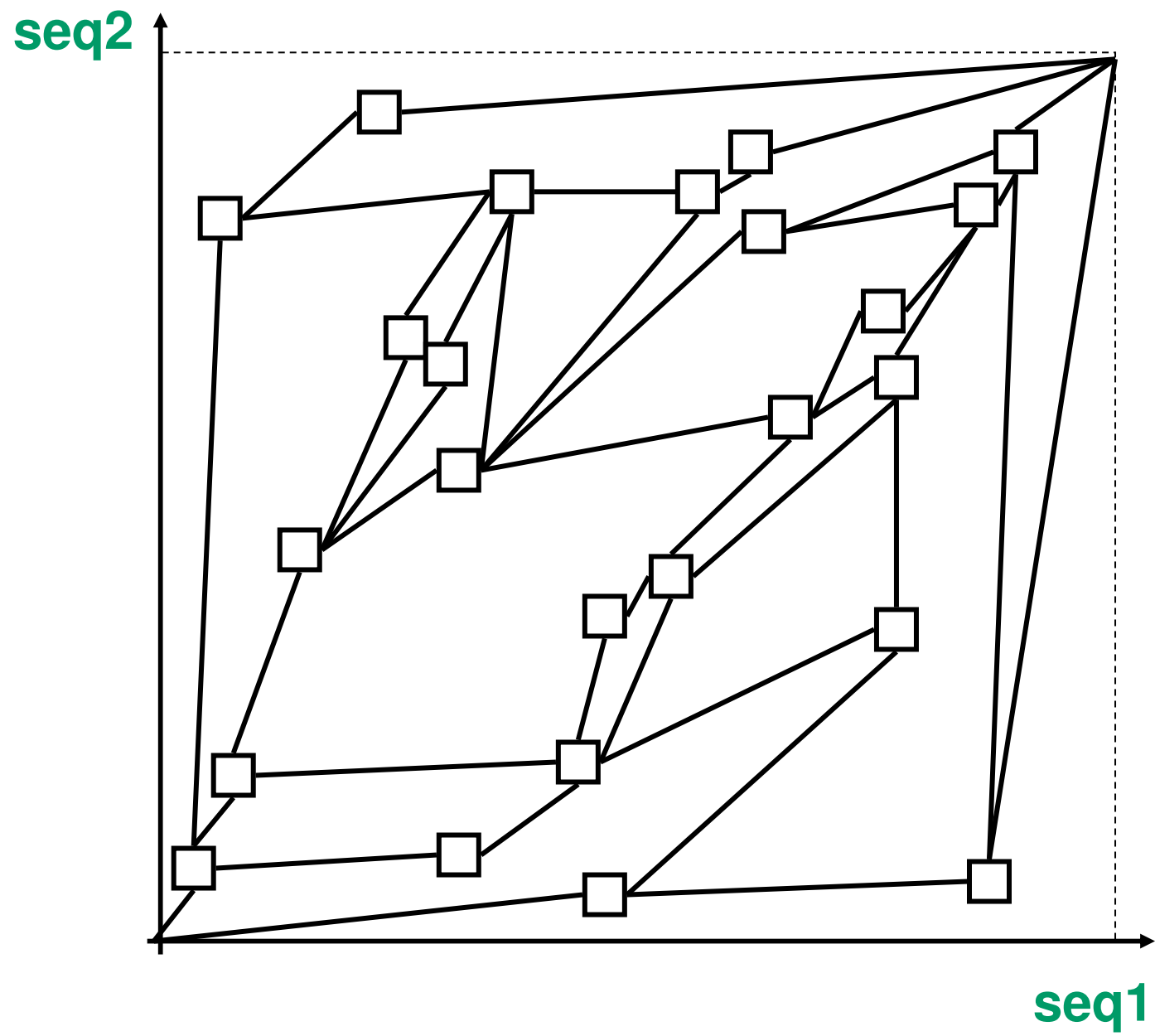


Reduces
TU-factor to *hT*

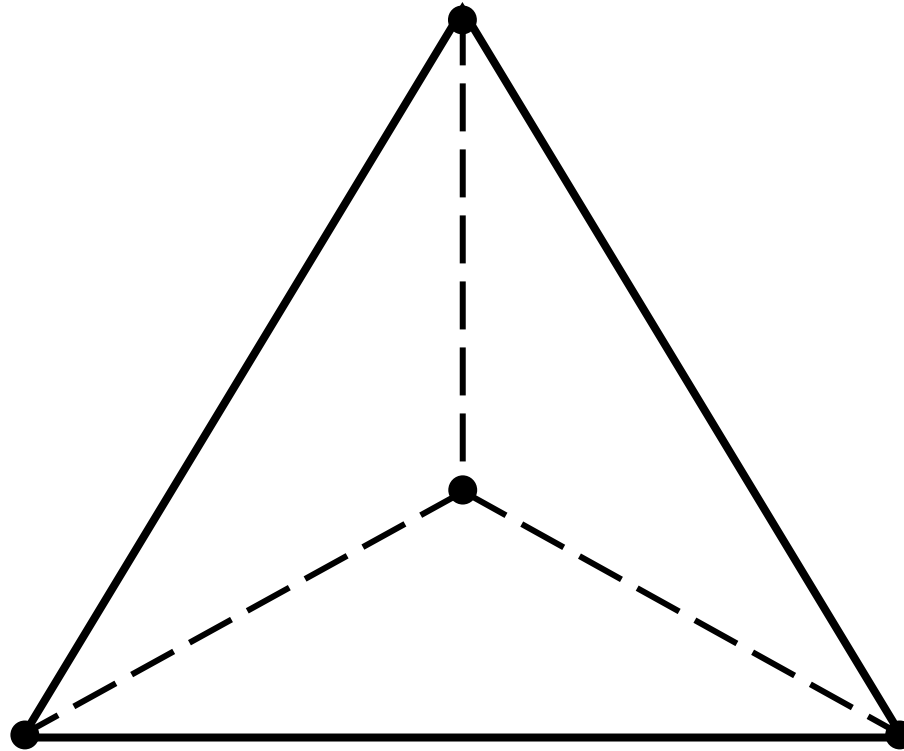




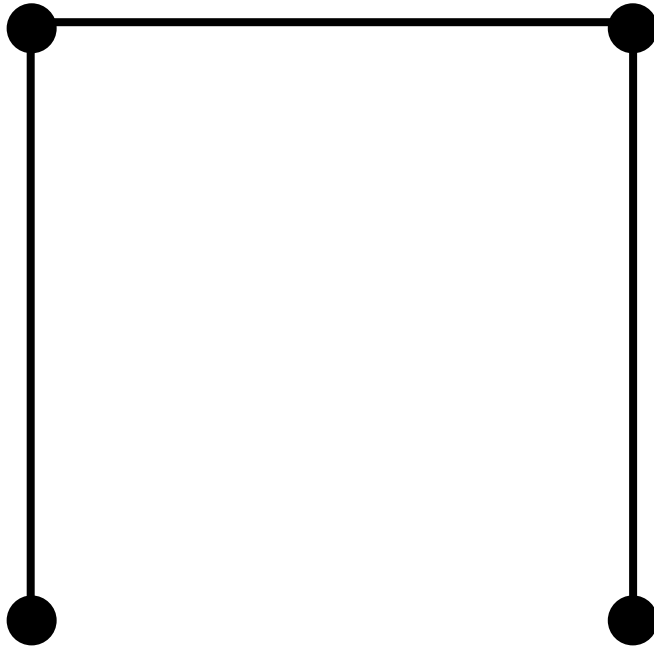




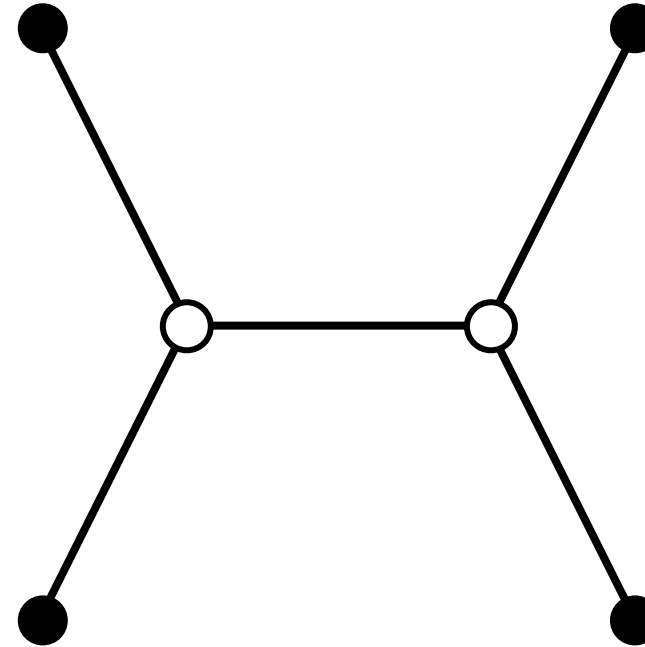
Steiner trees



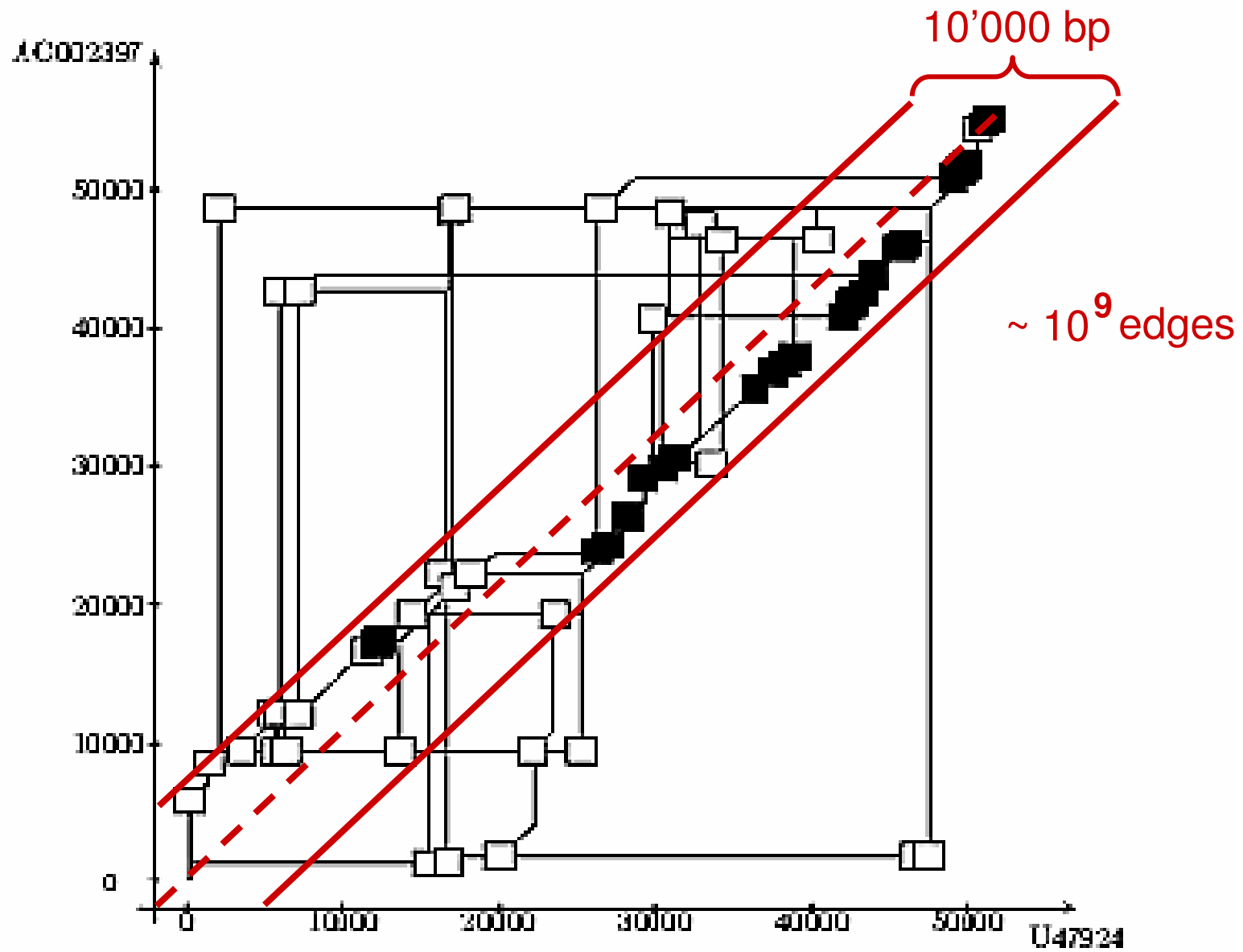
Steiner trees



Minimum spanning
tree



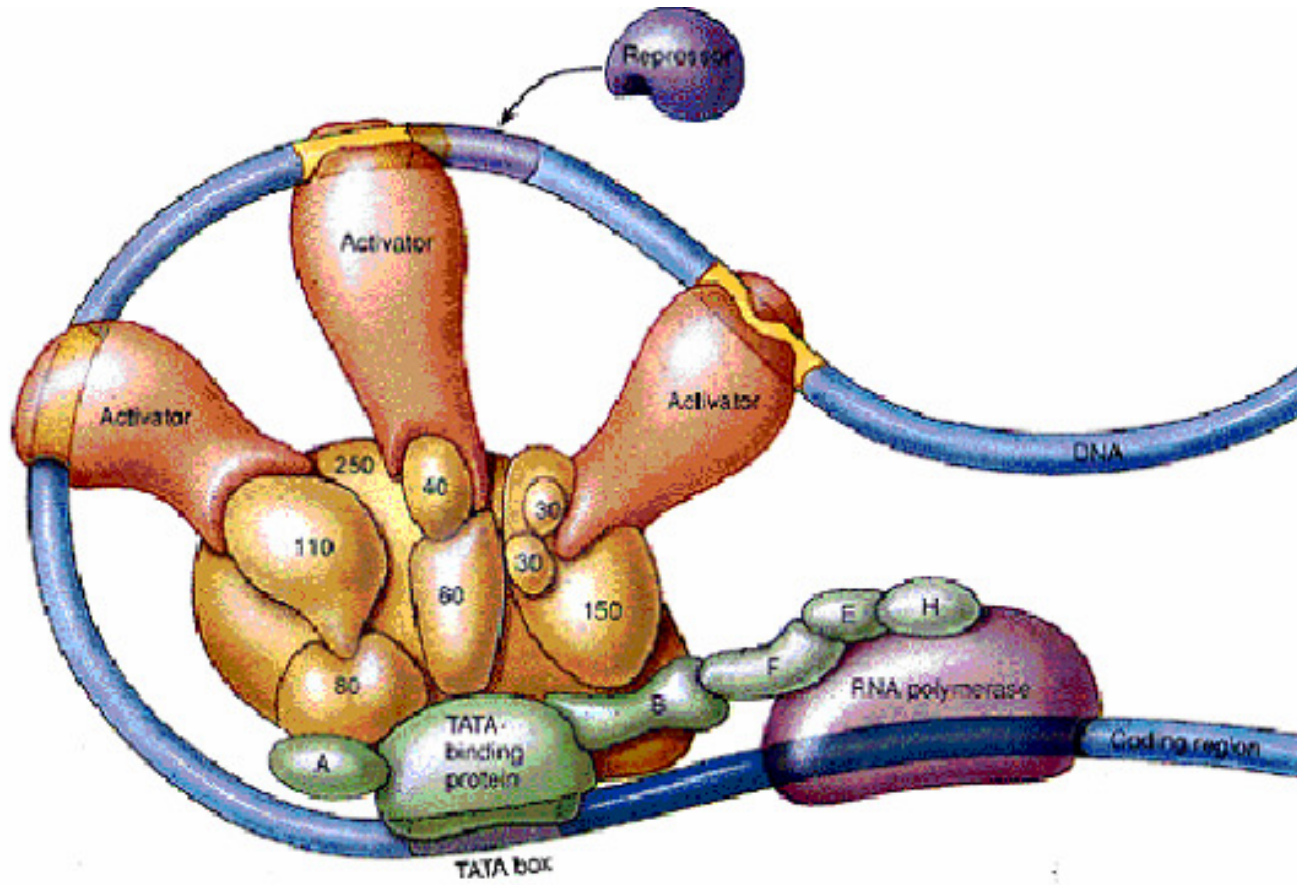
Steiner tree



2-approximation in $O(n^3)$

Regulatory regions

Regulatory complex



a

	Hepatic site C	CCAAT box
Mouse	NNNNAGCCTCAGGAACAGAGCTGATCCTTGA	AAGTTCCACATCGCCAGCAAAG
Rabbit	NNNN-GCCCTAGGGACGGAGCTGATCCTTGA	AAGTTCCACATGGCCAGGACCAG
Human	NNNNAGTCCCAGGGACAGAGCTGATCCTTGA	AAGTTCCACATTGCCAGGACCAG

Mouse	TAAGCAGTGGCAGGGCCAG-GCTGAGCTTATCAGTCTCCAGCCCAGCCCCTGCCACAC
Rabbit	GGAGCAGTGA
Human	TGAGCAGCAACAGGGCCAGGGCTGGGCTTATCAGCCTCCAGCCCAGACCCTGGCTGCAG

TATA box

Mouse	ACATATATAGACCAGGGAAGAAGAGCTGGACACCC-
Rabbit	ACATAAATAGGCCAGGGGCCA---GCTGGCCGCAGG
Human	ACATAAATAGGCCCTGCAAGA---GCTGGCTGC---

b

	Hepatic site C	CCAAT box
Mouse	AGCCTCAGGAACA-GAGC-TGATCCTTGA	AAGTTCCACATCGCCAGCAAAGTA
Rabbit	-GCCCTAGGGACG-GAGC-TGATCCTTGA	AAGTTCCACATGGCCAGGACCAGGG
Human	AGTCCCAGGGACA-GAGC-TGATCCTTGA	AAGTTCCACATTGCCAGGACCAGTG
Chicken	CTCTCCCGGGCCGTGCGCACGATCCTTGA	ACGCGCCACATCGCCCGCGCCGGGA

Mouse	AGCAGTGGCAGGGC--CAG-GCTGAGCTTATCAGTCTCCAGCCCAGCCCCTGCCACAC
Rabbit	AGCAGTGA
Human	AGCAGCAACAGGGC--CAGGGCTGGGCTTATCAGCCTCCAGCCCAGACCCTGGCTGCAG
Chicken	GTGATTTCTTGGGCTGCGGCCTG-GCTTATCTGGTGC

TATA box

Mouse	ACATATATAGACCAGGGAAGAAGAGCTGGACACCC-
Rabbit	ACATAAATAGGCCAGGGGCCA---GCTGGCCGCAGG
Human	ACATAAATAGGCCCTGCAAGA---GCTGGCTGC---
Chicken	-CATAAATAGCGGCGGCGGGA---ACCGGGCTCAC-

How many genes do we have???

27'462

Acknowledgements

- Lior Pachter, UC Berkeley
- Simon Cawley, Affymetrix