# Hidden Markov Models and Bioinformatics
## - Docentföreläsning

Marina Alexandersson

Göteborg, Dec 9, 2004

# Some genetic history

1850　　　　　　　　　1900　　　　　　1950　　　　　　2000

**1859**



Charles Darwin

1859: Origin of the Species

# Some genetic history

1850　　　　　　　　1900　　　　　　1950　　　　　2000

**1865**



Gregor Mendel

1865: Gregor Mendel's peas

# Some genetic history

1850　　　　　　　　　1900　　　　　　　1950　　　　　　　2000

**1869**



Courtesy of Herm Courvoisier, Portrait-Sammlung, University of Basel.
Noncommercial, educational use only.

Friedrich Miescher

1869: DNA isolated

# Some genetic history

1850           1900           1950           2000

**1902**

Walter Sutton      Theodor Boveri

1902: Chromosome theory of heredity

# Some genetic history

1850     1900     1950     2000

**1909**

1909: the word 'gene' introduced

Willhem Johannsen

# Some genetic history

1850          1900          1950          2000

**1941**



George Beadle



Edward Tatum

1941: One gene, one enzyme

# Some genetic history

1850           1900           1950           2000

**1953**



1953: The DNA double helix

James Watson           Francis Crick

# Some genetic history

1900 1950 2000

**1955**



1955: 46 human chromosomes

ITWM

Fraunhofer **CHALMERS**
Research Centre
Industrial Mathematics

# Some genetic history

1850             1900             1950         2000

**1966**

C
C
A
T
G
A
C
A
T
C
T
G
A
A
C
A
G
A

M → Methionine

T → Threonine

S → Serine

E → Glutamic acid

Q → Glutamine

1966: Genetic code cracked

ITWM

Fraunhofer

CHALMERS
Research Centre
Industrial Mathematics

# Some genetic history

1850            1900            1950            2000

**1977**



Walter Gilbert      Frederick Sanger

1977: DNA sequencing

# Some genetic history

1850          1900          1950          2000

**1990**



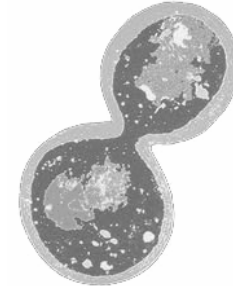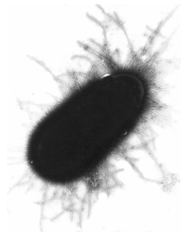1990: Human Genome Project

# Sequencing history

**1995: Two microbial genomes (1.8, 0.6)**

**1996: Saccharomyces cerevisiae (12)**

**1997: E. coli (4.6)**

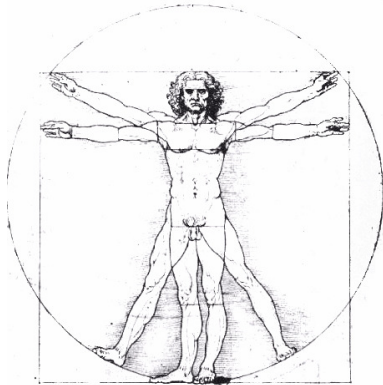**1998: C. elegans (97)**

**2000: Drosophila melanogaster (180)**

**2003: Homo sapiens (3,200)**

Fraunhofer

**CHALMERS**
Research Centre
Industrial Mathematics

**More to come...**

Fraunhofer ITWM

CHALMERS
Research Centre
Industrial Mathematics

# Whole genome analysis

- Gene finding
- Sequence alignment
- Regulatory region discovery

**ITWM**

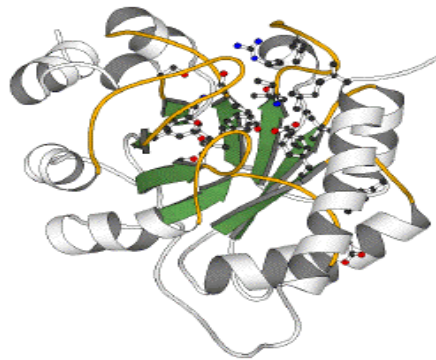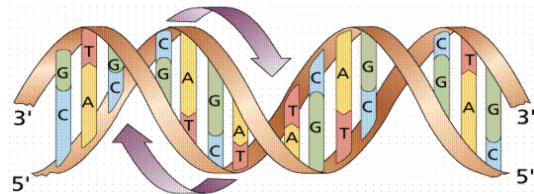Fraunhofer **CHALMERS**
Research Centre
Industrial Mathematics

# 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 methabolic pathways
- Trace disease genes
- ...

Fraunhofer

**CHALMERS**
Research Centre
Industrial Mathematics

# Gene finding

# Gene expression



**DNA**

↓ *transcription*

**RNA**

↓ *translation*
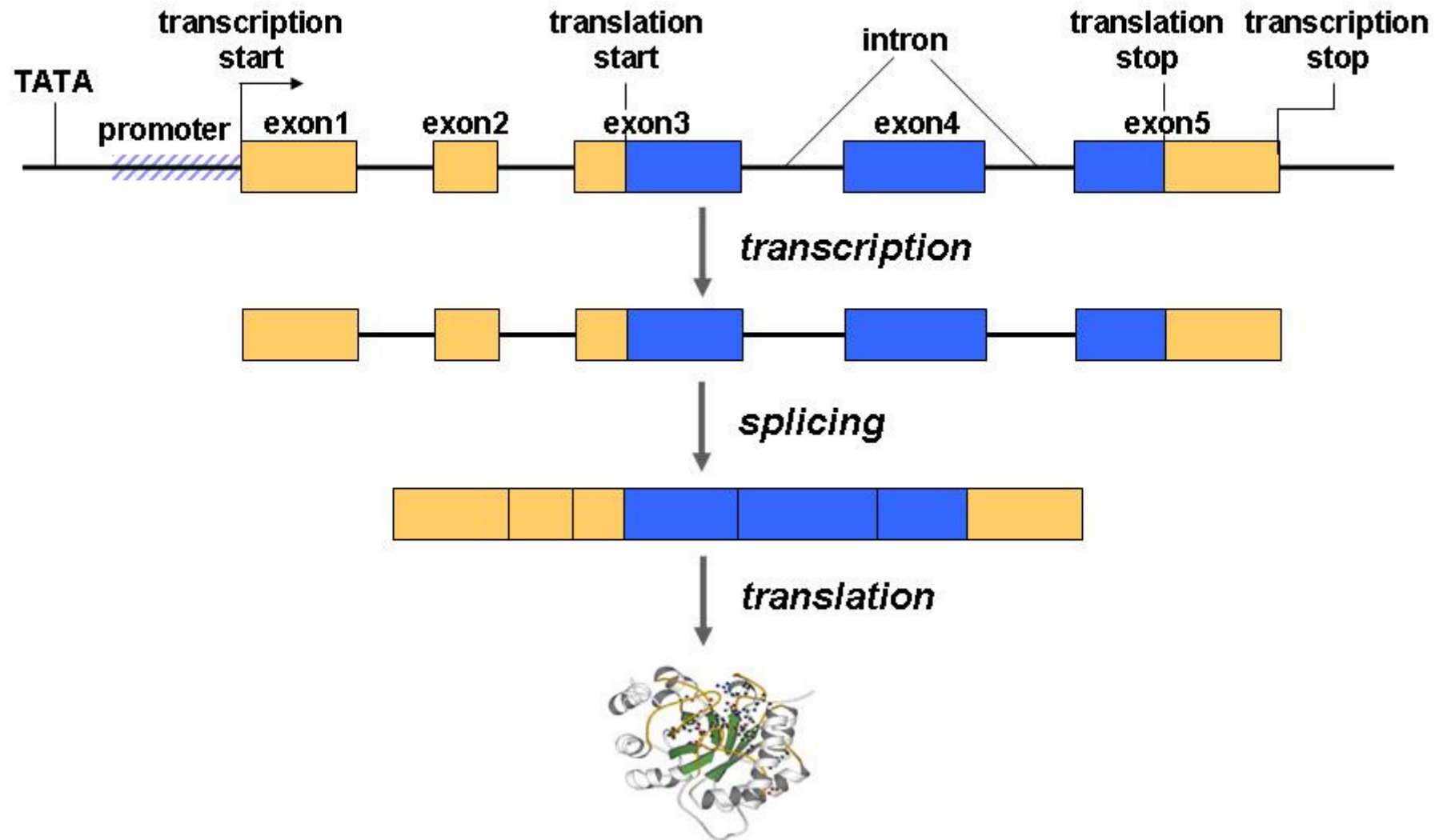
**Protein**

CCTGAGCCAACTATTGATGAA
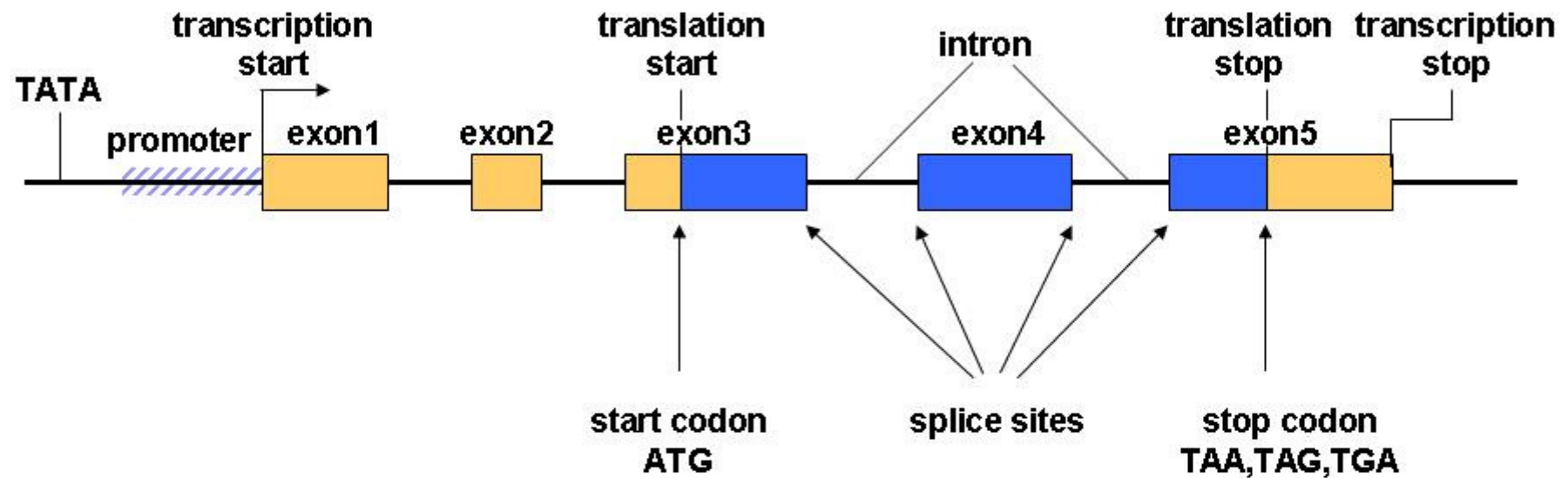
↓

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

↓

P**E**PT**I**DE

# Gene structure

# Finding genes

> HSCKIIBE, Human gene for casein kinase II subunit beta (EC 2.7.1.37).
ggggctgagatgtaaattagaggagctggagaggagtgcttcagagtttgggttgctttaagaaagggt
ggttccgaattctcccgtggttggagggccgaatgtgggaggagggaggataccagaggcagggaagga
gaacttgagctttactgacactgttctttttctagctgacgtgaagatgagcagctcagaggaggtgtc
ctggatttcctggttctgtgggctccgtggcaatgaattcttctgtgaagtgagttctcttcaacctcc
ctacttgccagcttcacatatcttcccaccagacgttccttcacatattccacttctacactgttctct
aaagcttttatgggagagagtgtaggtgaactagggagagacacaagtacttctgctgagttgggagtg
agaaacaagcacaacagatgcagttgtgttgatgataaggcatcacttagagcattttgcccaggtcaa
agatgaggatttgatatgggttccctcttggcttccatgtcctgacaggtggatgaagactacatcca
ggacaaatttaatcttactggactcaatgagcaggtccctcactatcgacaagctctagacatgatctt
ggacctggagcctggtgaggcaccctcagggttgttttgtgtgtgtgcgtgcactattttctcttcaa
atctctattcacttgcctgaatttttgccaaatttcctttggttctctgatttctttaaccccaaattca
tgctttattttgatcctccacctgactcttgtctagttttgtgacgtatatcacttgttctcatgtttt
tgcaagggtcagaagcccaggtttctgggtcccatgcccagatgttggatggggtaaggcccaaaagta
ggtgctaggcaaactgaatagcccgcagcccctggatatgggcagggcacctaggaaagctgaaaaaca
agtagttgcatttggccgggctgtggttcagatgaagaactggaagacaaccccaaccagagtgacctg
attgagcaggcagccgagatgctttatggattgatccacgcccgctacatccttaccaaccgtggcatc
gcccagatggtgaggcctctctgctcctacctgcctccttctgagcagtaagagacacaggttcctgca
gcaagaagtcatgtttaagccctgtttaaggaagctagctgagaagaggggaagaaccccagaacttgg
ccctgccctaatttggaagaaaggcaacacagaagtttgagagcccatctagtccagagaagggggcct
ctggacagagttggaaggagtgccgacagagttggtatgggttgggctgcgaagggagttgcctcttct
ttacatctacctgccaacccccttccattgtattcacctcagttggaaaagtaccagcaaggagactttg
gttactgtcctcgtgtgtactgtgagaaccagccaatgcttcccattGgtgagtgttgaagaagggaaa
ggaaagcaccgtgtggcagtcttatgggaaggagttggggctcaacacattggagcctgagtcctgagg
ggaggttaggtaggaataggggggatacctggcctgctgagtctggctgtctcccaggcctttcagacat
cccaggtgaagccatggtgaagctctactgccccaagtgcatggatgtgtacacacccaagtcatcaag
acaccatcacacggatggcgcctacttcggcactggtttccctcacatgctcttcatggtgcatcccga
gtaccggcccaagagacctgccaaccagtttgtgcccaggtagggagcagggagagtcattaagggtca
aaggaaaggcccaagatcccccagagagggggaggacagggcatggccctttcttgaggtctgcttctcc
cagaatcagggcatctccctgctgagtgactgtgggaaagttatttgattatctgtgcttgagttacct
tattgtagaatgttcttgagctgagaagttgggaaccacgaggctttagctctgagcaggtccatagag
gagctcaggtggggaggtgggaatgcaggtgactggcagggcctggatggggctcatgctgctgcctct
ctgacctctgccctggcctaggctctacggtttcaagatccatccgatggcctaccagctgcagctcca
agccgccagcaacttcaagagcccagtcaagacgattcgctgattccctcccccacctgtcctgcagtc
tttgtcttttcctttctttttttgccacccttttcaggaaccctgtatggttttttagtttaaattaaagga
gtcgttatcgtggtgggaatatgaaataaagtagaagaaaaggccatgagctagtctgctggtgcttgc
ggaaggggggtggagcgtggccatggaaatcgggctccacggcccagggatgg

> HSCKIIBE, Human gene for casein kinase II subunit beta (EC 2.7.1.37).
ggggctgagatgtaaattagaggagctggagaggagtgcttcagagtttgggttgctttaagaaagggt
ggttccgaattctcccgtggttggagggccgaatgtgggaggagggaggataccagaggcagggaagga
gaacttgagctttactgacactgttctttttctagctgacgtgaagatgagcagctcagaggaggtgtc
ctggatttcctggttctgtgggctccgtggcaatgaattcttctgtgaagtgagttctcttcaacctcc
ctacttgccagcttcacatatcttcccaccagacgttccttcacatattccacttctacactgttctct
aaagcttttatgggagagagtgtaggtgaactagggagagacacaagtacttctgctgagttgggagtg
agaaacaagcacaacagatgcagttgtgttgatgataaggcatcacttagagcattttgcccaggtcaa
agatgaggattttgatatgggttccctcttggcttccatgtcctgacaggtggatgaagactacatcca
ggacaaatttaatcttactggactcaatgagcaggtccctcactatcgacaagctctagacatgatctt
ggacctggagcctggtgaggcaccctcagggttgttttgtgtgtgtgcgtgcactatttttctcttcaa
atctctattcacttgcctgaattttgccaaatttcctttggttctctgatttctttaaccccaaattca
tgctttattttgatcctccacctgactcttgtctagttttgtgacgtatatcacttgttctcatgtttt
tgcaagggtcagaagcccaggtttctgggtcccatgcccagatgttggatggggtaaggcccaaaagta
ggtgctaggcaaactgaatagcccgcagcccctggatatgggcagggcacctaggaaagctgaaaaaca
agtagttgcatttggccgggctgtggttcagatgaagaactggaagacaaccccaaccagagtgacctg
attgagcaggcagccgagatgctttatggattgatccacgcccgctacatccttaccaaccgtggcatc
gcccagatggtgaggcctctctgctcctacctgcctccttctgagcagtaagagacacaggttcctgca
gcaagaagtcatgtttaagccctgtttaaggaagctagctgagaagaggggaagaaccccagaacttgg
ccctgccctaatttggaagaaaggcaacacagaagtttgagagcccatctagtccagagaaggggggcct
ctggacagagttggaaggagtgccgacagagttggtatgggttgggctgcgaagggagttgcctcttct
ttacatctacctgccaacccttccattgtattcacctcagttggaaaagtaccagcaaggagactttg
gttactgtcctcgtgtgtactgtgagaaccagccaatgcttcccattGgtgagtgttgaagaagggaaa
ggaaagcaccgtgtggcagtcttatgggaaggagttggggctcaacacattggagcctgagtcctgagg
ggaggttaggtaggaataggggggatacctggcctgctgagtctggctgtctcccaggcctttcagacat
cccaggtgaagccatggtgaagctctactgccccaagtgcatggatgtgtacacacccaagtcatcaag
acaccatcacacggatggcgcctacttcggcactggtttccctcacatgctcttcatggtgcatcccga
gtaccggcccaagagacctgccaaccagtttgtgcccaggtagggagcagggagagtcattaagggtca
aaggaaaggcccaagatcccccagagagggggaggacagggcatggccctttcttgaggtctgcttctcc
cagaatcagggcatctccctgctgagtgactgtgggaaagttatttgattatctgtgcttgagttacct
tattgtagaatgttcttgagctgagaagttgggaaccacgaggctttagctctgagcaggtccatagag
gagctcaggtggggaggtgggaatgcaggtgactggcagggcctggatggggctcatgctgctgcctct
ctgacctctgccctggcctaggctctacggtttcaagatccatccgatggcctaccagctgcagctcca
agccgccagcaacttcaagagcccagtcaagacgattcgctgattccctcccccacctgtcctgcagtc
tttgtcttttcctttcttttttgccacccttcaggaaccctgtatggtttttagtttaaattaaagga
gtcgttatcgtggtgggaatatgaaataaagtagaagaaaaggccatgagctagtctgctggtgcttgc
ggaagggggtggagcgtggccatggaaatcgggctccacggcccagggatgg

# Approaches to genefinding

- Homology
  - BLAST, Procrustes

- Ab initio (or de novo)
  - Genscan, Genie

- Hybrids
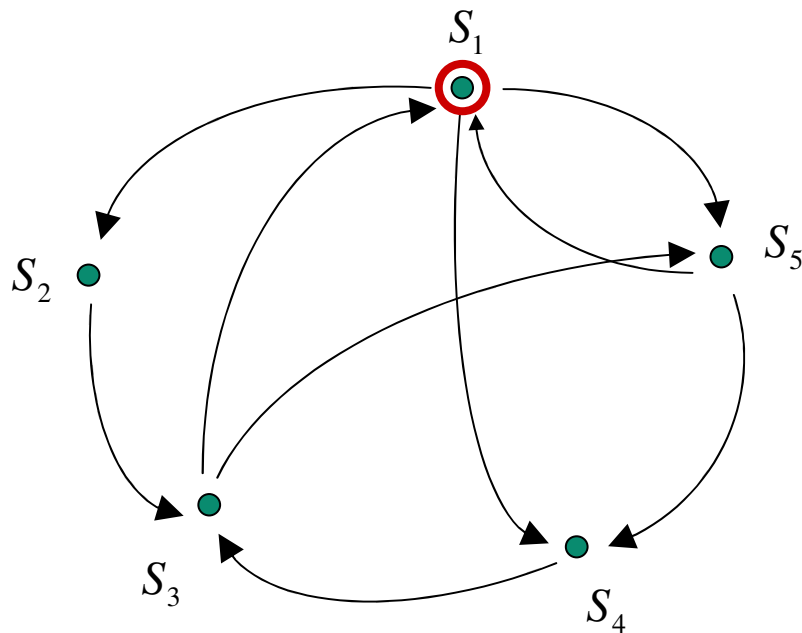  - GenomeScan, GenieEST, Twinscan, SGP, SLAM

# Hidden Markov Models

Fraunhofer **CHALMERS**
Research Centre
Industrial Mathematics

# A discrete process

A random process, jumping between a finite number of states



State sequence:

$S_1$ $S_2$ $S_3$ $S_1$ $S_4$ $S_3$ $S_5$ $S_1$ $S_5$ $S_4$

# A Markov process

The process is *Markov* if the next state only depends on the current state and not the history.

**General description**

↓

$$\Pr(X_t = S_j \mid X_{t-1} = S_i, X_{t-2} = S_k, \ldots) =$$

$$\Pr(X_t = S_j \mid X_{t-1} = S_i)$$

↑

**Markov description**

# A *hidden* Markov Model

A standard hidden Markov model is comprised of two interrelated processes:

- *Hidden process:* a Markov chain on the state space, generating a state sequence hidden from the observer.

- *Observed process:* generating output through random functions associated with each state.

Fraunhofer

**CHALMERS**
Research Centre
Industrial Mathematics

# A simple Hidden Markov Model (HMM)



$$b_A(i) = 1/6$$

$$P_{AB} = 1 - P_{AA}$$

$$P_{BB}$$

$$P_{AA}$$

$$P_{BA} = 1 - P_{BB}$$

$$b_B(i) = 1/4$$

Initial distribution:

$$\pi = (\pi_A, \pi_B)$$

# 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

Fraunhofer

**CHALMERS**
Research Centre
Industrial Mathematics

# The HMM algorithms

The forward variables:

$$\alpha_t(i) = \mathrm{Pr}(\text{obs. up to } t, \text{ ending in state } i \text{ at time } t)$$

The backward variables:

$$\beta_t(i) = \mathrm{Pr}(\text{obs. after } t \mid \text{ending in state } i \text{ at time } t)$$

The Viterbi variables:

$$\delta_t(i) = \max_{X_1,\ldots,X_t} \mathrm{Pr}(\text{obs. up to } t, \text{ ending in state } i \text{ at time } t)$$

# The hidden process:

Let $\{X_t\}_{t=1}^{T}$ denote the Markov process assuming values in the state space $S_1,...,S_N$ .

The process begins in a state determined by the initial distribution $\{\pi_i\}_{i=1}^{N}$ and evolves through the state space according to transition probabilities

$$a_{ij} = \Pr(X_{t+1} = j \mid X_t = i)$$

# The observed process:

Let $\{Y_t\}_{t=1}^{T}$ be the observed process generating output as random functions of the state $X_t$ according to some output function

$$b_j(Y_t \mid Y_1, ..., Y_{t-1}) = \Pr(Y_t \mid X_t = j, Y_1, ..., Y_{t-1})$$

# The joint probability

The joint probability of the two interrelated processes becomes

$$\Pr(X_{t+1} = j, Y_{t+1} \mid X_t = i, X_1, ..., X_{t-1}, Y_1, ..., Y_{t-1}) =$$

$$= \Pr(X_{t+1} = j \mid X_t = i)\,\Pr(Y_{t+1} \mid X_{t+1}, Y_1, ...Y_t) =$$

$$= \underbrace{a_{ij}}_{\text{transition probability}} \underbrace{b_j(Y_{t+1} \mid Y_1, ..., Y_t)}_{\text{output distribution}}$$

**transition probability**

**output distribution**

# Gene finding
## - Generalized HMMs

**Fraunhofer**

**CHALMERS**
Research Centre
Industrial Mathematics

# The Genscan state space

# State duration times

$$\text{Pr(leaving state)} = p$$

$$\text{Pr(staying in state)} = 1 - p$$

$$\text{Pr(output of exactly } r \text{ in state)} = (1-\text{p})^{\text{r-1}} \, p$$

The geometric distribution

# Observed duration times

# Generalized HMMs

When in state $X_l$ the duration $d_l$ is chosen from a generalized length distribution

$$f_{X_l}(d_l) = \text{Pr}(\text{state duration} = d_l \mid X_l)$$

Now the indices for the observed and the hidden process may differ, and we introduce partial sums

$$p_l = \sum_{k=1}^{l} d_k, \quad p_0 = 0$$

Fraunhofer

**CHALMERS**
Research Centre
Industrial Mathematics

ITWM

# Generalized HMMs, cont.

- In state $X_l$ choose state duration $d_l$.
- Generate output $Y_{p_{l-1}+1}, ..., Y_{p_l}$ according to

$$b_{X_l}(Y_{p_{l-1}+1}, ..., Y_{p_l} \mid Y_1, ..., Y_{p_{l-1}})$$

**Hidden process:**   $X_1$    $X_2$    $X_3$

**Observed process:**   $Y_1, ..., Y_{p_1}$   $Y_{p_1+1}, ..., Y_{p_2}$   $Y_{p_2+1}, ..., Y_{p_3}$

# Generalized HMMs, cont.

Assume that we observe a sequence of outputs $Y_1, ..., Y_T$ from a sequence of hidden states $X_1, ..., X_L$ with durations $d_1, ..., d_L$ (assume $p_L = T$). The joint probability of hidden and observed data becomes

$$\Pr(Y_1^T, X_1^L, d_1^L) =$$

$$= \pi_{X_1} f_{X_1}(d_1) b_{X_1}(Y_1^{p_1}) \prod_{l=2}^{L} a_{X_{l-1}, X_l} \underbrace{f_{X_l}(d_l)}_{\text{duration distribution}} \underbrace{a_{X_{l-1}, X_l}}_{\text{transition probability}} f_{X_l}(d_l) b_{X_l}(Y_{p_{l-1}+1}^{p_l} \mid Y_1^{p_{l-1}})$$

duration distribution

transition probability

output distribution

# Sequence alignment
## *- Pair HMMs*

Fraunhofer **CHALMERS**
Research Centre
Industrial Mathematics

# 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 TTCAATCTCAGCAACGCCATCATGGGCAGTGGAATTCTGGGGCTCGCCTA
```

Fraunhofer

**ITWM**

**CHALMERS**
Research Centre
Industrial Mathematics

# Sequence comparisons. Why?

- Are the sequences related?

- What regions are related?

- How evolutionary distant are they?

- Info about the evolutionary process.

Fraunhofer

**CHALMERS**
Research Centre
Industrial Mathematics

# Pair HMMs for alignment



**M** = (mis)match
**X** = insert seq1
**Y** = insert seq2

# Pair HMMs



**Output sequence:**
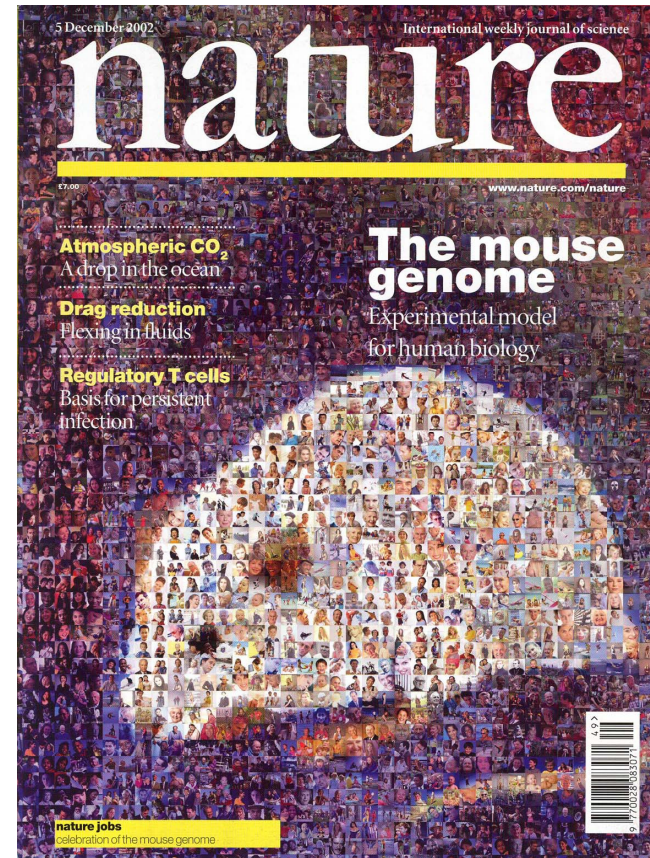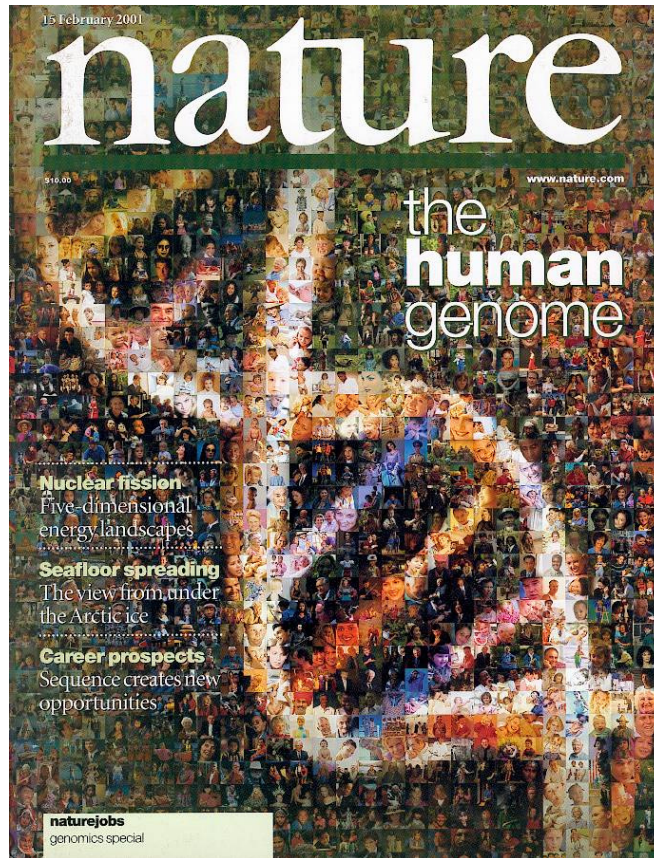
ATCG- - G
AC- GTCA

**Observed sequences:**

ATCGG
ACGTCA

Fraunhofer

**CHALMERS**
Research Centre
Industrial Mathematics

# Comparative gene finding
## - Generalized Pair HMMs

# Comparing human and mouse

**Why mouse?**

Colored By Fayebe

**Human**
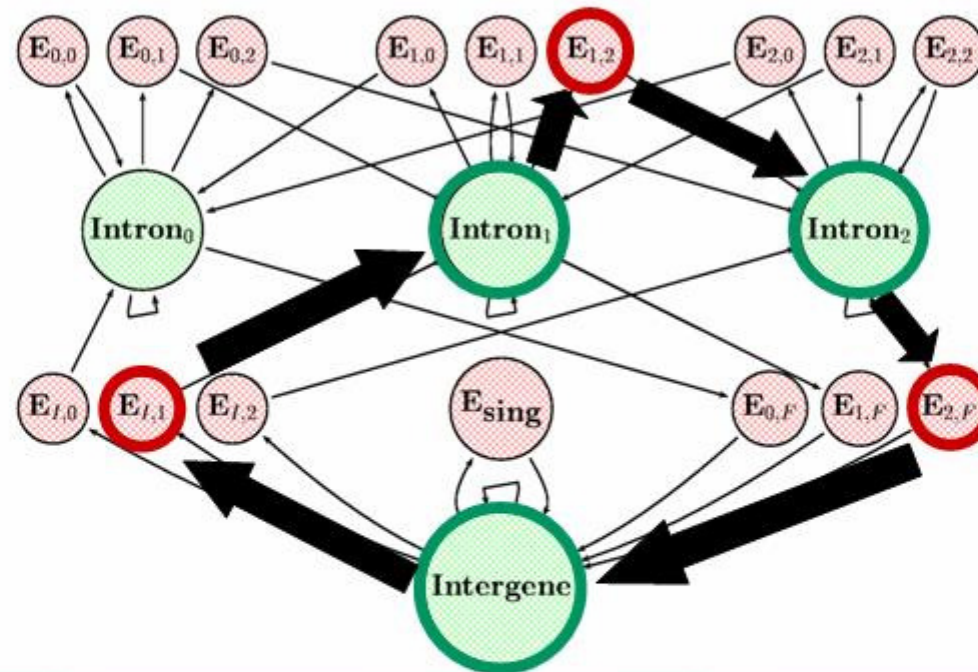
**Mouse**

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

Apolipoprotein AI

Seq1: human
Seq2: macaque

Seq1: human
Seq2: pig

Seq1: human
Seq2: rabbit

Seq1: human
Seq2: mouse

Seq1: human
Seq2: rat

Seq1: human
Seq2: chicken

Exon
CNS

# Comparison of 1196 orthologous gene pairs
## (Makalowski et al., 1996)

- Sequence identity
  - exons: 84.6%
  - protein: 85.4%
  - introns: 36%
  - 5' UTRs: 67%
  - 3' UTRs: 69%

Fraunhofer **CHALMERS**
Research Centre
Industrial Mathematics

# Comparative gene finding

# Generalized Pair HMMs

# Generalized Pair HMMs

- Same hidden process $X_1,...,X_L$ on state space $S_1,...,S_N$
- Two output sequences: $Y_1,...,Y_T$ and $Z_1,...,Z_U$
- Two sets of durations: $d_1,...,d_L$ and $e_1,...,e_L$
- Two sets of partial sums:

$$p_l = \sum_{k=1}^{l} d_k, \quad p_0 = 0, \quad p_L = T \quad \text{and}$$

$$q_l = \sum_{k=1}^{l} e_k, \quad q_0 = 0, \quad q_L = U$$

# Generalized Pair HMMs, cont.

- In state $X_l$ choose state durations $(d_l, e_l)$ from some joint distribution $f_{X_l}(d_l, e_l)$.

- Output $Y_{p_{l-1}+1}, ..., Y_{p_l}$ and $Z_{q_{l-1}+1}, ..., Z_{q_l}$ jointly generated from

$$b_{X_l}(Y_{p_{l-1}+1}^{p_l}, Z_{q_{l-1}+1}^{q_l} \mid Y_1^{p_{l-1}}, Z_1^{q_{l-1}})$$

**Hidden process:** $\qquad X_1 \qquad\qquad X_2 \qquad\qquad X_3$

**Observed process:** $\quad Y_1, ..., Y_{p_1} \qquad Y_{p_1+1}, ..., Y_{p_2} \qquad Y_{p_2+1}, ..., Y_{p_3}$

$$Z_1, ..., Z_{q_1} \qquad Z_{q_1+1}, ..., Z_{q_2} \qquad Z_{q_2+1}, ..., Z_{q_3}$$

# Generalized Pair HMMs, cont.

The joint probability of hidden and observed data, now $Y_1, ..., Y_T$ becomes

$$\Pr(Y_1^T, Z_1^U, X_1^L, d_1^L, e_1^L) =$$

$$= \pi_{X_1} f_{X_1}(d_1, e_1) b_{X_1}(Y_1^{p_1}, Z_1^{q_1}) \prod_{l=2}^{L} a_{X_{l-1}, X_l} f_{X_l}(d_l, e_l) b_{X_l}(Y_{p_{l-1}+1}^{p_l}, Z_{q_{l-1}+1}^{q_l} \mid Y_1^{p_{l-1}}, Z_1^{q_{l-1}})$$

*duration distribution*

*transition probability*

*output distribution*

# Reducing computational complexity

# Computational complexity

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

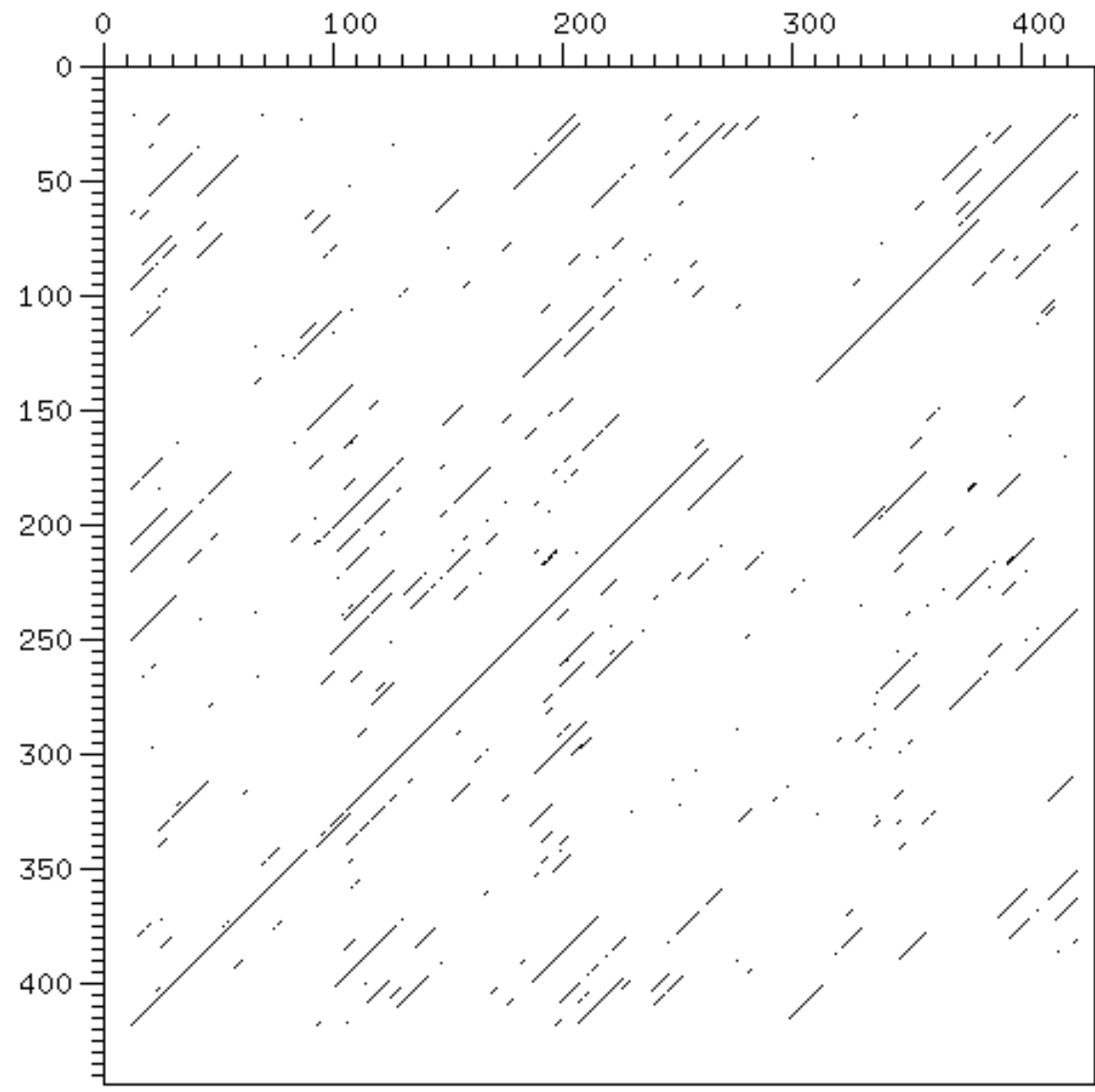$N =$ no. states $\qquad T =$ length seq1

$D =$ max duration $\qquad U =$ length seq2
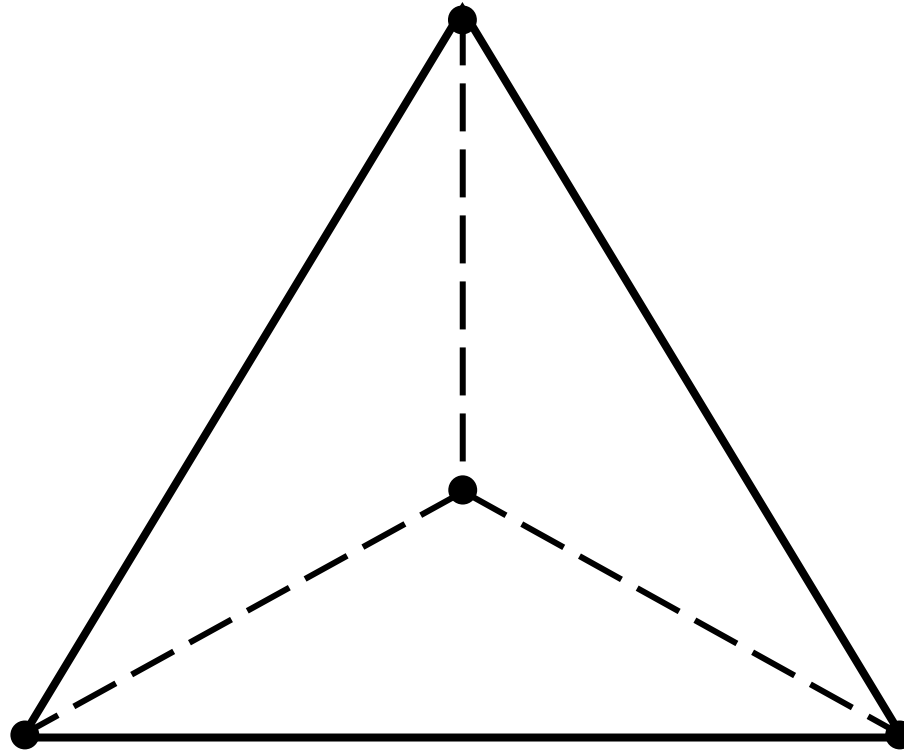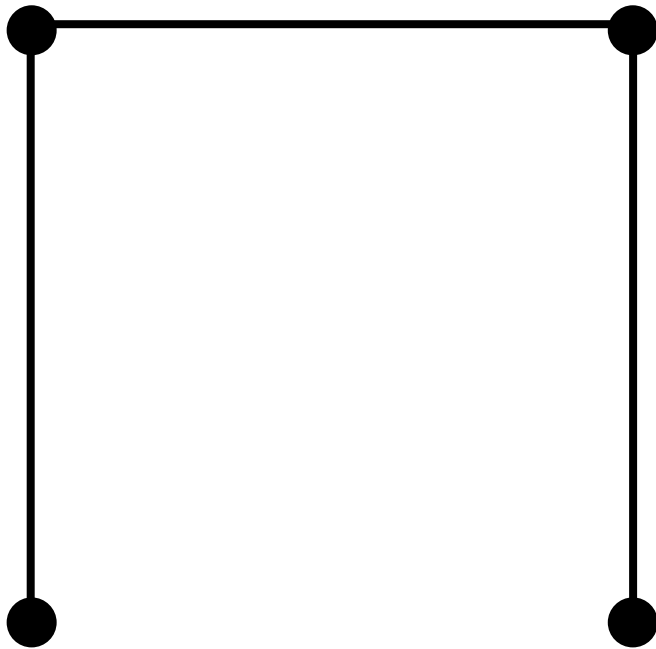
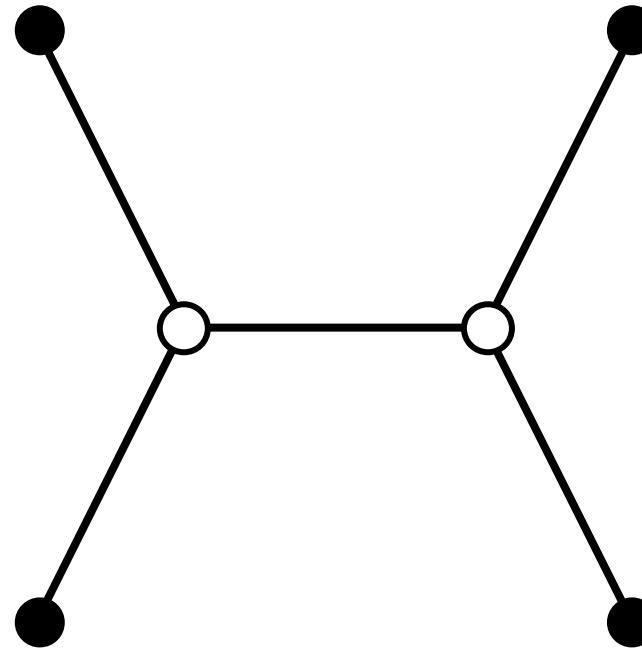# Approximate alignment



**Reduces**

*TU* -factor to *hT*

# Steiner trees

# Steiner trees



**Minimum spanning tree**

**Steiner tree**

10'000 bp

~ $10^9$ edges

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