# One Step Mutation (OSM) matrices

joint work with

Tanja GESELL
PhD student
+43 1 79044 4588
tanja.gesell@mfpl.ac.at

Steffen KLAERE
Codspot
+43 1 79044 4586
steffen.klaere@mfpl.ac.at

CIBIV MFPL
Center for Integrative
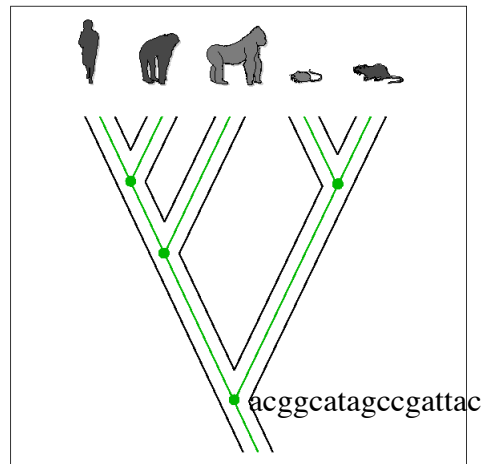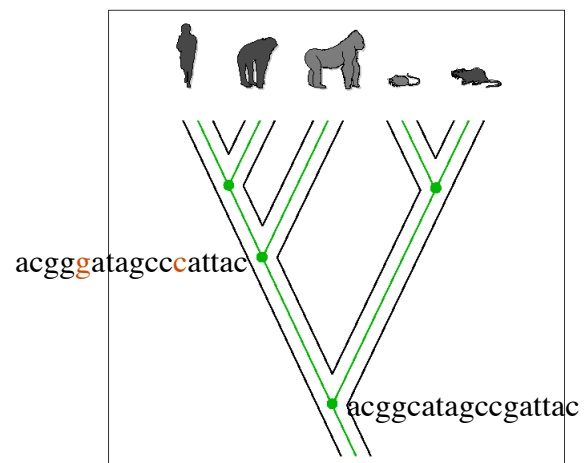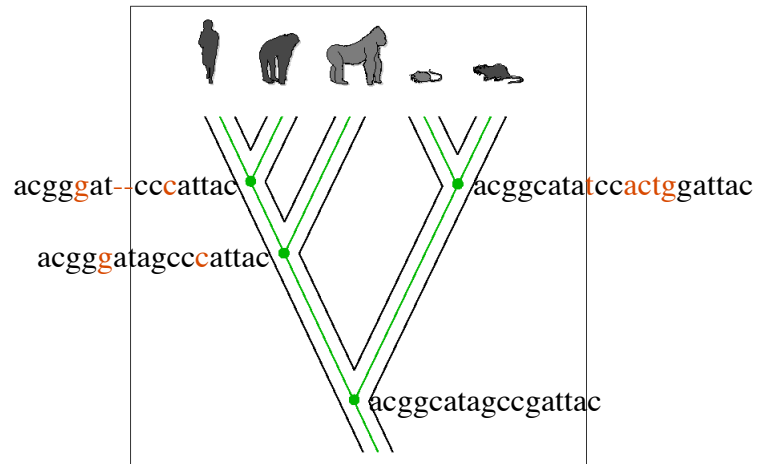Bioinformatics Vienna

---

## Sequence Evolution

CIBIV MFPL
Center for Integrative
Bioinformatics Vienna

Sequence Evolution

acggcatagccgattac



Sequence Evolution

acgggatagcccattac

acggcatagccgattac

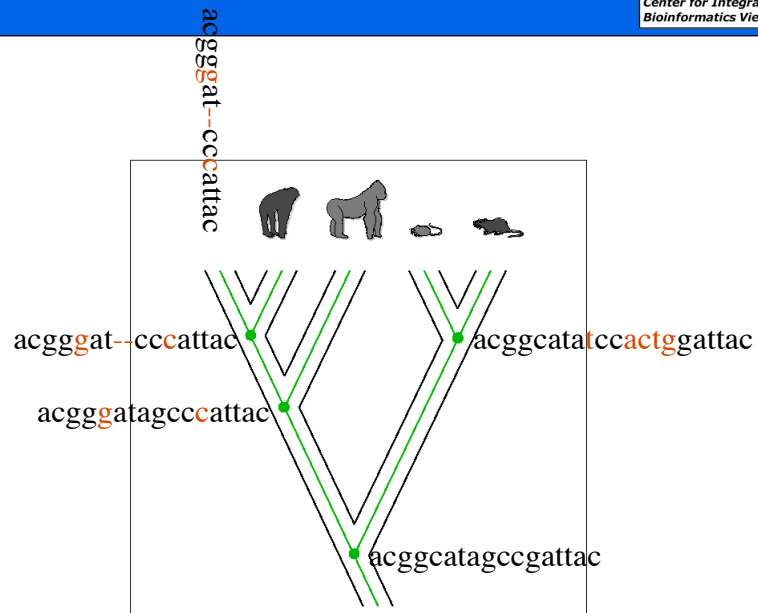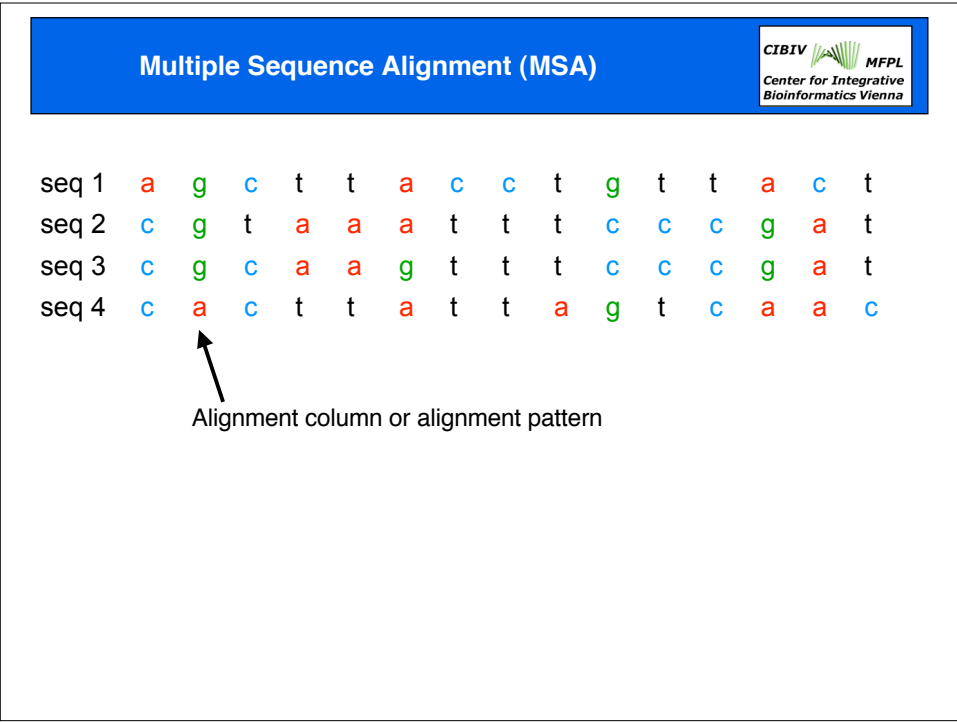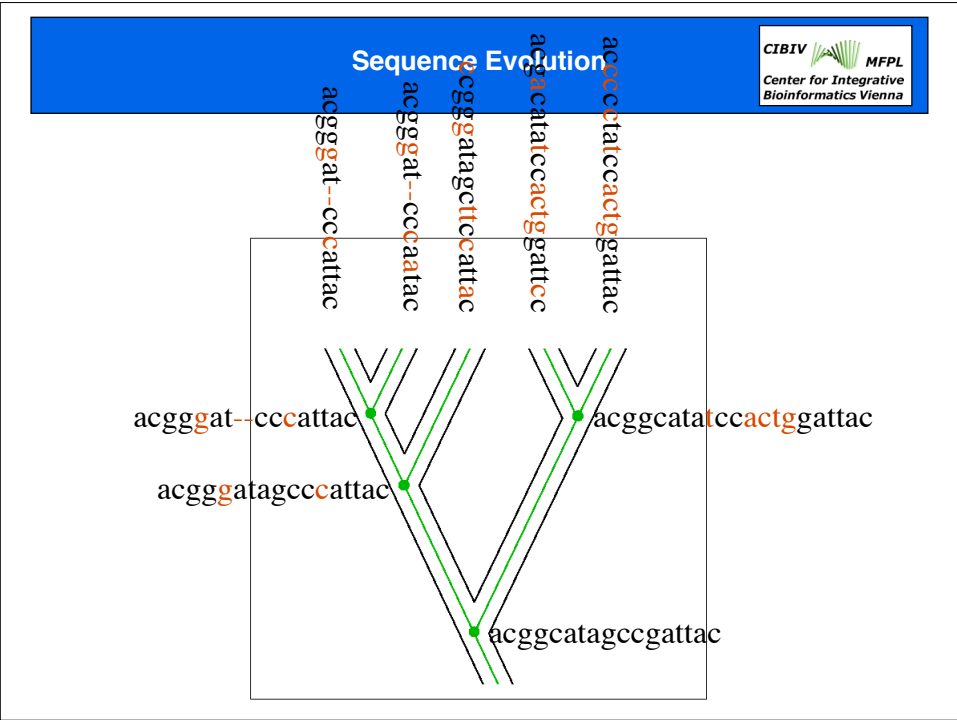**Sequence Evolution** — CIBIV / MFPL Center for Integrative Bioinformatics Vienna

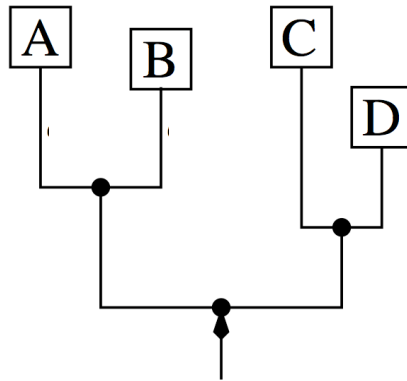**Multiple Sequence Alignment (MSA)** — CIBIV / MFPL Center for Integrative Bioinformatics Vienna

```
seq 1   a  g  c  t  t  a  c  c  t  g  t  t  a  c  t
seq 2   c  g  t  a  a  a  t  t  t  c  c  c  g  a  t
seq 3   c  g  c  a  a  g  t  t  t  c  c  c  g  a  t
seq 4   c  a  c  t  t  a  t  t  a  g  t  c  a  a  c
```

Alignment column or alignment pattern

**Example: Binary Alphabet {R, Y}**

CIBIV / MFPL
Center for Integrative
Bioinformatics Vienna



**Binary Alphabet {R, Y}**

CIBIV / MFPL
Center for Integrative
Bioinformatics Vienna

pattern    R    R    R    R

A    B        C    D

Binary Alphabet {R, Y}

CIBIV MFPL
Center for Integrative
Bioinformatics Vienna

permutation matrix $\sigma_A$



Binary Alphabet {R, Y}

CIBIV MFPL
Center for Integrative
Bioinformatics Vienna

permutation matrix $\sigma_{AB}$

7

**Internal branches**

CIBIV    MFPL
Center for Integrative
Bioinformatics Vienna

matrix multiplication
$\sigma_A \, \sigma_B = \sigma_{AB}$



**One Step Mutation Matrix**

CIBIV    MFPL
Center for Integrative
Bioinformatics Vienna

A   B   C   D

0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5



8

## Examples of OSM-Graphs

## Branch Lengths:

**Total branch length**

$$\Delta = d_A + d_B + d_C + d_D + d_{AB} + d_{CD}$$

**relative edge length**

$$p_{edge} = \frac{d_{edge}}{\Delta}$$



9

**Some Formalisms:**

CIBIV / MFPL
Center for Integrative
Bioinformatics Vienna

relative branch lengths

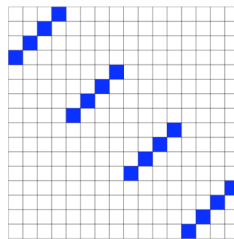$$p_A + p_B + p_C + p_D + p_{AB} + p_{CD} = 1$$

used to assign mutation probabilities

$$p_{edge} \cdot \sigma_{edge}$$

*general permutation matrix*

A

B

C

D

$p_A$

$p_B$

$p_C$

$p_D$

$p_{AB}$

$p_{CD}$

---



**Constructing the OSM:**

CIBIV / MFPL
Center for Integrative
Bioinformatics Vienna

A

B

C

D

$p_A$

$p_B$

$p_C$

$p_D$

$p_{AB}$

$p_{CD}$

$$\mathbf{M}_T = \sum_{edge} p_{edge} \cdot \sigma_{edge}$$

$$\mathbf{M}_T = \sum_{edge} p_{edge} \cdot \sigma_{edge}$$

**One substitution**

$$\mathbf{M}_T^k = \left( \sum_{edge} p_{edge} \cdot \sigma_{edge} \right)^k$$

***k* substitutions**

$\sigma_{CD} \times \sigma_{AB} \times \sigma_C \times \sigma_B$

$$\underset{k}{\mathrm{Min}}\left\{\mathbf{M}_T^k(i,j) > 0 \,\middle|\, k \in N\right\}$$

describes the minimal number of mutations to move from pattern $i$ to $j$

MP: For a tree $T$ and pattern $j$ compute:

$$\underset{k}{\mathrm{Min}}\left\{\mathbf{M}_T^k(R...R, j) > 0 \text{ or } \mathbf{M}_T^k(Y...Y, j) > 0 \,\middle|\, k \in N\right\}$$

---

We assume that the number of substitutions is Poisson distributed with parameter $\Delta$. Then we compute, the expected OSM as

$$\overline{\mathbf{M}_T} = \sum_{k=0}^{\infty} \frac{\exp(-\Delta)\Delta^k(\mathbf{M}_T)^k}{k!}$$
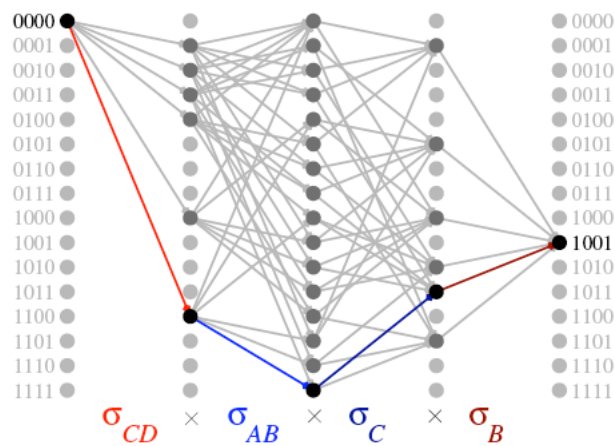
$$\overline{\mathbf{M}_T} = \exp(-\Delta) \cdot \exp(\Delta \cdot \mathbf{M}_T)$$

$$\overline{\mathbf{M}_T} = \exp(-\Delta) \cdot \mathbf{H}_{2^n} \cdot \exp(\Delta \cdot \mathbf{D}_T) \cdot \mathbf{H}_{2^n}$$

where $\qquad \mathbf{D}_T = \mathbf{H}_{2^n} \cdot \mathbf{M}_T \cdot \mathbf{H}_{2^n}$

and $\qquad \mathbf{H}_{2^n} = \underbrace{\mathbf{H}_2 \otimes \mathbf{H}_2 \otimes \ldots \otimes \mathbf{H}_2}_{n \text{ times}} \qquad \mathbf{H}_2 = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$

$$\overline{\mathbf{M}_T} = \exp(-\Delta) \cdot \mathbf{H}_{2^n} \cdot \exp(\Delta \cdot \mathbf{D}_T) \cdot \mathbf{H}_{2^n}$$

The likelihood of a tree $T$ with branch length $\Delta$, given an alignment of length $L$ is then

$$\Pr(T, \Delta) = \prod_{i=1}^{L} \overline{\mathbf{M}_T}(\{R...R, Y...Y\}, pattern(i))$$

$$\overline{\mathbf{M}_T} = \exp(-\Delta) \cdot \mathbf{H}_{2^n} \cdot \exp(\Delta \cdot \mathbf{D}_T) \cdot \mathbf{H}_{2^n}$$

From the above formula, we can **analytically** compute the posterior probability of the number of mutations that have occurred on a fixed tree.

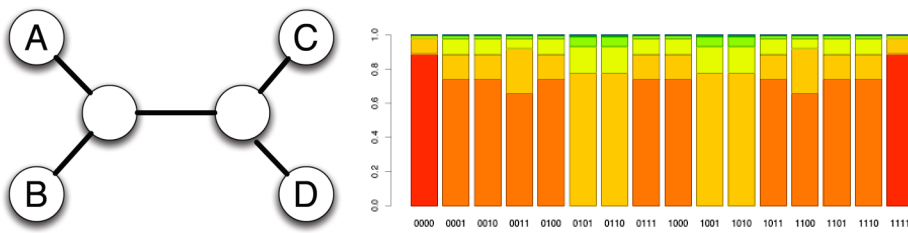$$\Pr(k \text{ mutations} \mid \text{pattern}) = \frac{\exp(-\Delta)\Delta^k}{k!} \frac{\left(\mathbf{M}_T(R,...,R,\text{pattern})\right)^k}{\overline{\mathbf{M}_T}(R...,R,\text{pattern})}$$

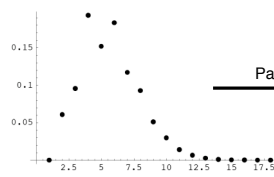similar work by Rasmus Nielsen, John Huelsenbeck, Jonathan Bollback (2002, 2003, 2005)

Posterior probabilities:clock-like tree

$$\text{ppd}[k\,|\,\mathbf{a}] = \frac{\exp[-\Delta]\Delta^{k}\big(\pi_0\mathbf{M}_T^{k}(\mathbf{0},\mathbf{a}) + \pi_1\mathbf{M}_T^{k}(\mathbf{1},\mathbf{a})\big)}{\pi_0\overline{\mathbf{M}}_T(\mathbf{0},\mathbf{a}) + \pi_1\overline{\mathbf{M}}_T(\mathbf{1},\mathbf{a})}$$
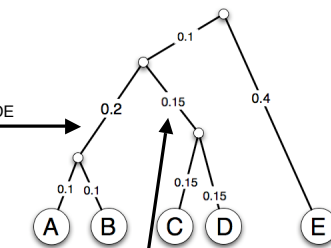
$\Delta = 1.0$



Posterior probabilities: five Taxa Tree

Pattern: ABICDE
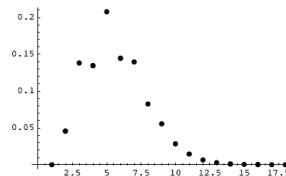
alignment patterns

Pattern ABEICD

## Summary and Outlook

Developed an evolutionary model that describes the action of a single substitution on an alignment pattern.

This leads to a tree-topology mediated random walk on the space of words of length $n$.

Maximum Parsimony and Maximum Likelihood are "extreme" cases within this framework.

Practical Aspect: Analytical formula for the posterior probabilities of the number of substitutions for a pattern.

**Open Questions:**
•Connection between OSM and Hadamard transform (Hendy, Penny 1989) and its generalization, the Fourier calculus on evolutionary trees (Szekely, Steel, Erdös 1993).

•Other type of substitution distributions?

•Computational issues

---

## The real stuff

```
..GUCAUAGAGGGUGAGAAUCCCGUG..
..GCCGGAGAGGGUGACAGCCCCAUC..
..CCCGUGGACGGUGUGAGGCCGGUA..
..GUGAUACAGGGUGACAACCCCGUA..
..ACCAGAGAAGGUGAAAGUCCUGUA..
..GCGAUACAGGGUGACAGCCCCGUA..
```

## The real stuff

```
..GUCAUAGAGGGUGAGAAUCCCGUG..
..GCCGGAGAGGGUGACAGCCCCAUC..
..CCCGUGGACGGUGUGAGGCCGGUA..
..GUGAUACAGGGUGACAACCCCGUA..
..ACCAGAGAAGGUGAAAGUCCUGUA..
..GCGAUACAGGGUGACAGCCCCGUA..
```

Observed pattern count

$$\mathbf{O}\left(d_1,\ldots,d_{4^n}\right)$$

---

Maximum
likelihood etc.

## The real stuff

```
..GUCAUAGAGGGUGAGAAUCCCGUG..
..GCCGGAGAGGGUGACAGCCCCAUC..
..CCCGUGGACGGUGUGAGGCCGGUA..
..GUGAUACAGGGUGACAACCCCGUA..
..ACCAGAGAAGGUGAAAGUCCUGUA..
..GCGAUACAGGGUGACAGCCCCGUA..
```

Observed pattern count

$$\mathbf{O}\left(d_1,\ldots,d_{4^n}\right)$$

Maximum likelihood etc.

$$\mathrm{E}\left(p_1,\ldots,p_{4^n}\right)$$

$$\hat{T}$$

---

## The real stuff

```
..GUCAUAGAGGGUGAGAAUCCCGUG..
..GCCGGAGAGGGUGACAGCCCCAUC..
..CCCGUGGACGGUGUGAGGCCGGUA..
..GUGAUACAGGGUGACAACCCCGUA..
..ACCAGAGAAGGUGAAAGUCCUGUA..
..GCGAUACAGGGUGACAGCCCCGUA..
```

Observed pattern count

$$\mathbf{O}\left(d_1,\ldots,d_{4^n}\right)$$

Maximum likelihood etc.

$$\mathrm{E}\left(p_1,\ldots,p_{4^n}\right)$$

$$\hat{T}$$

**OSM**