# *Background*

Inferring the evolutionary relationships among species is one of the oldest and most basic tasks of evolutionary research.



**Since 19th**
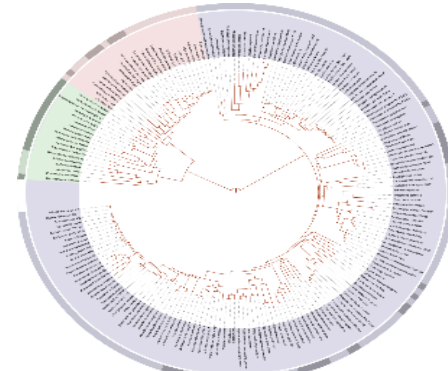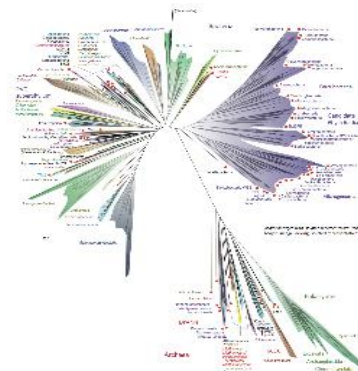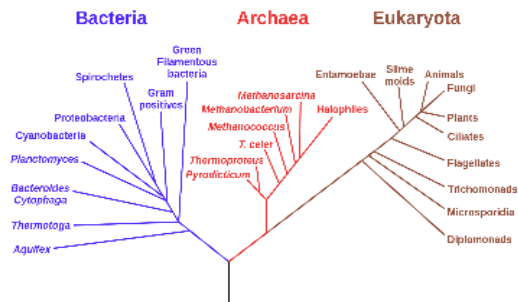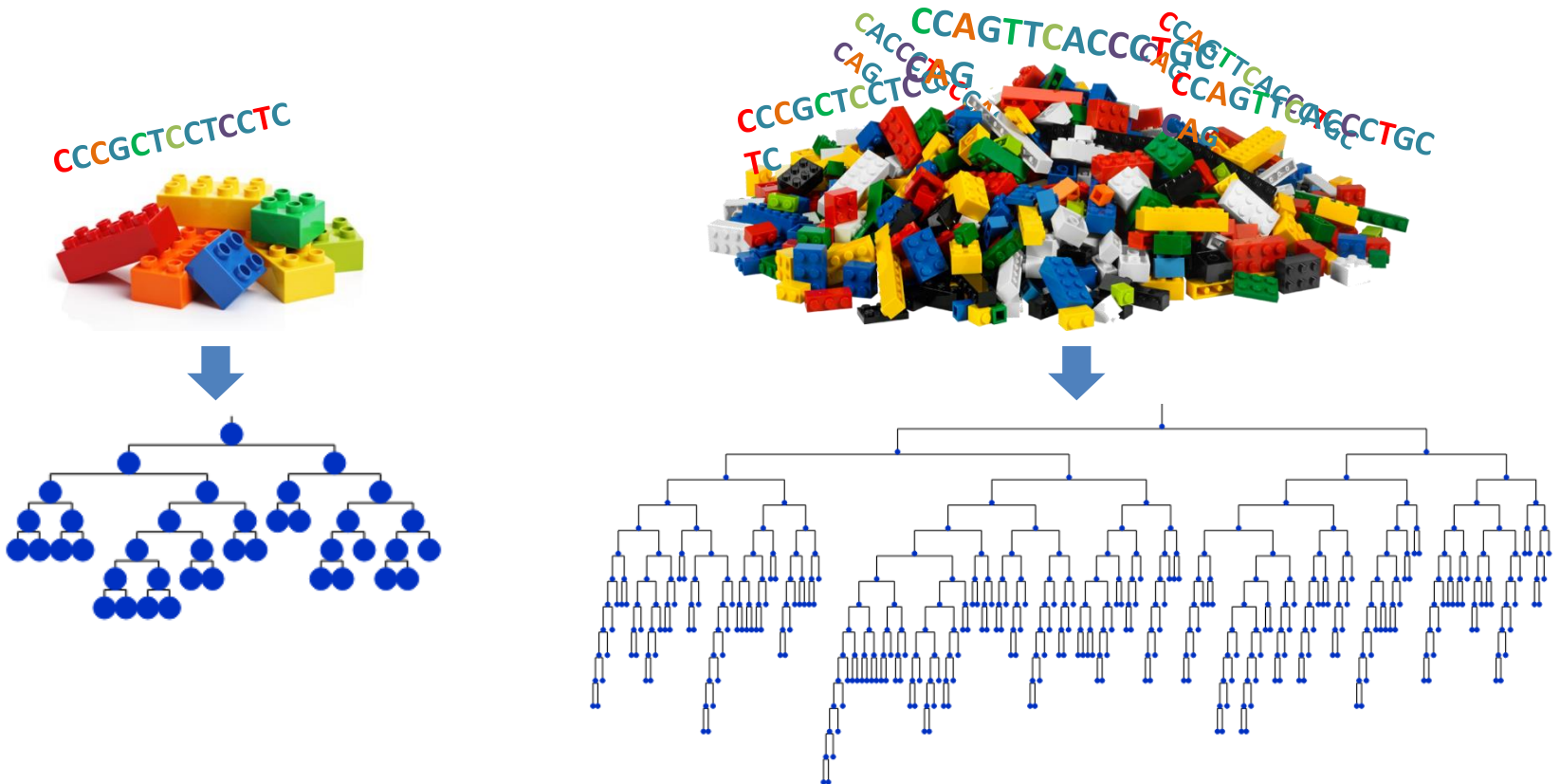
# *Background*

The massive accumulation of sequence data
should provide more accurate phylogenies with
the hope to resolve the tree of life.
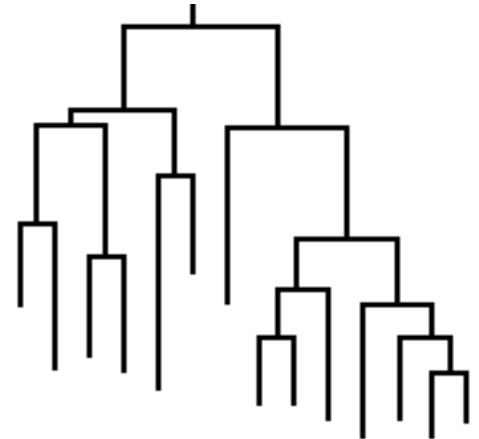
# Phylogeny reconstruction with *one* gene

Sequence1    -TCAGGA-TGAAC----
Sequence2    ATCACGA-TGAACC---
Sequence3    ATCAGGAATGAATCC--
Sequence4    -TCACGATTGAATCGC-
Sequence5    -TCAGGAATGAATCGCM

input

tree building methods

output

# Phylogeny reconstruction with *many* genes

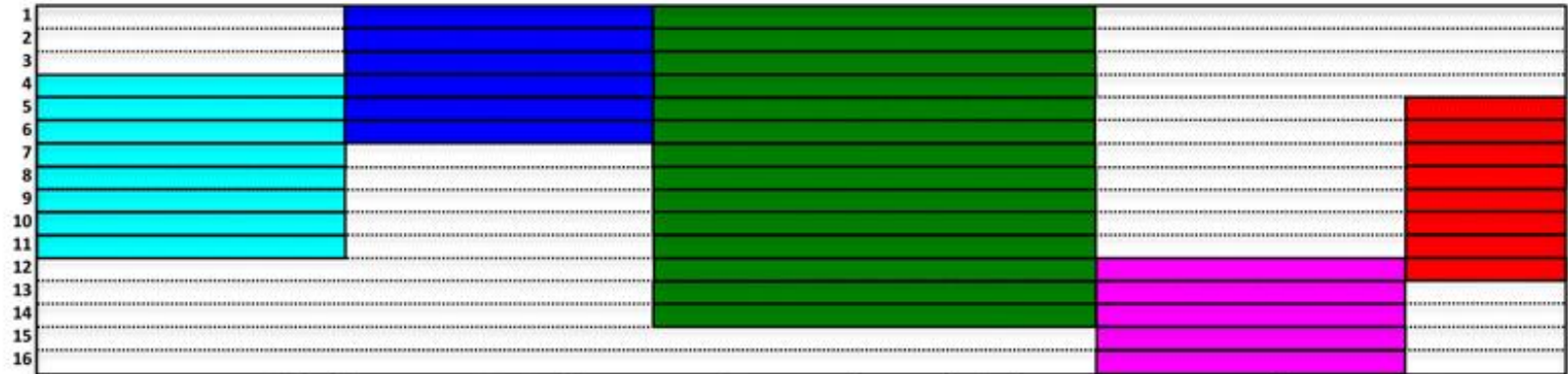**What if we want to use many loci?**
We concatenate genes into a single sequence

gene1

gene2

gene3

GTTCAAAATCACTGCCCGCTCCTCCTCGC **+** CAATGTGAAAGCTGGTG **+** GGCTTGGCGCGACAAAAGCTCCACCTA

**=**

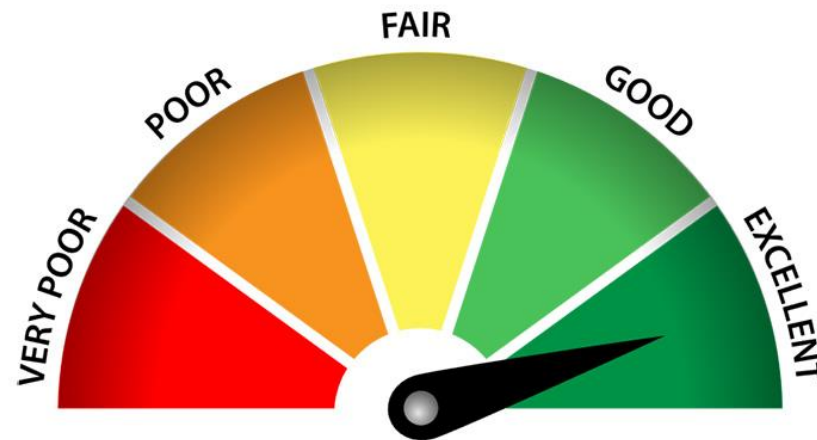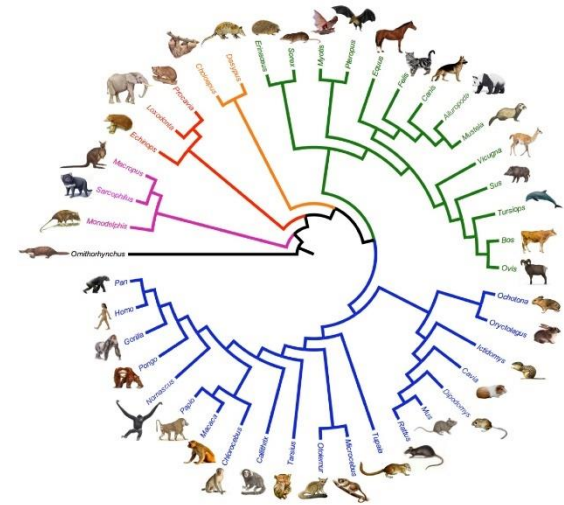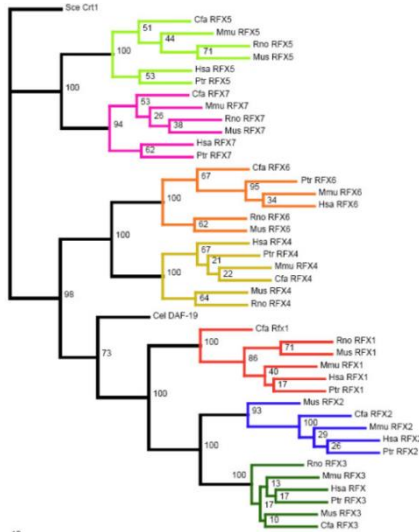GTTCAAAATCACTGCCCGCTCCTCCTCGCCAATGTGAAAGCTGGTGGGCTTGGCGCGACAAAAGCTCCACCTA

# *Super*matrix **method**

# How can we asses the quality of the tree?

# Assessing phylogeny quality

If the tree is correct, we expect the species that are part of the same **taxonomic rank** to form a **monophyletic** group in the tree.

A clade is *monophyletic* if it consists of an ancestral species and all its descendants.
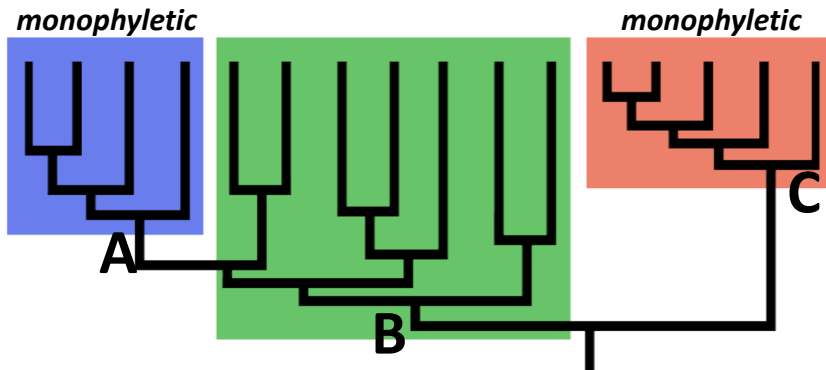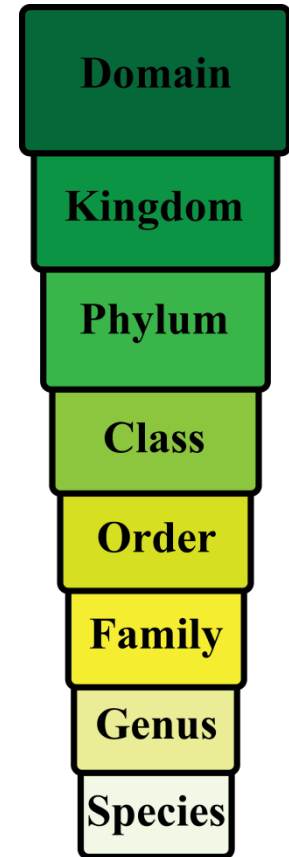


*taxonomic hierarchy*
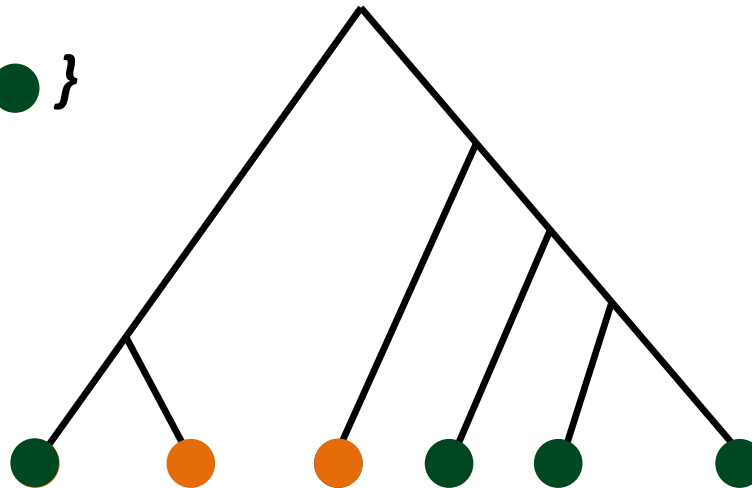
# Assessing phylogeny quality

If the tree is correct, we expect the species that are part of the same **taxonomic rank** to form a **monophyletic** group in the tree.

*Linnaean taxonomy*

For example:

*Genus1 = { ●, ●, ●, ● }*
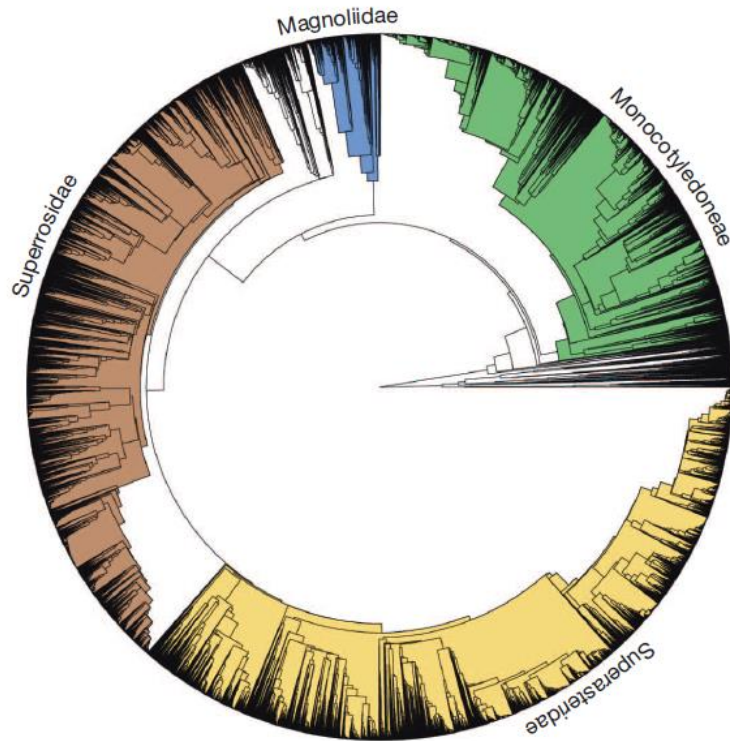*Genus2 = { ●, ● }*

# An example of a large phylogeny



*A phylogeny of land plant species*
*32,223 species*
*Based on 7 loci*

Zanne et al., nature 2014

# However, many genera are far from being monophyletic



*Genus:* *Salvia (99 species)*

*211 species*

The average monophyletic score for >1000 genera is only **0.66**

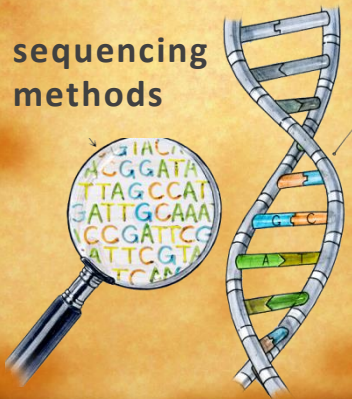*Monophyletic score:*
*99/211 = 0.47*

# Who should we blame of distorting the trees?



WANTED!

LOOKING FOR

sequencing methods

WANTED!

LOOKING FOR

**MSA**
Species1: AT—G—CACATG-CA----CATT
Species2: ATGG-CAC—CCCAAC--GT----
Species3: ATGG-----AA—CCATC—TCA
Species4: --ATG—--AT—GC--GCG—TTA
Species5: ATGGC------A----GC-----TTA
Species6: AT—G—CACATG-CA----CATT
Species7: AG-----TACCC-------AACTGG
Species8: A-GGC-T-----CCC----ACA-A-G
Species9: ---TGGCAAC-CC--CCTACAAT
Species10: A---GCAT-------CCCTA---TC
Species11: -----CATATC-----TACAACTA
Species12: ATGGC-------------C—GCAA
Species13: AT---CCA-CCA-CT-CA-CTAG
Species14: ATG----GCCAT-----GCC---G-
Species15: --TGG--ACAT-TAG--GC—AG
Species16: A--GGCC----GCA-AGA--GGT

WANTED!

LOOKING FOR

tree building methods

# Who should we blame?

## Hypothesis:

**WANTED!**

**LOOKING FOR**

**MSA**

Species1: AT—G—CACATG-CA----CATT
Species2: ATGG-CAC—CCCAAC--GT----
Species3: ATGG-----AA—CCATC—TCA
Species4: --ATG—-AT—GC--GCG—TTA
Species5: ATGGC-------A----GC-----TTA
Species6: AT—G—CACATG-CA----CATT
Species7: AG-----TACCC--------AACTGG
Species8: A-GGC-T-----CCC----ACA-A-G
Species9: ---TGGCAAC-CC--CCTACAAT
Species10: A---GCAT--------CCCTA---TC
Species11: -----CATATC-----TACAACTA
Species12: ATGGC-------------C—GCAA
Species13: AT---CCA-CCA-CT-CA-CTAG
Species14: ATG----GCCAT-----GCC---G-
Species15: --TGG--ACAT-TAG--GC—AG
Species16: A--GGCC----GCA-AGA--GGT

**SUSPECT**

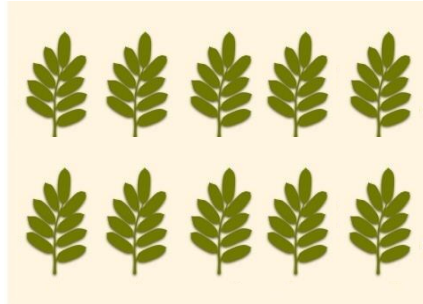Input sequences are too far apart from each other → methods of MSA construction produce highly erroneous alignments → many of the large trees are far from being true.

# Aligning of sequences

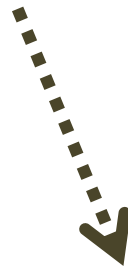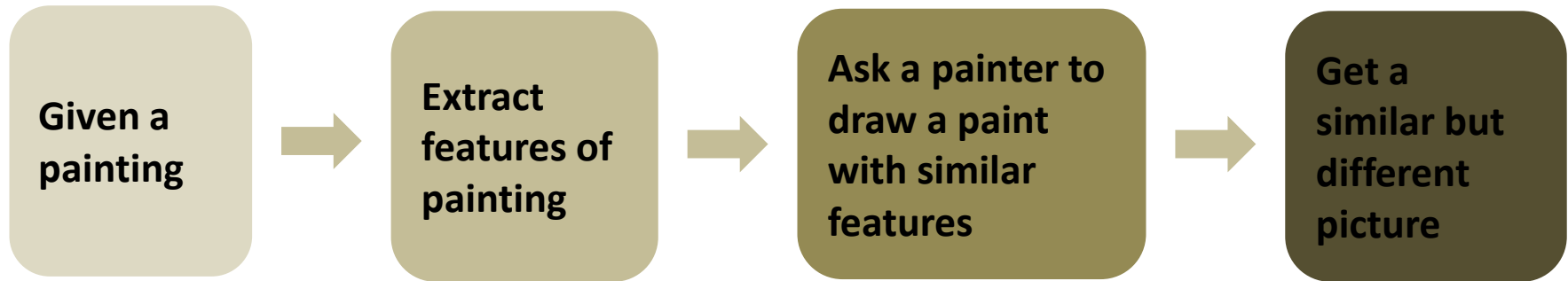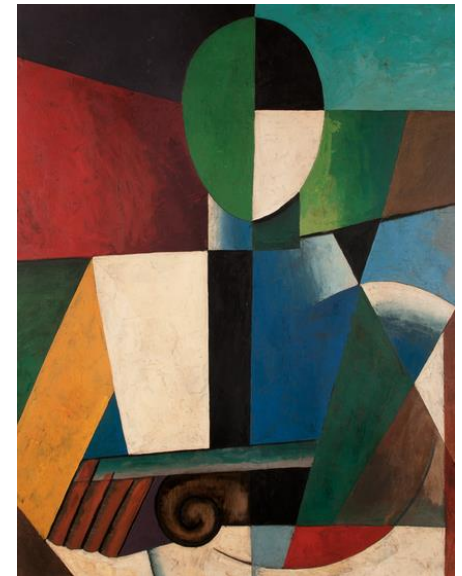|  | | |
|---|---|---|
| **sequences evolves** | slowly | rapidly |
| **aligning is** | easy | hard |
| **good for solving** | ancient divergence events | recent divergence events |

# How to prove that the MSA is guilty?



By simulating MSAs of different sizes (in terms of number of species) and test our hypothesis that the size of the MSA affect the quality of the tree.
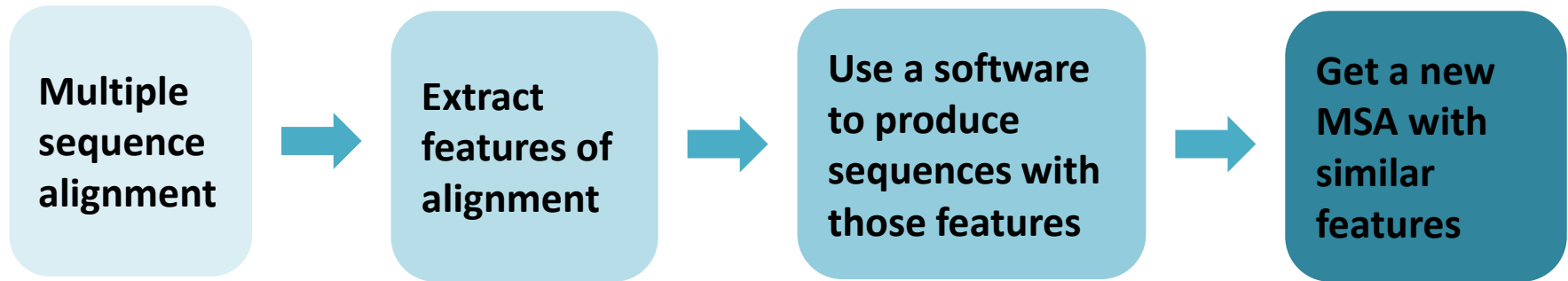
# Understanding the concept of simulations



| Given a painting | → | Extract features of painting | → | Ask a painter to draw a paint with similar features | → | Get a similar but different picture |

- Color scale, shades
- Shapes
- Number of shapes
- Objects
- And more.

# What does it mean to simulate sequences?

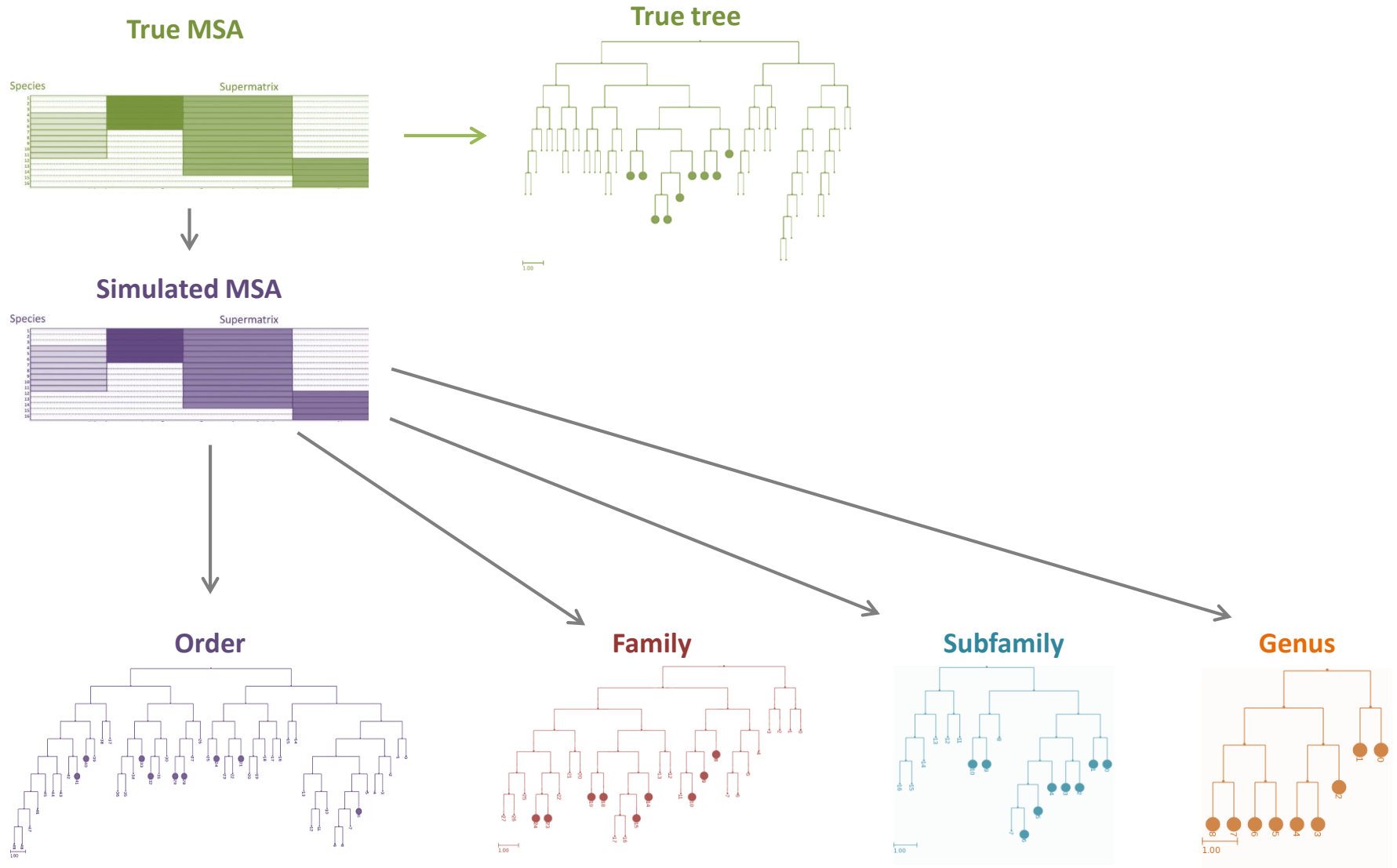**Multiple sequence alignment** → **Extract features of alignment** → **Use a software to produce sequences with those features** → **Get a new MSA with similar features**

- Length of alignment 10
- Number of gaps 1
- Average lengths of gaps 3
- Nucleotides frequencies A: 30 %
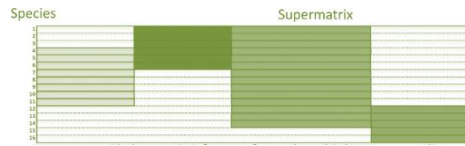- Substitutions rates (A -> C, G -> T, …)
- And more.
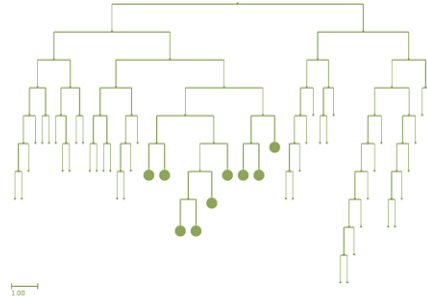
ATTGACCTGA
||| |||||
AT - - -CCTGA
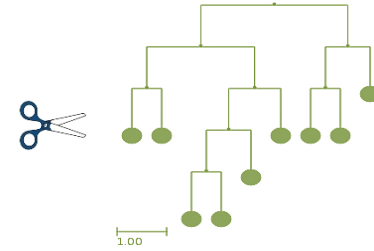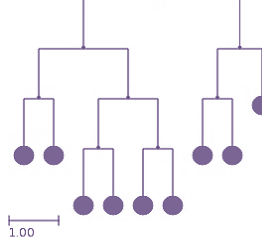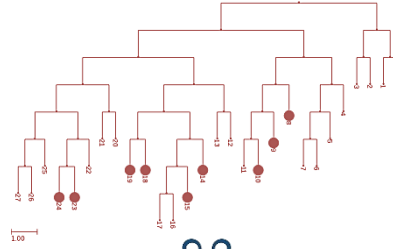
# The effect of MSA size on tree quality

True MSA

Simulated MSA

True tree

Genus

Order

Family

Subfamily

Genus

Genus

Genus

Genus

Order
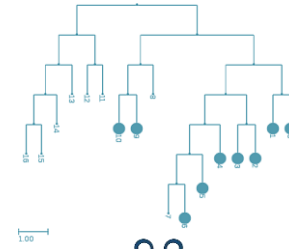
Family
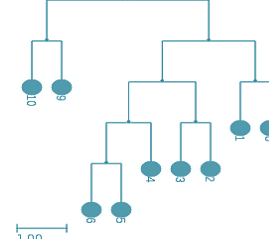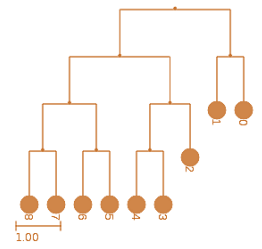
Subfamily

Genus

Genus

Genus

Genus

distance 1

distance 2

distance 3

distance 4

True Genus

True tree

d1  d2  d3  d4

1.00

**Order** ×100

**Family** ×100

**Subfamily** ×100

**Genus**

**Genus**

**Genus**

**Genus**

average d1

average d2

average d3

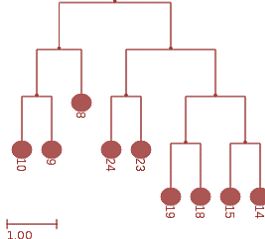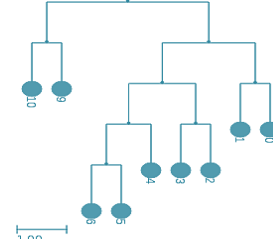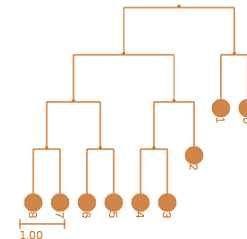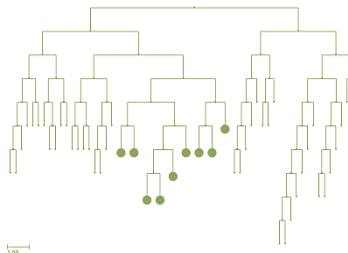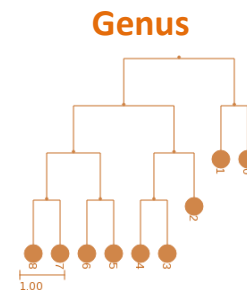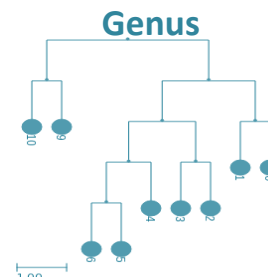average d4

**True tree**

**True Genus**

d1 > d2 > d3 > d4

# How can we solve the problem?

## Locus 1

```
ATGGCATATCCCATACACTAGGAT--CAA-GT-A
AT---GCTTACCCATTTCAACTTGGCT-ACAAGAT
ATGGCCA-----ACTCCCAACTAGGCTTTCAAGTT
ATGGC-ATCC--CACACAATTAGGATTCCAAGAT
ATGGCCTACCCATTCCAACTTG—TCTAC-AGA--
```

→

## Locus 1

```
ATGGCATATCCCATACACTAGGAT--CAA-GT-A
AT---GCTTACCCATTTCAACTTGGCT-ACAAGAT
ATGGCCA-----ACTCCCAACTAGGCTTTCAAGTT
ATGGC-ATCC--CACACAATTAGGATTCCAAGAT
ATGGCCTACCCATTCCAACTTG—TCTAC-AGA-
```

## Locus 2

```
TATCCCGGG - - -AAGCTAATT -AATGGAGGAATTTCAAGTATATT
TAT - - - -GGGCA - - - - - - - - - -ATGGA - - AATTTCAA - - -TAT
TATCCCTGG - - -AAAGCTAAT -CAATGGAGGAATTTCAAGTATAT
TAA - - - -GGGCA - - - - - - - - - -ATGGA - -AATTGCAA - - -TAT
TAT - - - - GCGCA - - - - - - - - - -ATGCA - - AATTTCAA - - -TAT
```

→

## Locus 2a                                    Locus 2b

```
TATCCCGGGAAGCTAATTAATGGAGGAATTTCAAGTATATT-----------------------------------------
TATCCCTGGAAAGCTAATCAATGGAGGAATTTCAAGTATAT-----------------------------------------
---------------------------------------------------TATGGGCAATGGAAATTTCAATATT
---------------------------------------------------TAAGGGCAATGGAAATTGCAATAT
---------------------------------------------------TATGCGCAATGCAAATTTCAATATT
```

## *Supermatrix:*

```
ATGGCATATCCCATACAACTAGGATTCCAAGTTATCCCGGGAAGCTAATTAATGGAGGAATTTCAAGTATATT-----------------------------------------
ATGGCTTACCCATTTCAACTTGGCTTACAAGATATCCCTGGAAAGCTAATCAATGGAGGAATTTCAAGTATAT-----------------------------------------
ATGGCCAACCACTCCCAACTAGGCTTTCAAGT---------------------------------------------------TATGGGCAATGGAAATTTCAATATT
ATGGCACATCCCACACAATTAGGATTCCAAGA---------------------------------------------------TAAGGGCAATGGAAATTGCAATAT
ATGGCCTACCCATTCCAACTTGGTCTACAAGA---------------------------------------------------TATGCGCAATGCAAATTTCAATATT
```