# Mega

## Phylogeny

## Approach

- **Reconstructing the "Tree of Life" is one of the enduring goals of evolutionary biology.**



**Since 19th**

# Background

- The massive accumulation of sequence data should provide more accurate phylogenies with the hope to resolve the tree of life.
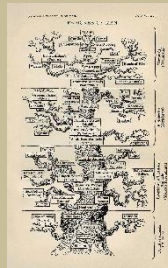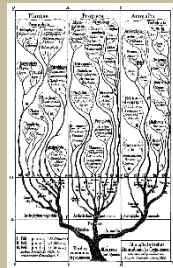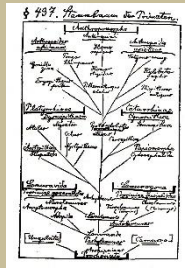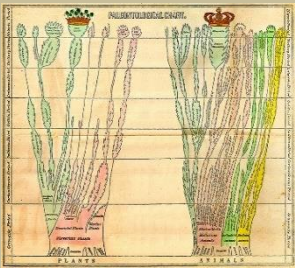- Recently, the Open Tree of Life was presented.



Synthesis of phylogeny and taxonomy into a comprehensive tree of life

# Current approaches for reconstructing mega phylogenies

- **Supertree:** compile many individual trees with partially overlapping taxa into a single large phylogeny.

- **Supermatrix:** gather sequence data for many loci into a single multiple sequence alignment, which is then used to reconstruct a large phylogeny.

+ twist

Species1:    ATGGCATATCCCATACAACTAGGATTCCAAGATGCAACC
Species2:    ATGGCACACCCAACGCAACTAGGTTTCAAGGACGCGGC
Species3:    ATGGCCAACCACTCCCAACTAGGCTTTCAAGACGCCTCC
Species4:    ATGGCACATGCAGCGCAAGTAGGTCTACAAGACGCTAC
Species5:    ATGGCACATGCAGCGCAAGTAGGTCTACAAGACGCTAC
Species6:    ATGGCACATGCAGCGCAAGTAGGTCTACAAGACGCTAC
Species7:    ATGGCTTACCCATTTCAACTTGGCTTACAAGACGCTACC
Species8:    ATGGCATACCCCCTACAAATAGGCCTACAAGATGCAAC

The reconstructed trees can only be as good as their underlying multiple sequence alignment (garbage in -> garbage out).

*Slowly evolving regions:*    to align –

good for solving          divergence events,
cannot be used to solve          events.


*Rapidly evolving regions:*    to align –

good for solving          event, cannot be used
to solve          divergence events.

Designed by **Vecteezy**

Species1:     AT-----CATA-CCATA-----A—GATTC-AG-AT--GC-AA-CC
Species2:     ATGG-CAC—CCCAAC------CTAGGTT----GGA-CGC-GGC-
Species3:     ATGG-----AA—CCATC—CAA---AGGC---------------G------
Species4:     --ATG—--AT—GC--GCG—AGTA—GG-----A—G----CTAC
Species5:     ATGGC-------A----GC00-----TAG--TCTACAAGAC-----TAC
Species6:     AT—G—CACATG-CA----CA-AAG------TCCA--GACGTAC-
Species7:     AG-----TACCC-------AACTG----GCTCAA--GA---GCTACC
Species8:     A-GGC-T-----CCC-----ACA-A—AGGC-CTACA-GATGCAAC-
Species9:     ---TGGCAAC-CC--CCTACAAAT---------TAC---AGA--CAAC
Species10:    A---GCAT--------CCCTA---TA----CCT----ACA------GC---C
Species11:    -----CATATC-----TACAACTAG----TTCA---------TGCA--CC
Species12:    ATGGC--------------C—GCAAC--TATTTCA-AGGACG---GC
Species13:    AT---CCA-CCA-CT-CA-CTAG-----CTTT-----GAC—GCC-CC
Species14:    ATG----GCCAT-----GCGC----GTAG------------------CTAC
Species15:    ---TGG--ACAT-AG--GC—AGTAG-GTC--CAAGA—CGC--C
Species16:    A--GGCC------GCA-AGA--GGT--TA--CAAG---ACGCTAC--
Species17:    ---GGCT---AC—CCA---CAA-----GGCTTA---AG--GCTA-CC
Species18:    AT--GG--CATA---C---RAAAT—AG-GC-C—---GATG--AC-
Species19:    AT---ATA------ATA----ACTAG----TTCCAA---TGC-A----C--
Species20:    A—TGG—CAC--CC---GCAACTA---GTCA---GAC—CG-GC-
Species21:    AT---C-----CAA--CACT----CAACTA—T-T-CA-C-------TCC-
Species22:    ---GGC------ATGCA---------GTAG----TACAA--ACGCTAC
Species23:    --ATGGC--TGCA-----GCA-AGTA------CTAC-----CG--TAC-
Species24:    ---GA--C—AT-G—G-C----------GGTCTAA---AC--CTA—C-
Species25:    ATG—C----C-CAT-----TTAA---TCTT—ACAAG-C---TAC-C
Species26:    --TGG--CA-----CCC------A-AT—AGG-CCTACA-----TGCAC
Species27:    -------CATA--TCCC---------------TTCCA----ATG-—AACC-
Species28:    A--TGGC-ACAC---AA--GC---CTGGT---C-AGGAC—G-GGC
Species29:    ---G---CC--AC----CTCCC—AAC-A----CTTC-----ACGC-TCC
Species30:    A-T-GG-C---ATGCA--GC--GCAAGG----TAC--GA--GCTAC

...

Species9999:  ATGGCA-TAC—------ACAA—TAG-CACAA--GATGC-AAC-
Species10000:ATGG----ATATCCCATACAACTAGGATTCCAA----GCAACC

## The problem

**Aligning rapidly evolving loci for highly diverged taxa**

leads to **poor** alignments and inaccurate trees.

*The problem*

- **Almost all MSA algorithms build a phylogeny during the estimation procedures.**

- **Those phylogenies are often based on pairwise distances.**

- **Pairwise distances:**
  - **- raw pairwise distances (uncorrected).**
  - **- model-corrected (e.g. Jukes-Cantor).**

ATGGCATATCCCATACAACTAGG
ATGGCATACACCCAACGCAACTAGG
ATGGCACACCCATCAACGCAACTAGC
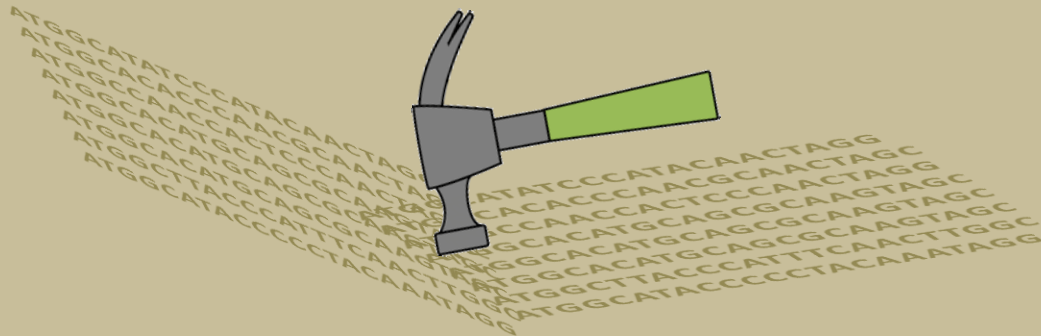ATGGCATCAACCACTCCCAACTAGG
ATGGCACACCGCAACTAGC

ATGGCATATCCCATACAACTAGG
ATGGCACACCCAACGCAACTAGC
ATGGCCAACCACTCCCAACTAGG
ATGGCACATGCAGCGCAAGTAGC
ATGGCACATGCAGCGCAAGTAGC
ATGGCACATGCAGCGCAAGTAGC
ATGGCTTACCCATTTCAACTTGGC
ATGGCATACCCCCTACAAATAGG

# Mega phylogeny approach

- **Separate sequences into subgroups of aligned sequences based on the degree of sequence saturation.**

## *Mega phylogeny approach*

- **If the most inclusive group of sequences is saturated, then the group is broken up into less inclusive groups using the next level in the taxonomic hierarchy.**

- **For example, if an "order" is found to be saturated, it would be broken into "families". Each smaller subset of sequences is then re-aligned and the saturation reassessed.**

- **This process continues iteratively.**

**Order**

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

**Saturated?**
**yes**

**Family1**
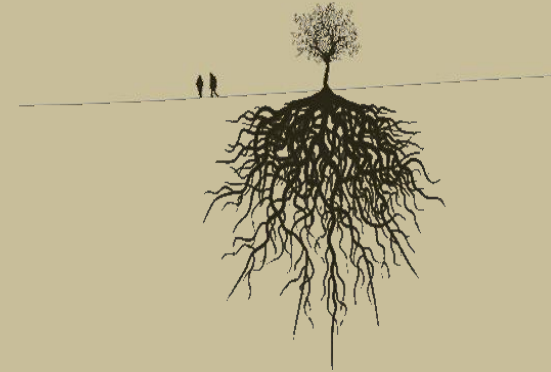
TATGGGCAATGGAAATTTCAATAT
TAAGGGCAATGGAAATTGCAATAT
TATGCGCA ATGCAAATTTCAATAT

**Family2**

TATCCCGGGAAGCTAATTAATGGAGGAATTTCAAGTATATT
TATCCCTGGAAAGCTAATCAATGGAGGAATTTCAAGTATAT

## *Detect saturation across a set of sequence data*

| d1 = | d2 = |
|---|---|
| raw pair-wise distance | Jukes-Cantor distance |

ACGGTCATGTATAC
TCTCTCCTGTGAAT
*d1 = 0.5*

ACGGTCATGTATAC
TCTCTCCTGTGAAT
*d2 = 0.8*

- **x = |d1 − d2|** *(x = 0.3)*

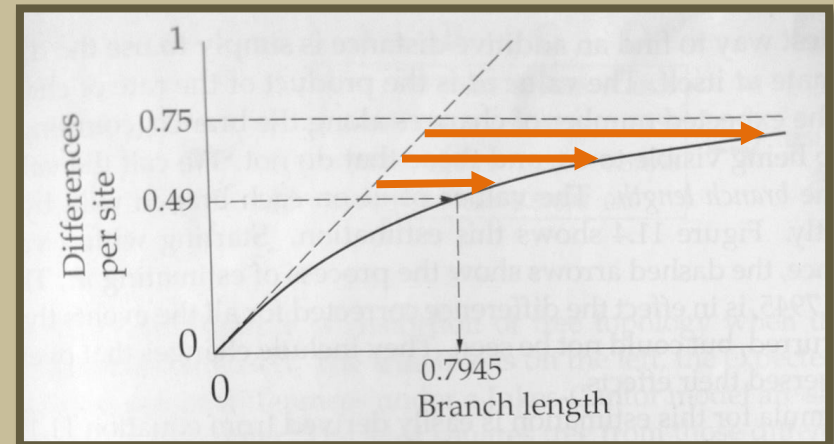- **Measure of dispersion is based on the median absolute deviation (MAD):**

**MAD = Med (|xi - Med (x)|)**

e.g.:
(9 ,6 ,4 ,**2** ,2 ,1 ,1)

**absolute deviations:**

(7 ,4 ,2 ,0 ,0 ,1 ,1)  => (7 ,4 ,2 ,**1** ,1 ,0 ,0)

*Detect saturation across a set of sequence data*

- The larger the MAD -> the larger the overall spread of distances.

- That is, above a certain value the assumed nucleotide substitution model is no longer adequately accounting for the rate variation exhibited by pair-wise distances among species.
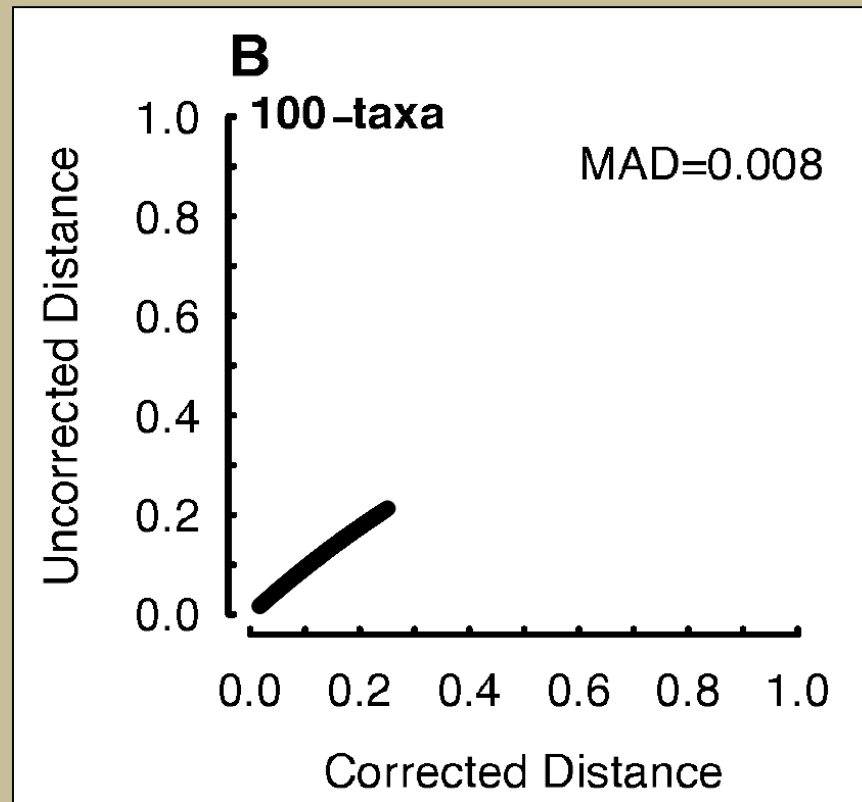
Jukes – Cantor model

WRONG!

**Large MAD    ---->**

$$Q = \begin{bmatrix} -3\alpha & \alpha & \alpha & \alpha \\ \alpha & -3\alpha & \alpha & \alpha \\ \alpha & \alpha & -3\alpha & \alpha \\ \alpha & \alpha & \alpha & -3\alpha \end{bmatrix}$$

## *Assessing saturation*
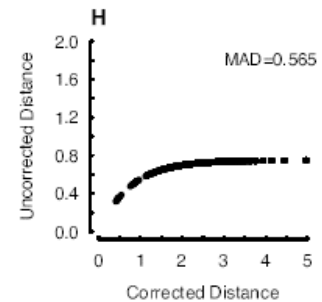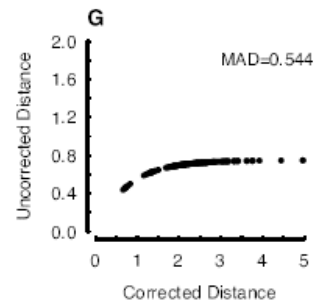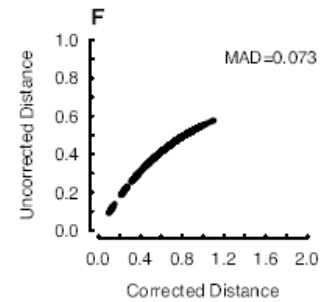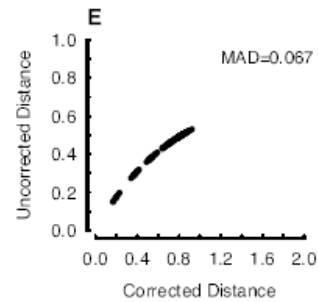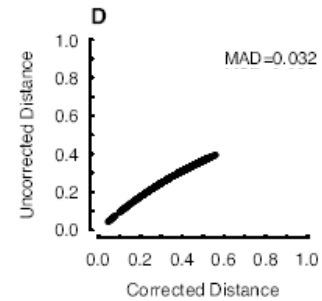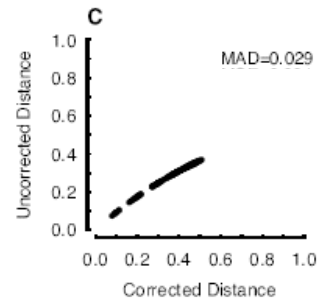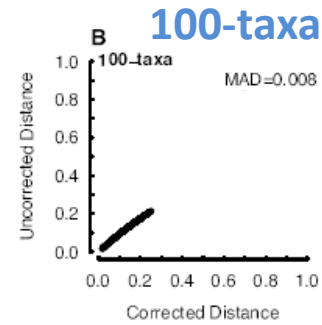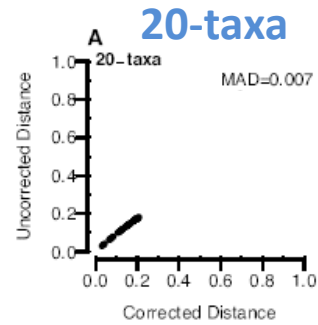
- **Determine a threshold for subdividing sequences into smaller alignments.**

**Assessing saturation**

**Tree length**

20-taxa

100-taxa

# One more slide:
# Our mega phylogeny approach

## Locus 1   ->

ATGGCATATCCCATACAACTAGGATTCCAAGT
ATGGCTTACCCATTTCAACTTGGCTTACAAGA
ATGGCCAACCACTCCCAACTAGGCTTTCAAGT
ATGGCACATCCCACACAATTAGGATTCCAAGA
ATGGCCTACCCATTCCAACTTGGTCTACAAGA

## Locus 1

ATGGCATATCCCATACAACTAGGATTCCAAGT
ATGGCTTACCCATTTCAACTTGGCTTACAAGA
ATGGCCAACCACTCCCAACTAGGCTTTCAAGT
ATGGCACATCCCACACAATTAGGATTCCAAGA
ATGGCCTACCCATTCCAACTTGGTCTACAAGA

## 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- ------------------------------------------
-------------------------------------------------------------------TATGGGCAATGGAAATTTCAATAT
-------------------------------------------------------------------TAAGGGCAATGGAAATTGCAATAT
-------------------------------------------------------------------TATGCGCAATGCAAATTTCAATATT

**Thanks.**