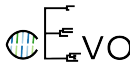# Computational Biology

Lecturers:
Tanja Stadler, Carsten Magnus & Tim Vaughan

Teaching Assistants:
Jūlija Pečerska, Jérémie Sciré,
Sarah Nadeau & Marc Manceau

Computational Evolution
Department of Biosystems Science and Engineering

HS 2019

# How to study evolution?

The easiest way to study something is by observation.

# How to study evolution?

The easiest way to study something is by observation.

- ▶ Wetlab
  - Very realistic;
  - Time-consuming and expensive;
  - Impossible (sometimes).

# How to study evolution?

The easiest way to study something is by observation.

- ▶ Wetlab
    - Very realistic;
    - Time-consuming and expensive;
    - Impossible (sometimes).
- ▶ Simulation
    - A virtual experiment in which we mimic a (biological) process on a computer to study its properties
    - Not necessarily realistic
    - Allows us to:
        - * generate data with given assumptions;
        - * test predictive properties of models.

# How to study evolution?

The easiest way to study something is by observation.

- ▶ Wetlab
    - Very realistic;
    - Time-consuming and expensive;
    - Impossible (sometimes).
- ▶ Simulation
    - A virtual experiment in which we mimic a (biological) process on a computer to study its properties
    - Not necessarily realistic
    - Allows us to:
        * generate data with given assumptions;
        * test predictive properties of models.

Today we will simulate evolution!

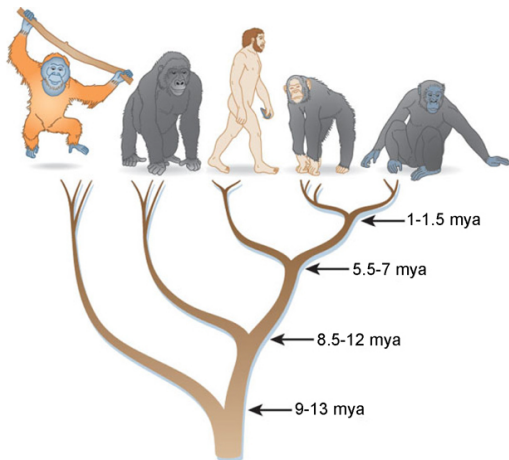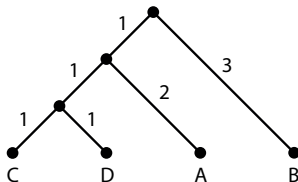# The tree of great apes

Figure adapted from [?]

# Storing trees: Newick format

- ▶ Format for tree representation
- ▶ To record a tree in Newick format:
    - Assign a label to each tip
    - Choose two tips that are a cherry (e.g. C and D)
    - Replace selected tips with a new tip of the form
      (tip1:branch1,tip2:branch2) (e.g. $(C : 1, D : 1)$)
        - Branch length to the new tip is the branch length to the
          cherry
    - Repeat until the full tree is rewritten
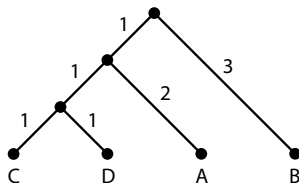- ▶ What is the Newick format for the rooted tree above?

# Storing trees: Newick format

- ▶ Format for tree representation
- ▶ To record a tree in Newick format:
    - Assign a label to each tip
    - Choose two tips that are a cherry (e.g. C and D)
    - Replace selected tips with a new tip of the form
      (tip1:branch1,tip2:branch2) (e.g. $(C : 1, D : 1)$)
        - Branch length to the new tip is the branch length to the
          cherry
    - Repeat until the full tree is rewritten
- ▶ What is the Newick format for the rooted tree above?
  $(((C : 1, D : 1) : 1, A : 2) : 1, B : 3);$

# Evolution Simulation Algorithm

**Steps**:

1. **Initialization of the starting sequence**:
   - Sample a starting nucleotide for each position in the sequence

2. **Iterative simulation** of sequence evolution, along all branches of the tree
   - Compute the transition probability matrix $P(t_b)$.
   - Sample a new nucleotide for each position in the sequence.

# Step 1: Initialization of the starting sequence

1a. Sample a starting nucleotide $n$

# Step 1: Initialization of the starting sequence

---

1a. Sample a starting nucleotide $n$

---

From the vector of equilibrium frequencies of nucleotides

|     | T    | C    | A    | G    |
|-----|------|------|------|------|
| Π   | 0.22 | 0.26 | 0.33 | 0.19 |

# Step 1: Initialization of the starting sequence

1a. Sample a starting nucleotide $n$

From the vector of equilibrium frequencies of nucleotides

|       | T    | C    | A    | G    |
|-------|------|------|------|------|
| $\Pi$ | 0.22 | 0.26 | 0.33 | 0.19 |

Knowing $\Pi$, how do we sample a nucleotide?

# Inverse transform method

▶ Sample $u$ from $U(0, 1)$;

# Inverse transform method

- Sample $u$ from $U(0,1)$;
- Transform $u$ into a sample from the desired distribution using the **CDF** == **C**umulative **D**istribution **F**unction $F_X(x) = P(X \leqslant x)$.

# Inverse transform method

- Sample $u$ from $U(0, 1)$;
- Transform $u$ into a sample from the desired distribution using the **CDF == C**umulative **D**istribution **F**unction $F_X(x) = P(X \leqslant x)$.

# Inverse transform method

- Sample $u$ from $U(0,1)$;
- Transform $u$ into a sample from the desired distribution using the **CDF == C**umulative **D**istribution **F**unction $F_X(x) = P(X \leqslant x)$.

# Inverse transform method

- Sample $u$ from $U(0, 1)$;
- Transform $u$ into a sample from the desired distribution using the **CDF** $==$ **C**umulative **D**istribution **F**unction $F_X(x) = P(X \leqslant x)$.

# Sampling discrete random variables

|     | T    | C    | A    | G    |
| --- | ---- | ---- | ---- | ---- |
| Π   | 0.22 | 0.26 | 0.33 | 0.19 |
| CDF | 0.22 | 0.48 | 0.81 | 1.00 |

# Sampling discrete random variables

|      | T    | C    | A    | G    |
|------|------|------|------|------|
| Π    | 0.22 | 0.26 | 0.33 | 0.19 |
| CDF  | 0.22 | 0.48 | 0.81 | 1.00 |



Sample $u$ from $U(0, 1)$.

# Sampling discrete random variables

|       | T    | C    | A    | G    |
| ----- | ---- | ---- | ---- | ---- |
| Π     | 0.22 | 0.26 | 0.33 | 0.19 |
| CDF   | 0.22 | 0.48 | 0.81 | 1.00 |



Sample $u$ from $U(0,1)$.

E.g. $u = 0.62$.

# Sampling discrete random variables

|       | T    | C    | A    | G    |
|-------|------|------|------|------|
| Π     | 0.22 | 0.26 | 0.33 | 0.19 |
| CDF   | 0.22 | 0.48 | 0.81 | 1.00 |



Sample $u$ from $U(0, 1)$.

E.g. $u = 0.62$.

Select nucleotide **A**.

# Sampling discrete random variables

|      | T    | C    | A    | G    |
|------|------|------|------|------|
| Π    | 0.22 | 0.26 | 0.33 | 0.19 |
| CDF  | 0.22 | 0.48 | 0.81 | 1.00 |



Sample $u$ from $U(0, 1)$.

# Sampling discrete random variables

|       | T    | C    | A    | G    |
|-------|------|------|------|------|
| Π     | 0.22 | 0.26 | 0.33 | 0.19 |
| CDF   | 0.22 | 0.48 | 0.81 | 1.00 |



Sample $u$ from $U(0, 1)$.

E.g. $u = 0.18$.

# Sampling discrete random variables

|     | T    | C    | A    | G    |
|-----|------|------|------|------|
| Π   | 0.22 | 0.26 | 0.33 | 0.19 |
| CDF | 0.22 | 0.48 | 0.81 | 1.00 |



Sample $u$ from $U(0, 1)$.

E.g. $u = 0.18$.

Select nucleotide **T**.

# Step 1: Initializing the starting sequence

1b. Place $n$ on the root node;

# Step 1: Initializing the starting sequence

1b. Place $n$ on the root node;
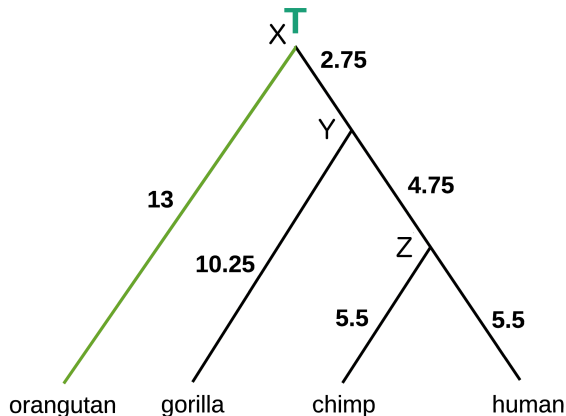
# Step 2a: Choose the next branch for simulation

Get a branch $b$ with a nucleotide at the start;
$t_b = \text{length}(b)$;
$n = $ nucleotide at start of branch $b$;

Get a branch $b$ with a nucleotide at the start;
$t_b = \text{length}(b)$;
$n = $ nucleotide at start of branch $b$;

# Step 2b-d: Sample the new nucleotide

$P(t_b) = e^{Qt_b}$;
Sample new nucleotide $n_{new}$ from row $n$ in $P(t_b)$;
Place $n_{new}$ at the end of branch $b$;

# Step 2b-d: Sample the new nucleotide

$P(t_b) = e^{Qt_b}$;
Sample new nucleotide $n_{new}$ from row $n$ in $P(t_b)$;
Place $n_{new}$ at the end of branch $b$;

To sample new nucleotide $n_{new}$ we will need the substitution rate matrix $Q$, and transition probability matrix $P$.

# Substitution rate matrix – TN93

$\Pi = (\pi_T, \pi_C, \pi_A, \pi_G)$ - equilibrium frequencies.

$\alpha_1, \alpha_2$ - transition ratios.

$\beta$ - transversion ratios.

$$
Q_{TN93} = \begin{array}{c} \\ T \\ C \\ A \\ G \end{array}
\begin{array}{c}
T \qquad\quad C \qquad\quad A \qquad\quad G
\end{array}
\left(
\begin{array}{cccc}
\cdot & \alpha_1\pi_C & \beta\pi_A & \beta\pi_G \\
\alpha_1\pi_T & \cdot & \beta\pi_A & \beta\pi_G \\
\beta\pi_T & \beta\pi_C & \cdot & \alpha_2\pi_G \\
\beta\pi_T & \beta\pi_C & \alpha_2\pi_A & \cdot
\end{array}
\right)
$$

The diagonals are set such that each row sums up to zero, e.g.
$q_{TT} = -(\alpha_1\pi_C + \beta\pi_A + \beta\pi_G)$.

# Substitution rate matrix – TN93

$\Pi = (0.22, 0.26, 0.33, 0.19)$
$\alpha_1 = 44.229, \ \alpha_2 = 21.781$
$\beta = 1$

$$Q_{TN93} = \begin{array}{c} \\ T \\ C \\ A \\ G \end{array} \begin{array}{cccc} T & C & A & G \\ \begin{pmatrix} -0.01957 & 0.01873 & 0.00054 & 0.00031 \\ 0.01584 & -0.01669 & 0.00054 & 0.00031 \\ 0.00036 & 0.00042 & -0.00752 & 0.00674 \\ 0.00036 & 0.00042 & 0.01170 & -0.01249 \end{pmatrix} \end{array}$$

Note: the matrix is scaled to 0.0135 substitutions per mya so that we get reasonable sequences.

# Transition probability matrix – TN93

$\Pi = (\pi_T, \pi_C, \pi_A, \pi_G)$ - equilibrium frequencies.

$\alpha_1, \alpha_2$ - transition ratios.

$\beta$ - transversion ratios.

$t_b$ - branch length

$$P(t_b) = e^{t_b Q_{TN93}(\alpha_1, \alpha_2, \beta, \Pi)}$$

# Substitution rate matrix – TN93

$\Pi = (0.22, 0.26, 0.33, 0.19)$
$\alpha_1 = 44.229,\ \alpha_2 = 21.781$
$\beta = 1$
$t_b$ - 13 mya

$$
P_{TN93}(13mya) = \begin{array}{c} \\ T \\ C \\ A \\ G \end{array}
\begin{array}{cccc}
T & C & A & G \\
\left( \begin{array}{cccc}
0.795 & 0.194 & 0.007 & 0.004 \\
0.164 & 0.824 & 0.007 & 0.004 \\
0.005 & 0.005 & 0.913 & 0.077 \\
0.005 & 0.005 & 0.134 & 0.856
\end{array} \right)
\end{array}
$$

# Sampling substitution times

We start with nucleotide **T**, so we are interested in row T:

$$
P_{TN93}(13mya) = \begin{array}{c} \\ T \\ C \\ A \\ G \end{array} \overset{\begin{array}{cccc} T & C & A & G \end{array}}{\begin{pmatrix} 0.795 & 0.194 & 0.007 & 0.004 \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{pmatrix}}
$$

# Sampling substitution times

We start with nucleotide **T**, so we are interested in row T:

$$P_{TN93}(13mya) = \begin{array}{c} T \\ C \\ A \\ G \end{array} \begin{pmatrix} 0.795 & 0.194 & 0.007 & 0.004 \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \end{pmatrix} \begin{array}{cccc} T & C & A & G \end{array}$$

Sample new nucleotide $n_{new}$ with the weights
$[p_{T \to T}, p_{T \to C}, p_{T \to A}, p_{T \to G}]$

# Getting the substitution

Sample $u$ from $U(0, 1)$.
E.g. $u = 0.81$.

# Getting the substitution

Sample $u$ from $U(0, 1)$.

E.g. $u = 0.81$.



Selected substitution is T $\rightarrow$ C.

# Step 2b-d: Sample the new nucleotide

$P(t_b) = e^{Qt_b}$;
Sample new nucleotide $n_{new}$ from row $n$ in $P(t_b)$;
Place $n_{new}$ at the end of branch $b$;

# Step 2b-d: Sample the new nucleotide

$P(t_b) = e^{Qt_b};$
Sample new nucleotide $n_{new}$ from row $n$ in $P(t_b)$;
Place $n_{new}$ at the end of branch $b$;

# Repeat step 2

**while** *not all branches are used* **do**
    Get a branch $b$ with a nucleotide at the start;
    $t_b = \text{length}(b)$;
    $n = $ nucleotide at start of branch $b$;
    $P(t_b) = e^{Qt_b}$;
    Sample new nucleotide $n_{new}$ from row $n$ in $P(t_b)$;
    Place $n_{new}$ at the start of the daughter branches of $b$;
**end**

# Evolution

# Evolution

# Evolution

# Evolution

# Evolution

# Evolution

# Evolution

# Evolution

# Evolution

# Evolution

# Evolution

# Evolution

# Exercise for today

1. Split into pairs;
2. Get the materials;
3. Evolve a character along the tree;

All of the characters together will produce an alignment.

# Using dice to generate random numbers

We will be using 10-sided dice for our random number generation.

# Using dice to generate random numbers

We will be using 10-sided dice for our random number generation.
Samples from Unif(0,1) with 2 decimal point precision:

1. Take 2 dice of different colours;
2. Assign a fixed decimal position to each of the dice (and keep it for the whole simulation);
   ▶ E.g. red is 1st position, blue is 2nd;
3. Roll the dice to get 2 numbers;
   ▶ E.g. red = 5, blue = 8;
4. Combine the numbers to get a sample;
   ▶ E.g. 0.58.

# Algorithm

$N$ = number of sites in the alignment;
$Q$ = substitution rate matrix;
**for** $i = 1$ **to** $N$ **do**
    Sample a nucleotide $n$ from the initial distribution;
    Add $n$ to the sequence of the root node;
**end**
**while** *not all branches are visited* **do**
    Get a branch $b$ with a sequence at the start;
    $t_b = \text{length}(b)$;
    $P(t_b) = e^{Q t_b}$;
    **for** $i = 1$ **to** $N$ **do**
        $n$ = nucleotide at position $i$ at the start of branch $b$;
        Sample new nucleotide $n_{new}$ from row $n$ in $P(t_b)$;
        Place $n_{new}$ at the end of sequences in the daughter
         branches of $b$;
    **end**
**end**