

big
trees



What's the problem with big trees?

Highly diverged taxa + rapidly evolving loci →
poor alignments →
inaccurate trees.

Species1: ATGGCATATCCCATTCAAGTGGATTCCAAGATGCAACC
Species2: ATGGCACACCTCAAGGCCAAGCTGGTTCAAGGACGCCGC
Species3: ATGGCCAACCACTCCGCTCGTGGCTTCAAGACGCCCTCC
Species4: ATGGCACATGCCCGCAAGCTGGTCTACAAGACGCTAC
Species5: ATGGCACATGCCCGCAAGCTGGTCTACAAGACGCTAC
Species6: ATGGCACATGCCCGCAAGCTGGTCTACAAGACGCTAC
Species7: ATGGCTTACCCATTCAAGTTGGCTTACAAGACGCTACC
Species8: ATGGCATACCCCTACAAATAGGCCTACAAGATGCAAC



Species1: AT-----CATA-CCATA-----A—GATTC-AG-AT--GC-AA-CC
Species2: ATGG-CAC—CCCACAC-----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: ---TGGCAACCC—GCTACAAAT-----TAC---AGA--CAAC
Species10: A---GCAT-----TACACGAG-----ACA-----GC---C
Species11: -----CATATC-----TACACGAG-----TGCA--CC
Species12: ATGGC-----GCAAGATGCA-----AGGACG---GC
Species13: A-----GCAACGCTACACGCTAC-----GA—GCC-CC
Species14: A-----GCAACGCTACACGCTAC-----CTAC
Species15: ---TGTAG—GC—AGTA—GCTCAA—GCA—CGC--C
Species16: A-----GCAAGATGCA-----GGT-----GCTACACGCTAC--
Species17: A-----GCAACGCTACACGCTAC-----GCTAA—GCTACACGCTAC
Species18: A-----GCAACGCTACACGCTAC-----TG—AC-
Species19: AT-----GATA—C—TAAAT—ACATA—T—TCAA—TGC-A---C--
Species20: A—TGA—GAC—C—GCAACTA---GTCA—SAC—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—AAC-
Species28: A—TGGC—ACAC---AA—GC---CTGGT---C—AGGAC—G—GGC-
Species29: ---G—CC—AC---CTCCC—AAC-A---CTTC---ACGC-TCC

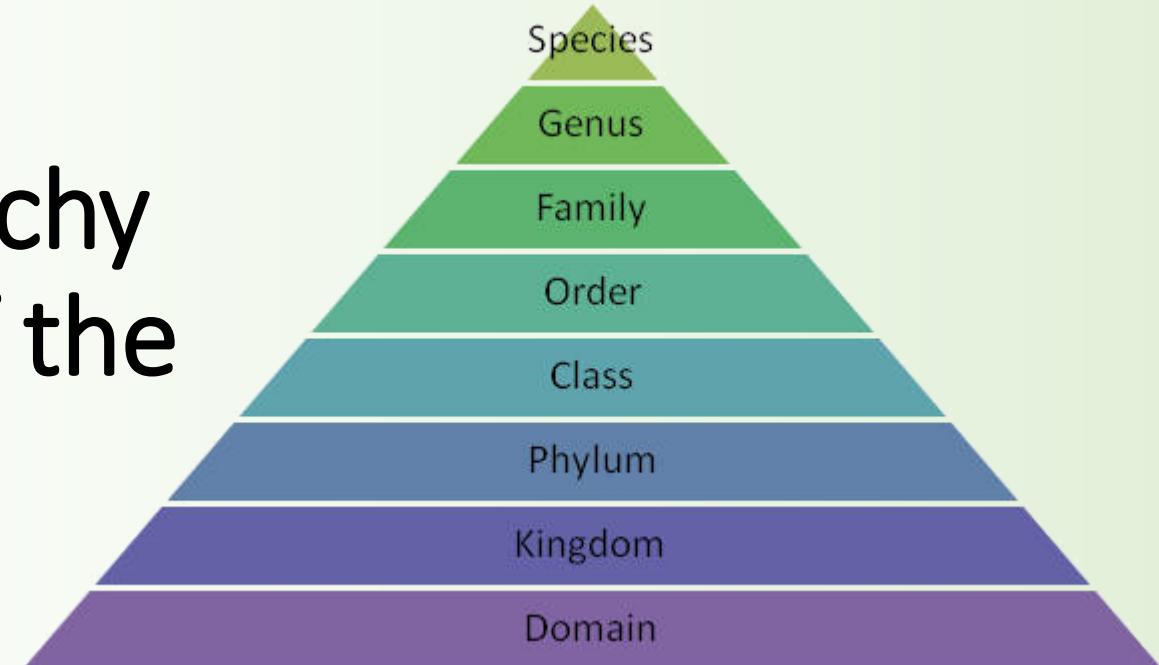


First approach:

Use the taxonomic hierarchy
to estimate the quality of the
alignment.



First approach:
Use the taxonomic hierarchy
to estimate the quality of the
alignment.



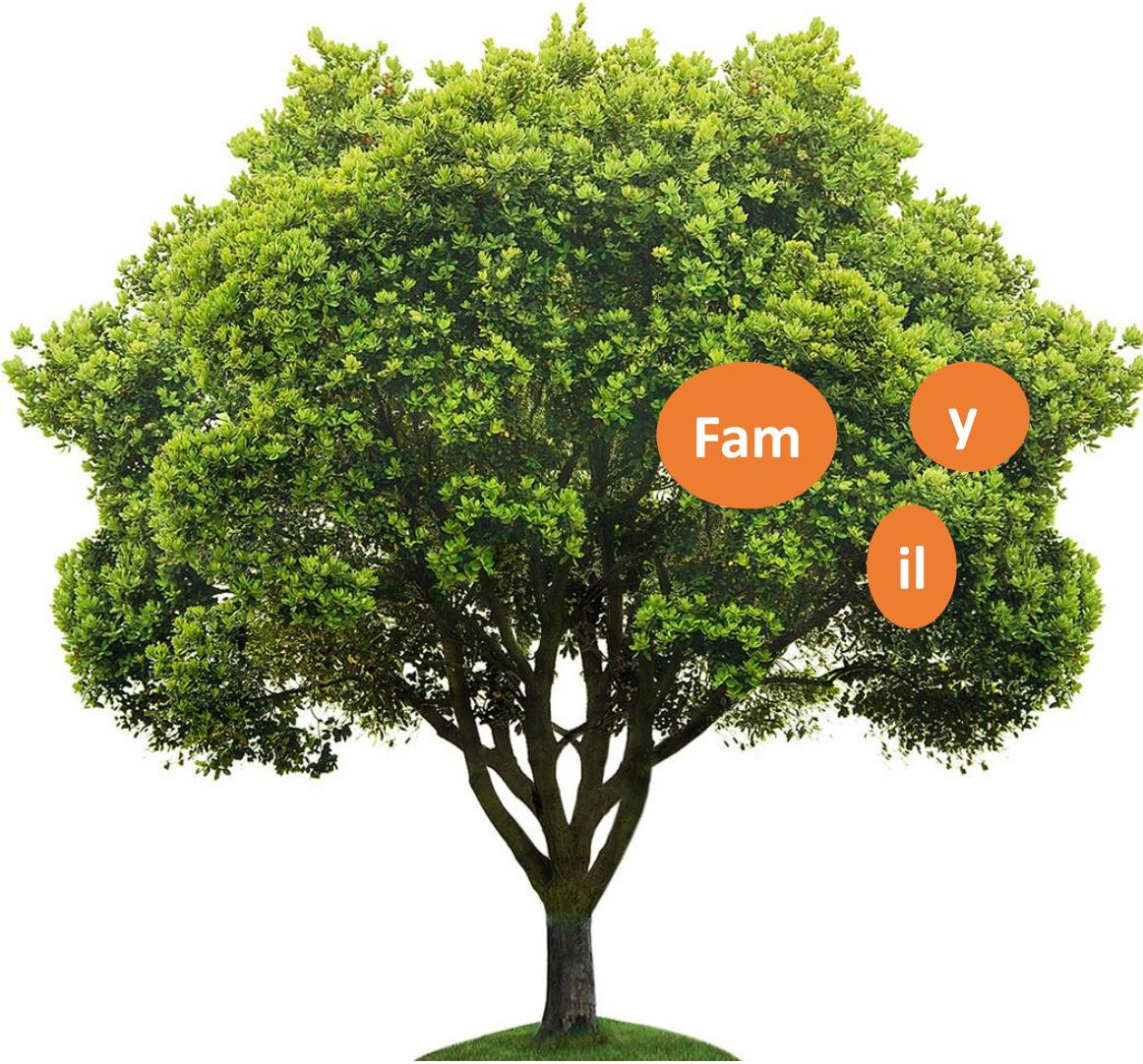


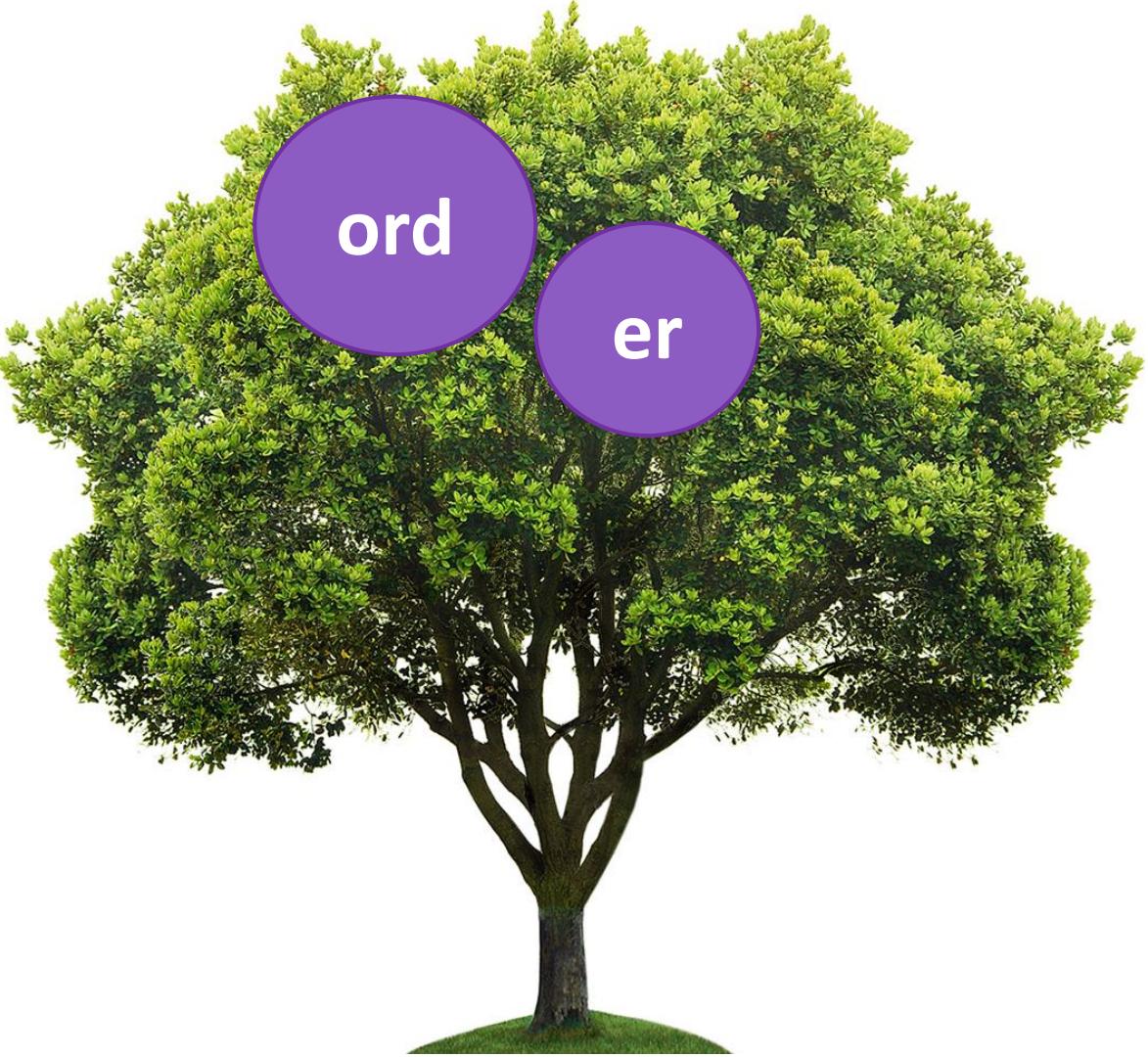


Family



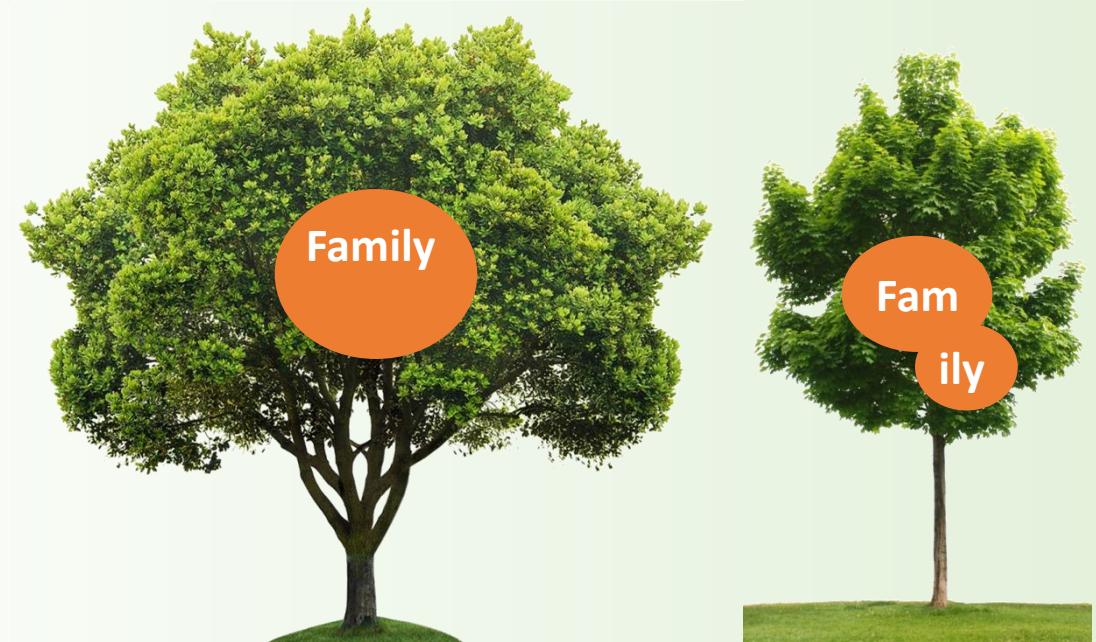








Results were a little unexpected
Go to plan B...





**Second approach:
Simulations.**

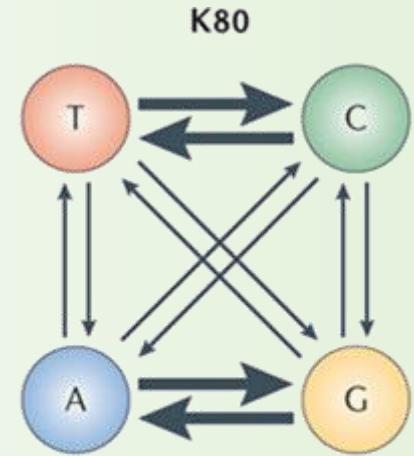


Second approach: Simulations.



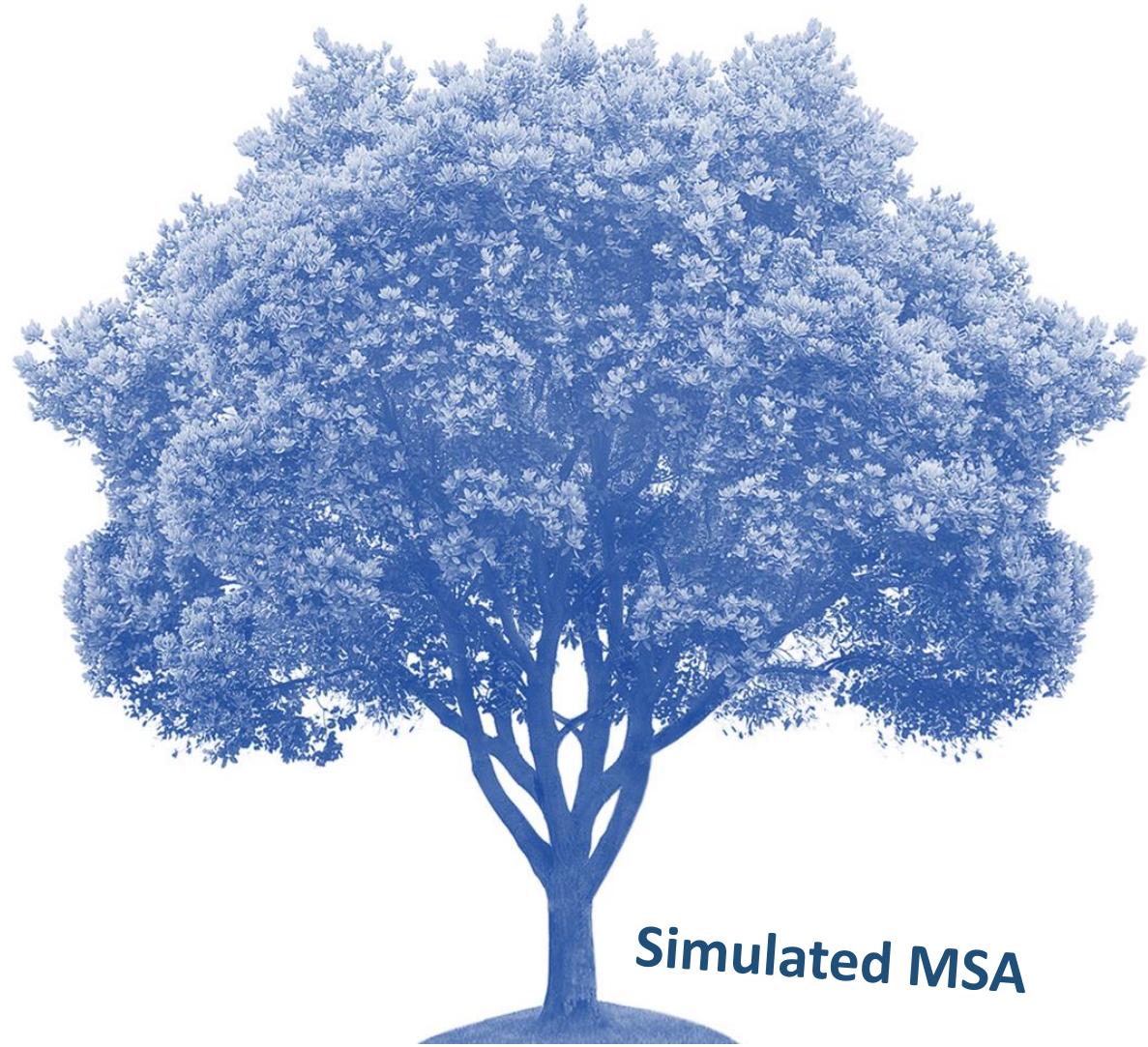
Characters of alignments:

- nucleotides frequencies
- substitutions model
- insertion-deletion parameters
(rate, length) - SpartaABC!



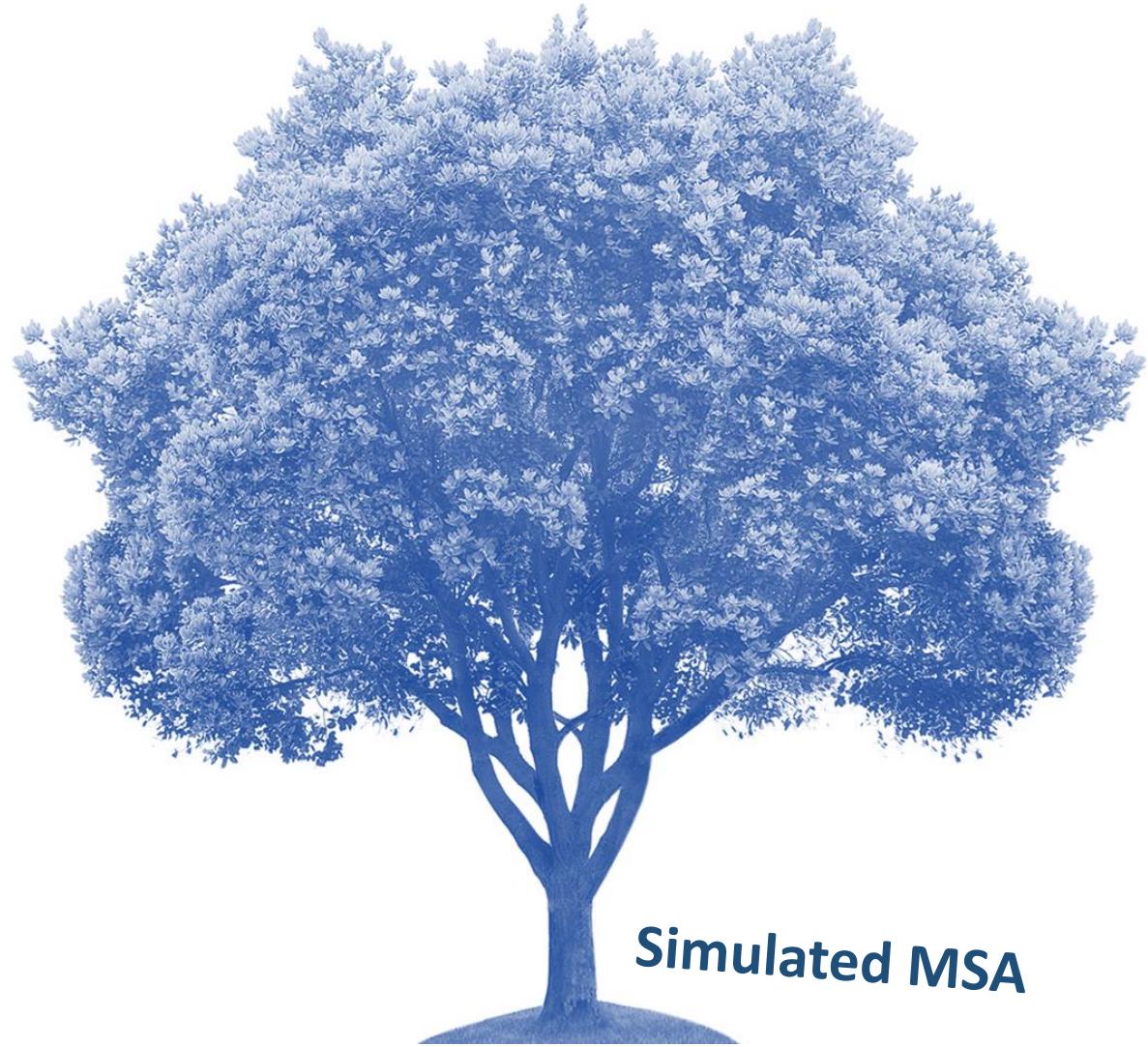


TRUE

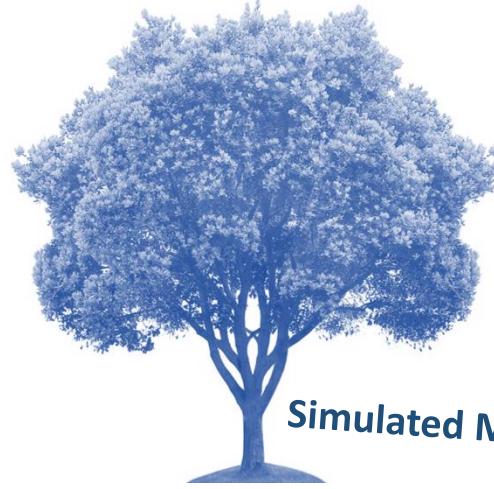
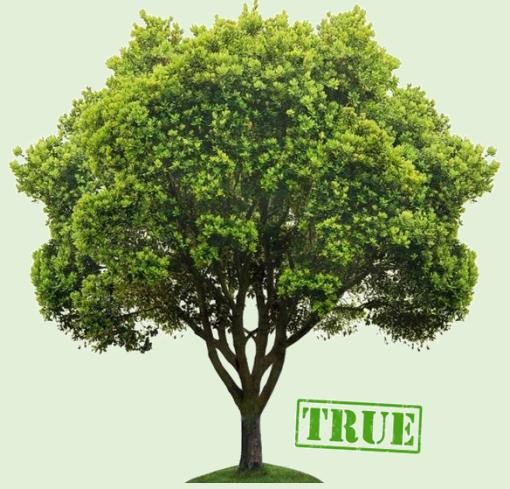


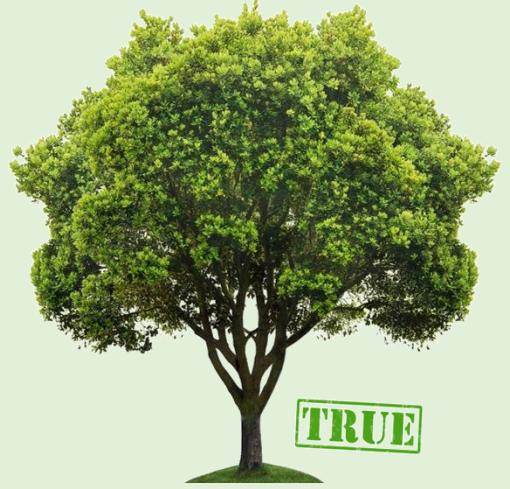


TRUE

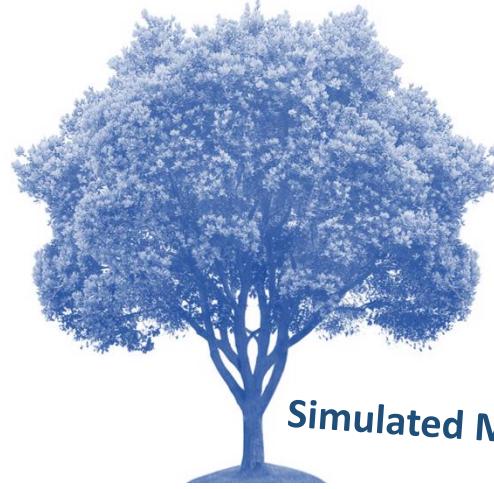


Simulated MSA

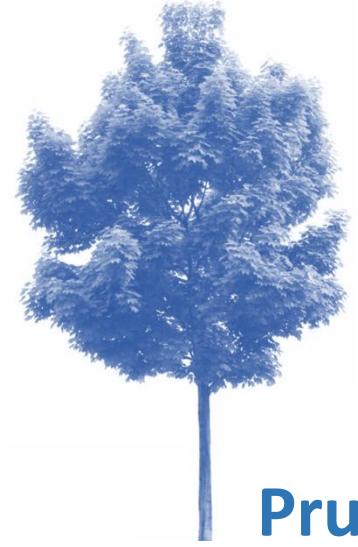




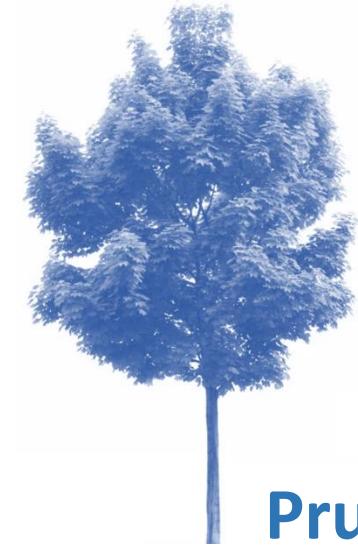
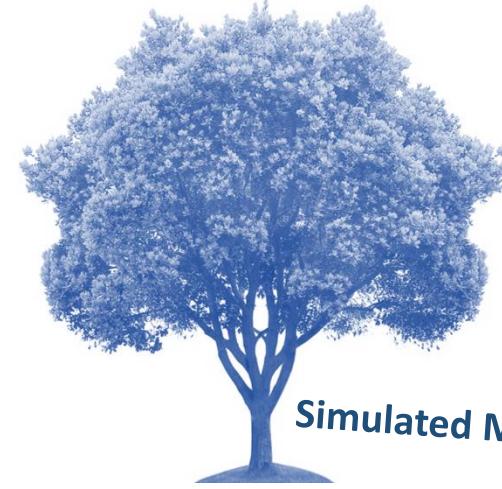
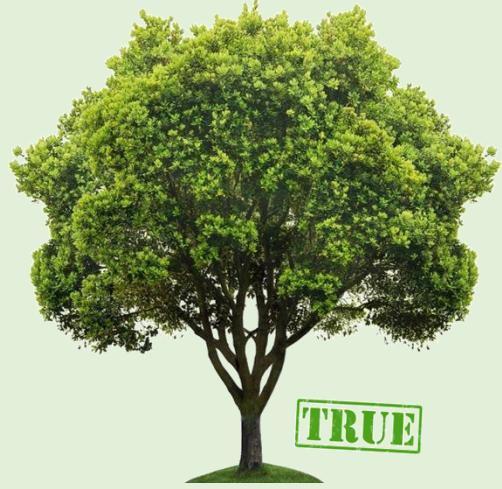
TRUE

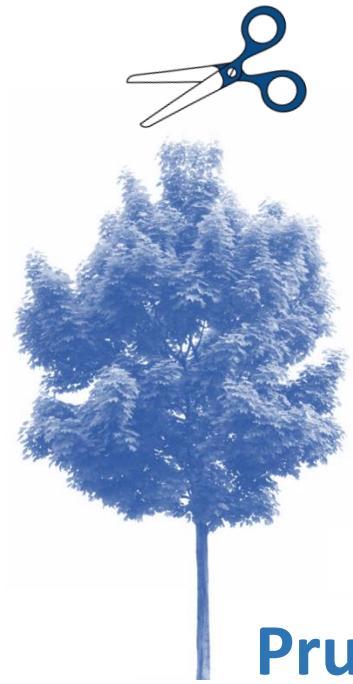
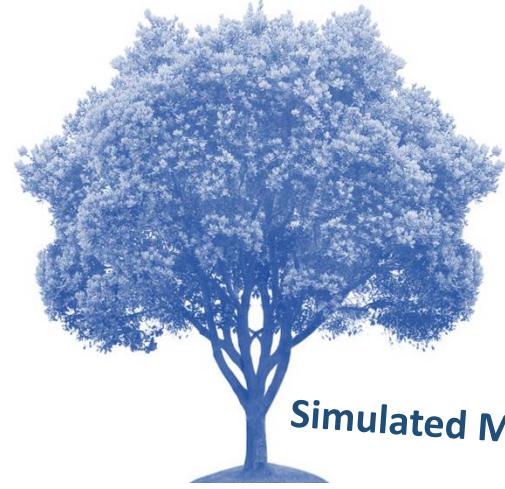
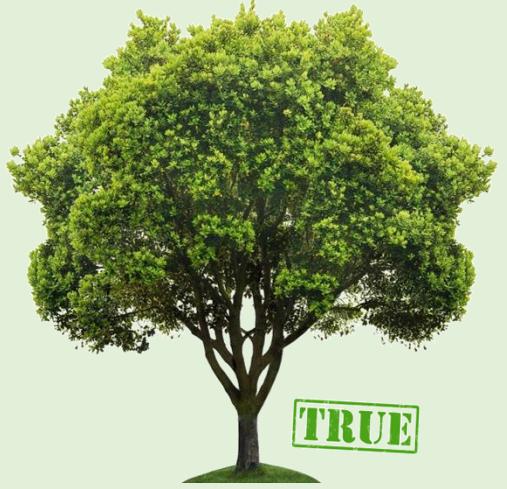


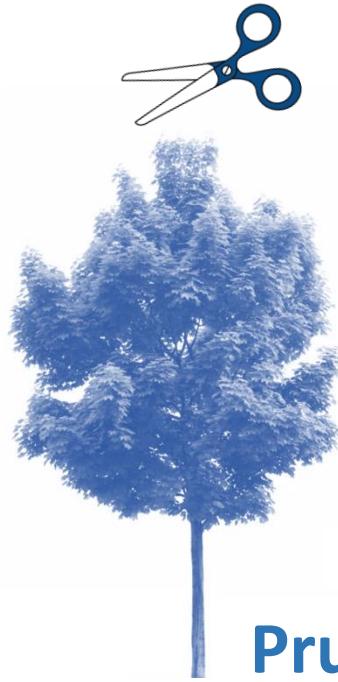
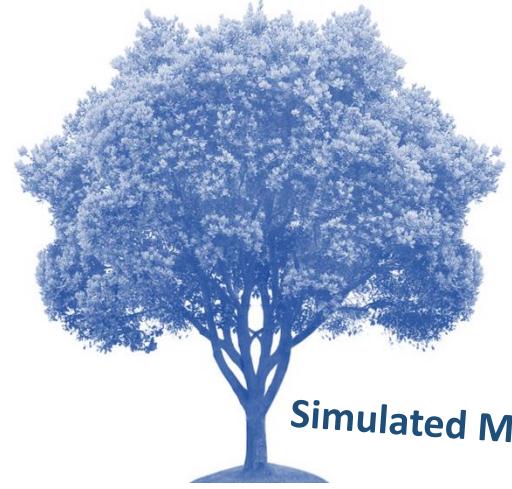
Simulated MSA

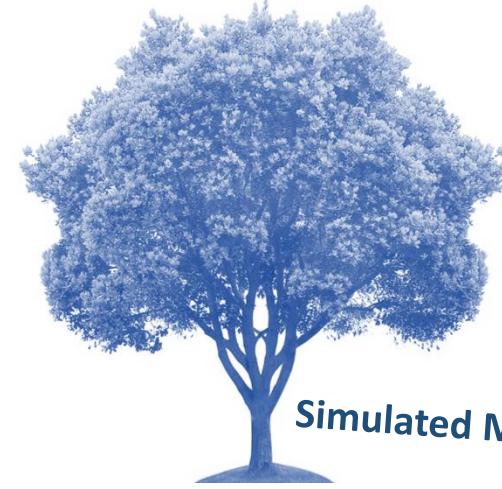
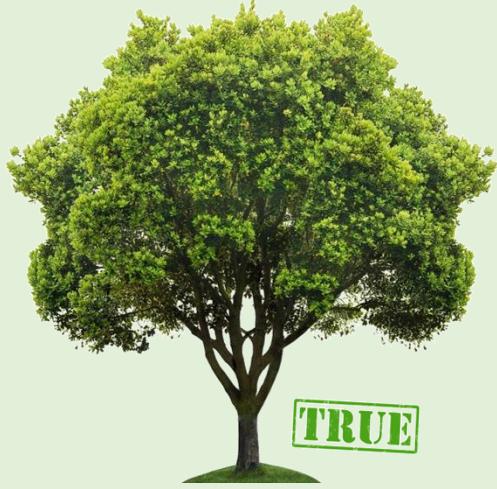


Pruned





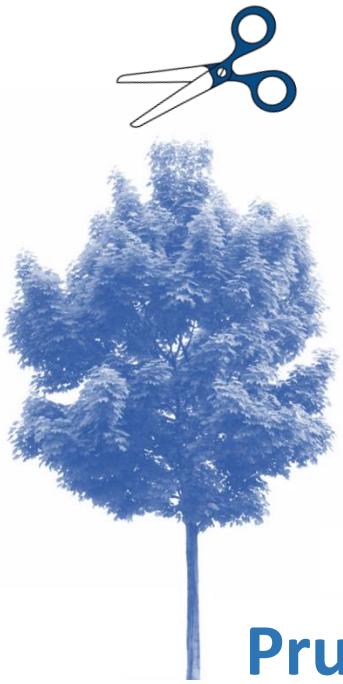


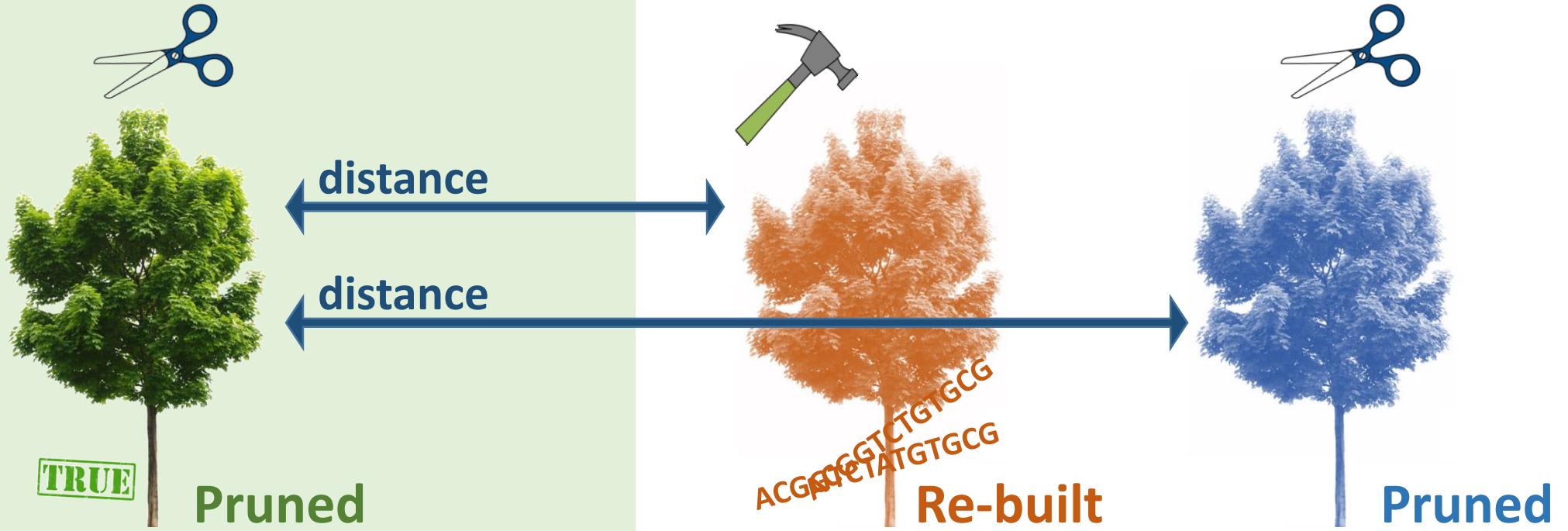


distance

distance

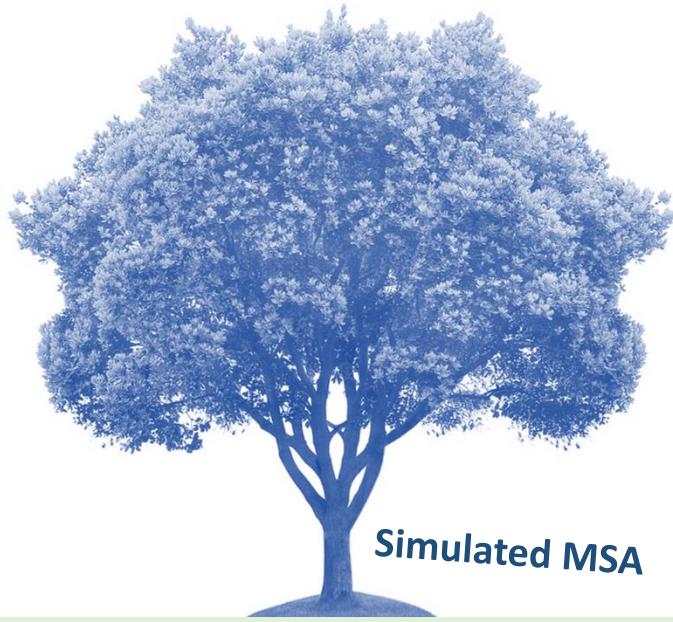
Two horizontal double-headed arrows, one above the other, indicating the distance between the "Pruned" tree on the left and the "Re-built" tree on the right.





Which of the distances we
expect to be larger?

~ 15,000 species



order



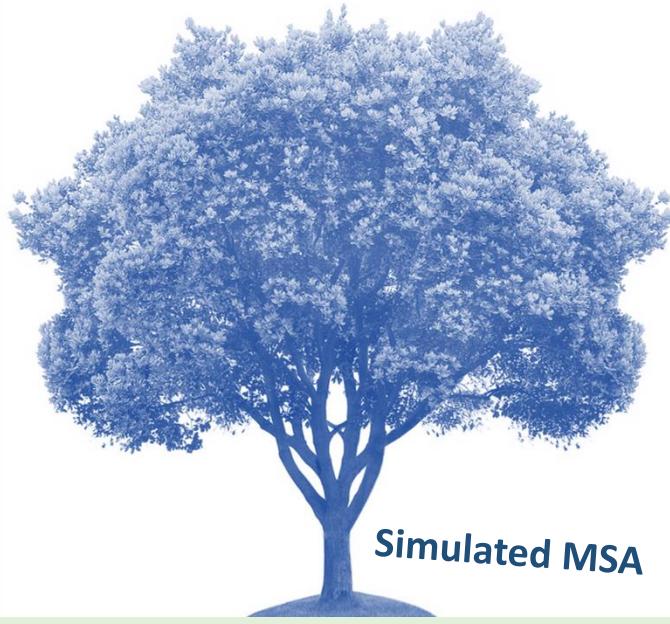
family



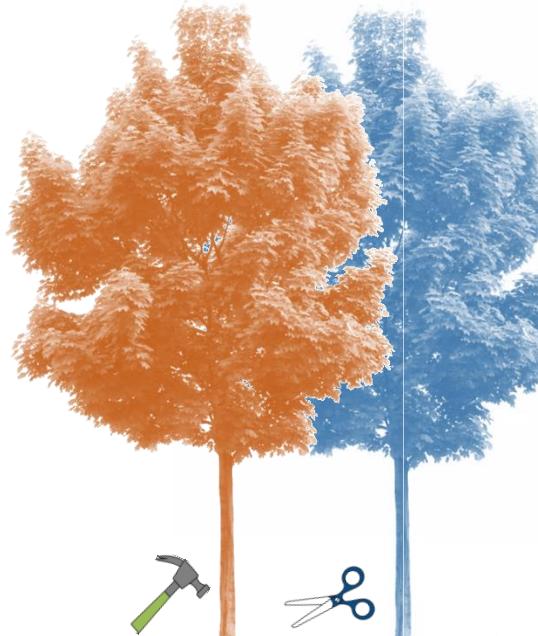
sub-
family



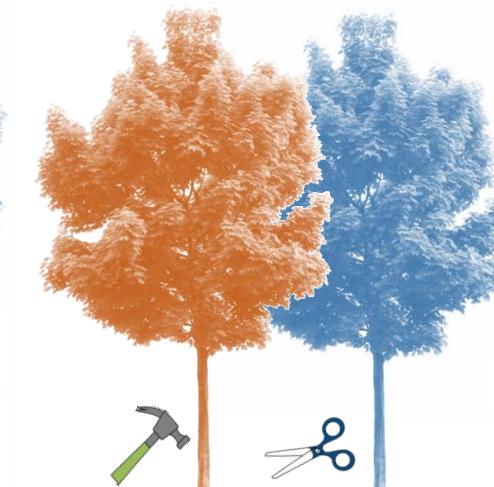
~ 15,000 species



order



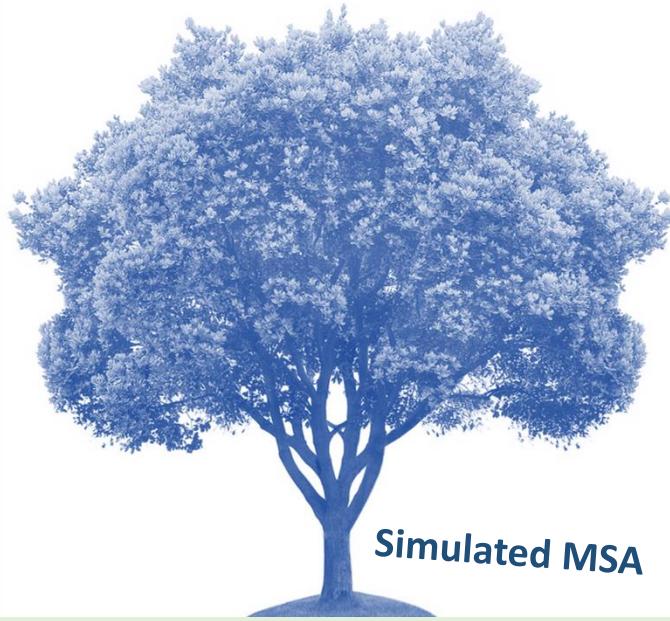
family



sub-
family



~ 15,000 species



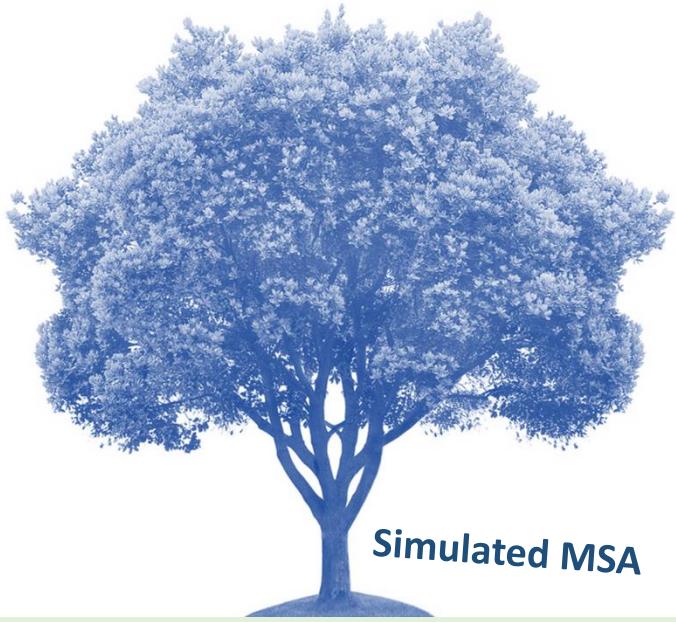
order



family



~ 15,000 species



~ 15,000 species



order



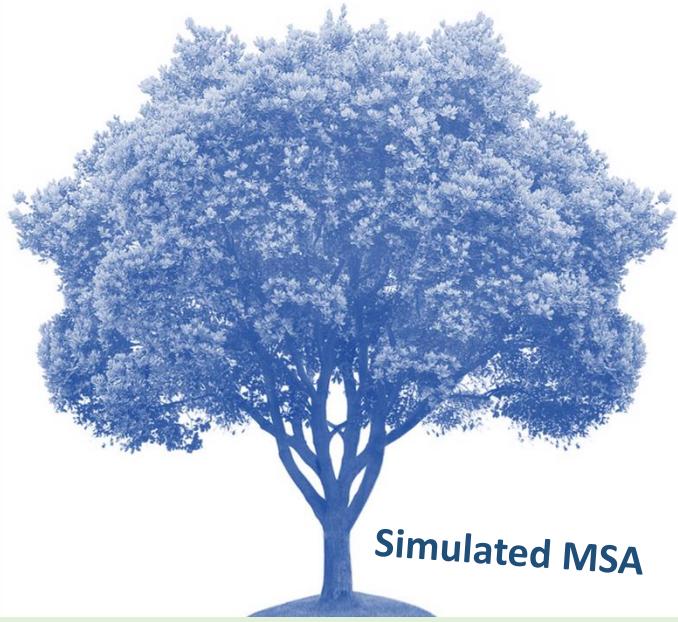
family



sub-
family



~ 15,000 species



~ 15,000 species



order



family



d1

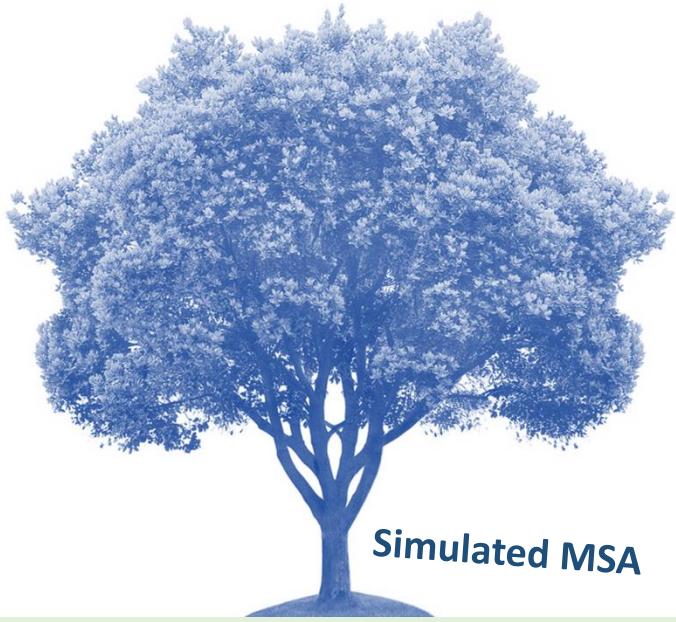
d2

d3

sub-
family



~ 15,000 species



~ 15,000 species



order



family



?

d1 d2 d3

d1

d2

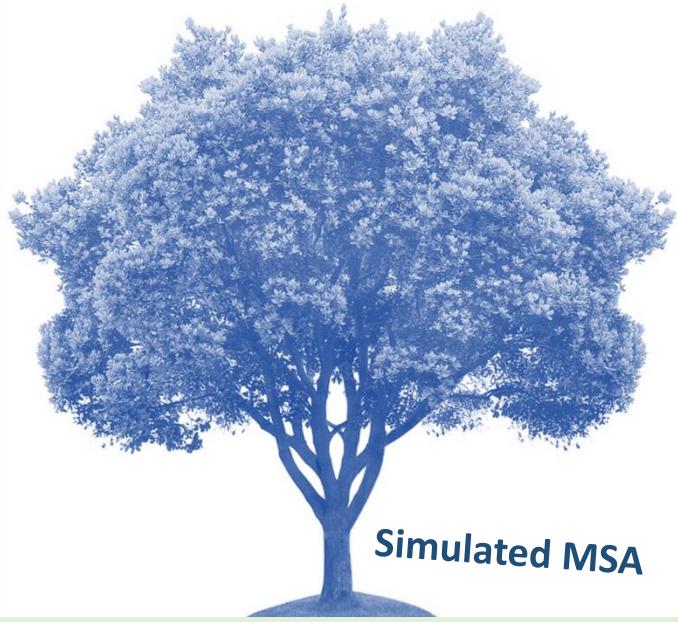
d3

sub-family



TRUE

$\sim 15,000$ species



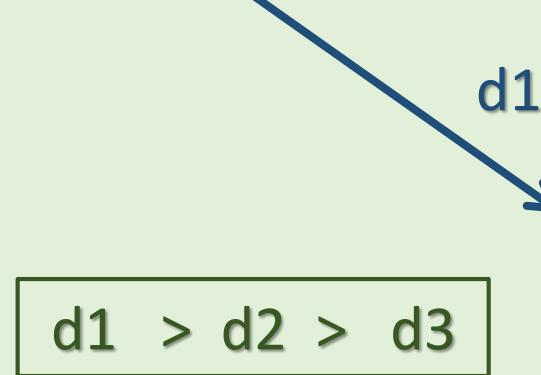
$\sim 15,000$ species



order



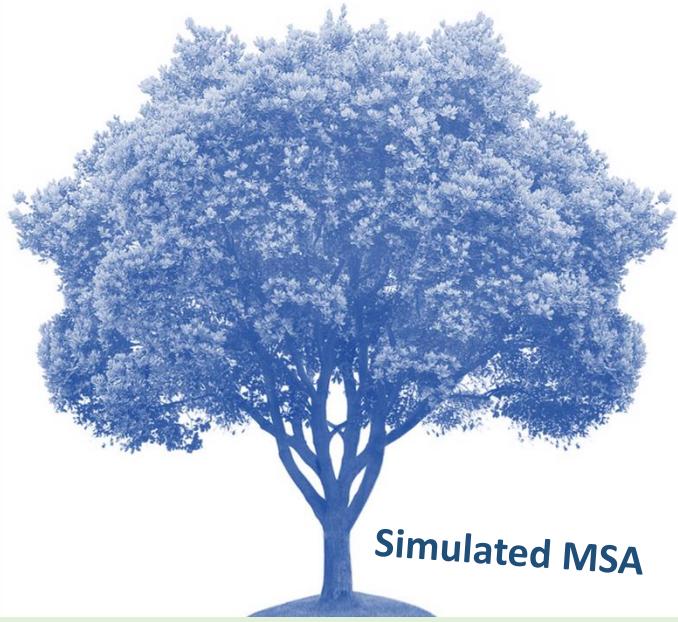
family



sub-
family

scissors
TRUE

~ 15,000 species



~ 15,000 species



order



family



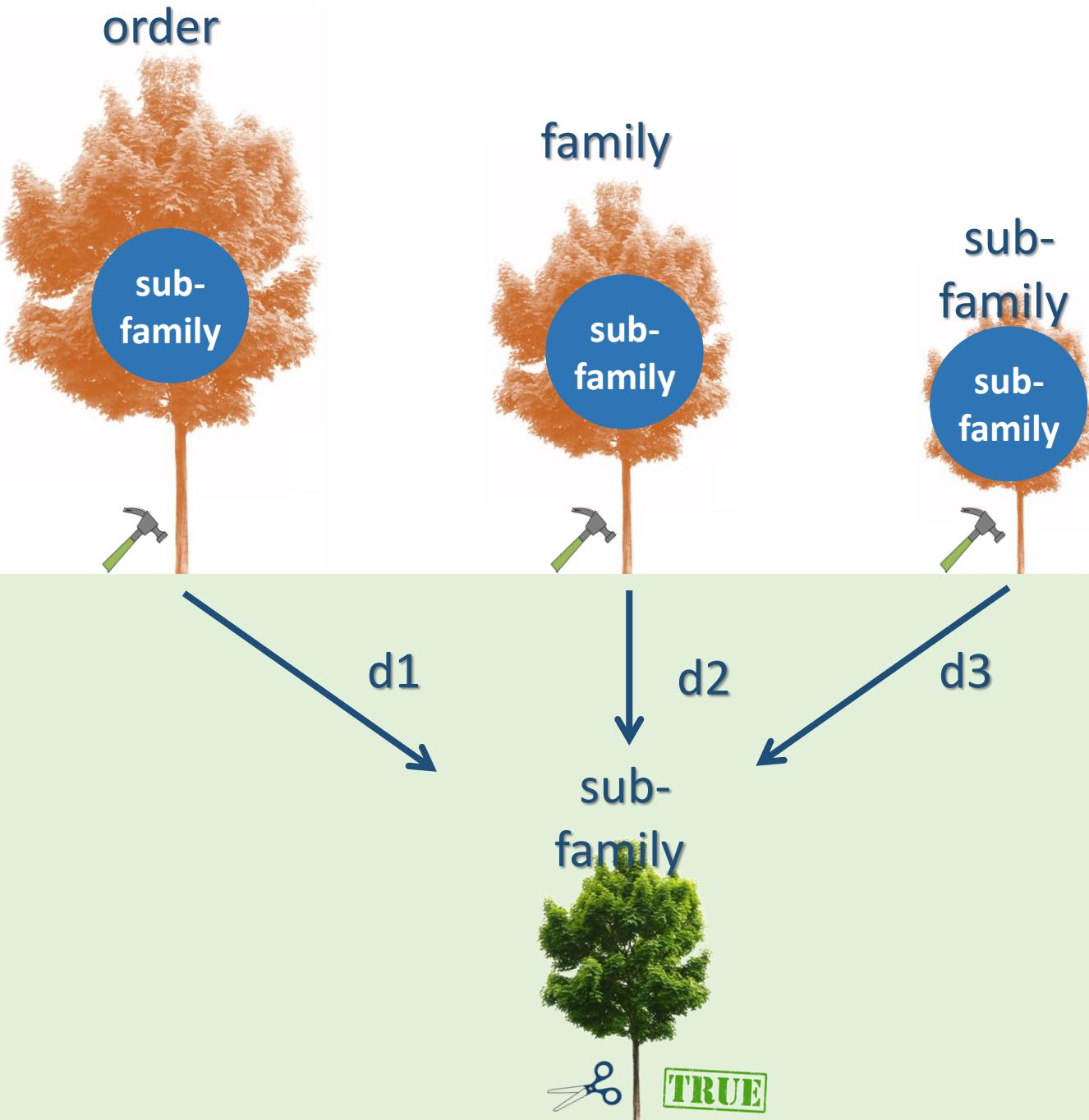
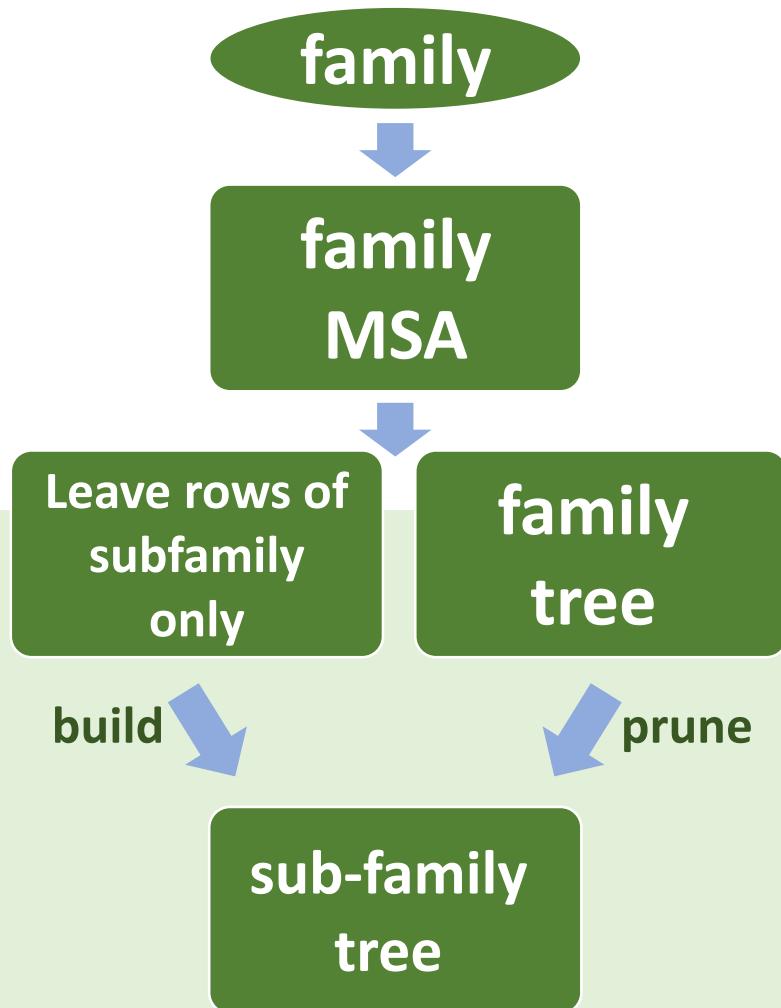
210

$d_1 > d_2 > d_3$

sub-family

scissors
TRUE

Two ways to create the
sub-family tree:





SpartaABC:

GENE	ORIGINAL_MSA_LENGTH	INDELIBLE_MSA_LENGTH
18S	1818	340
26S	3490	872
atpB	1509	779
ITS	1221	1253
matK	1996	1119
rbcL	1423	905
trnLtrnF	1316	1408
...		

A photograph of a large, leafy tree standing on a grassy hill. A dirt path leads up the hill towards the tree. The sky is filled with soft, white clouds.

Hope you could
see the forest...