# Introduction to Bioinformatics for Computer Scientists

## Lecture 6

# Plan for next lectures

- Today: Introduction to phylogenetics
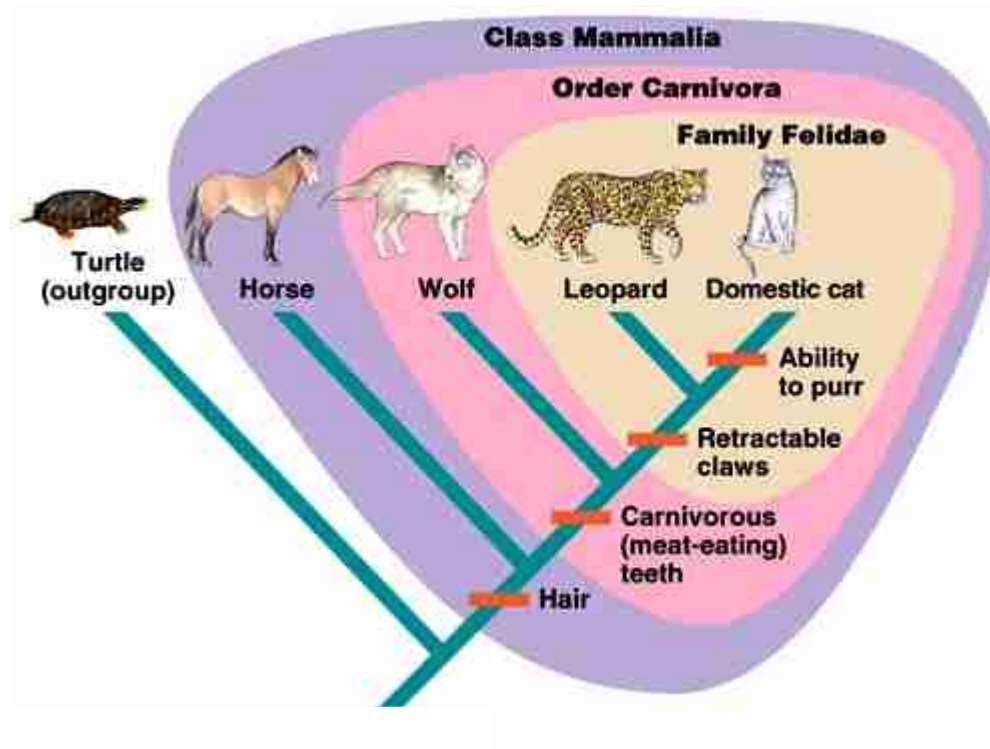- Lecture 7 (Alexis): Phylogenetic search algorithms

# The story so far

- Biological Terminology: RNA, DNA, genes, genomes, etc

- Pair-wise Sequence Alignment

- Sequence Comparison

- Genome Assembly

- Multiple Sequence Alignment

# The story so far
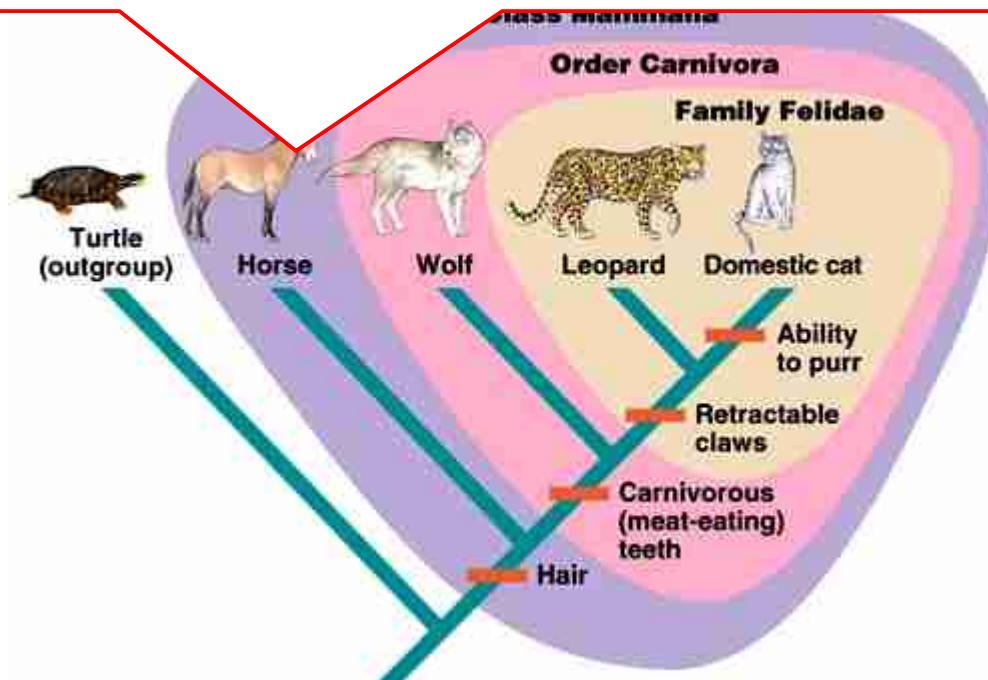
- Biological Terminology: RNA, DNA, genes, genomes, etc

- Pair-wise Sequence Alignment

- Sequence Comparison

- Genome Assembly

- Multiple Sequence Alignment

- Phylogenetic Inference

# A Taxonomy

# A Taxonomy

First systematic classification of living beings by Aristotele *384 -382* BC
Some terms still in use today, e.g., classification of animals into
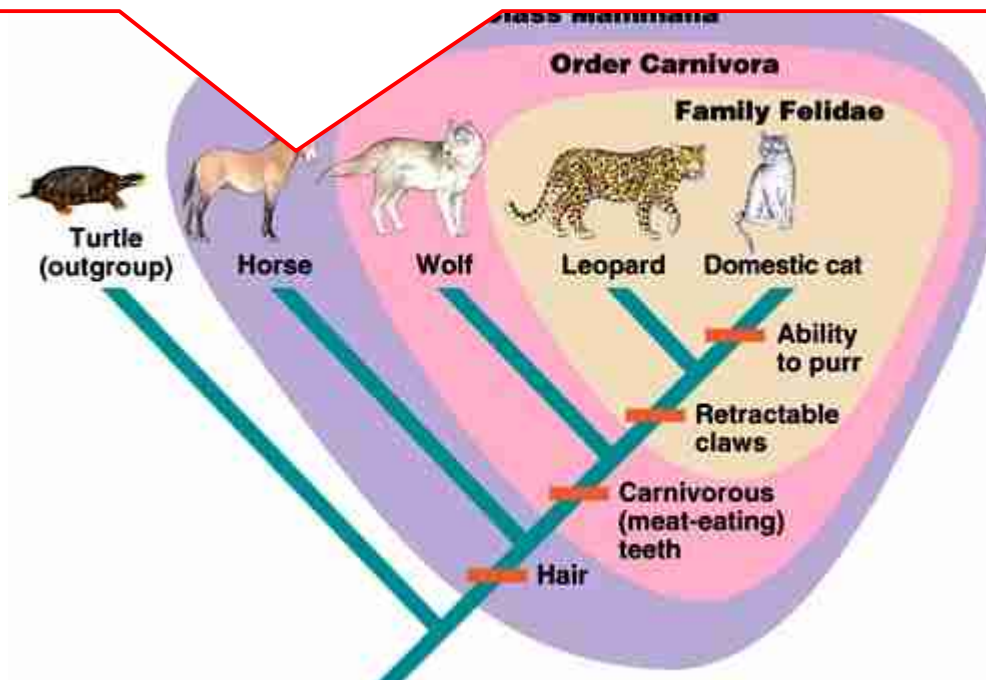*Vertebrates* versus *Invertebrates*

# A Taxonomy

First systematic classification of living beings by Aristotele *384 -382* BC
Some terms still in use today, e.g., classification of animals into
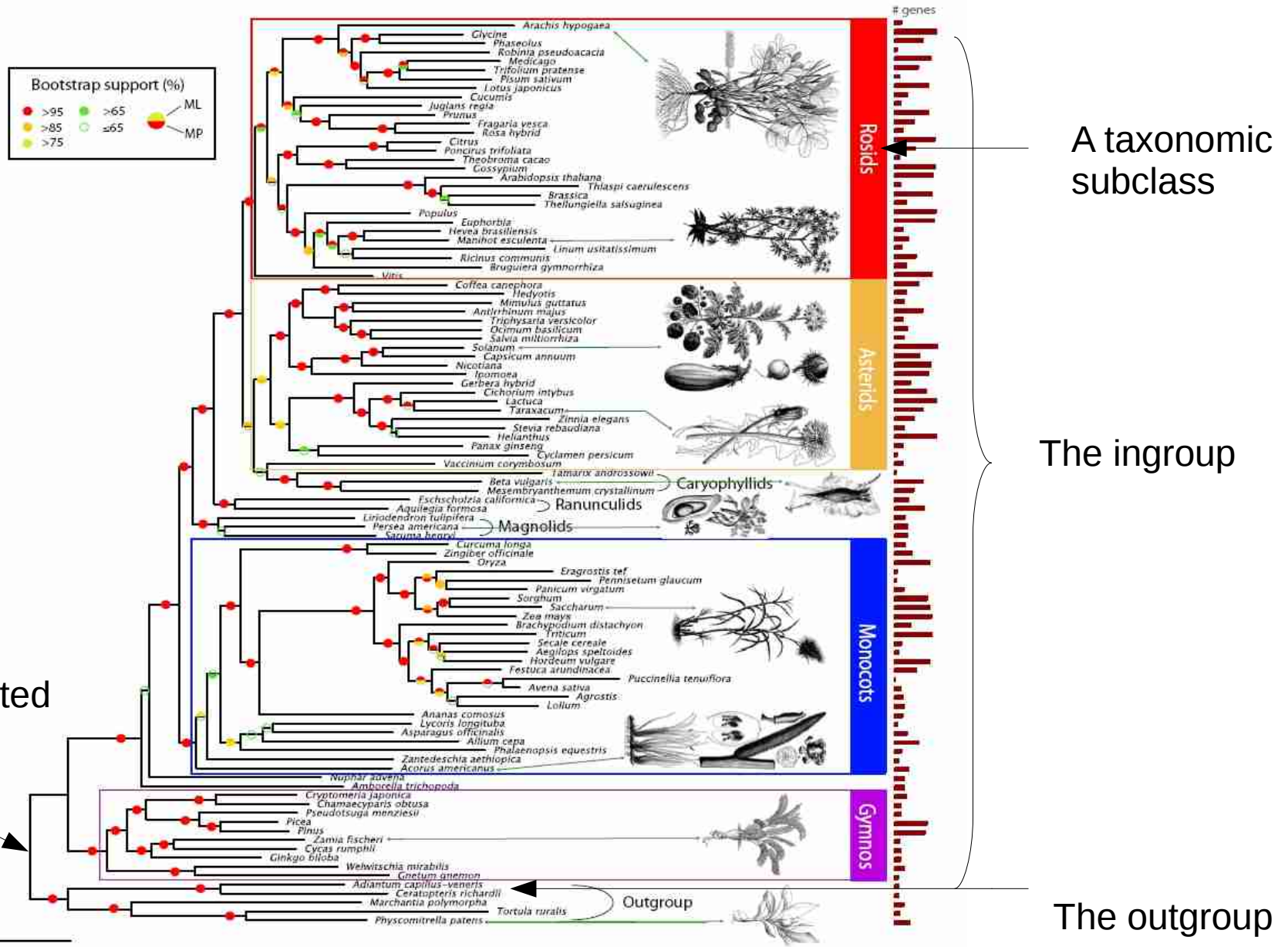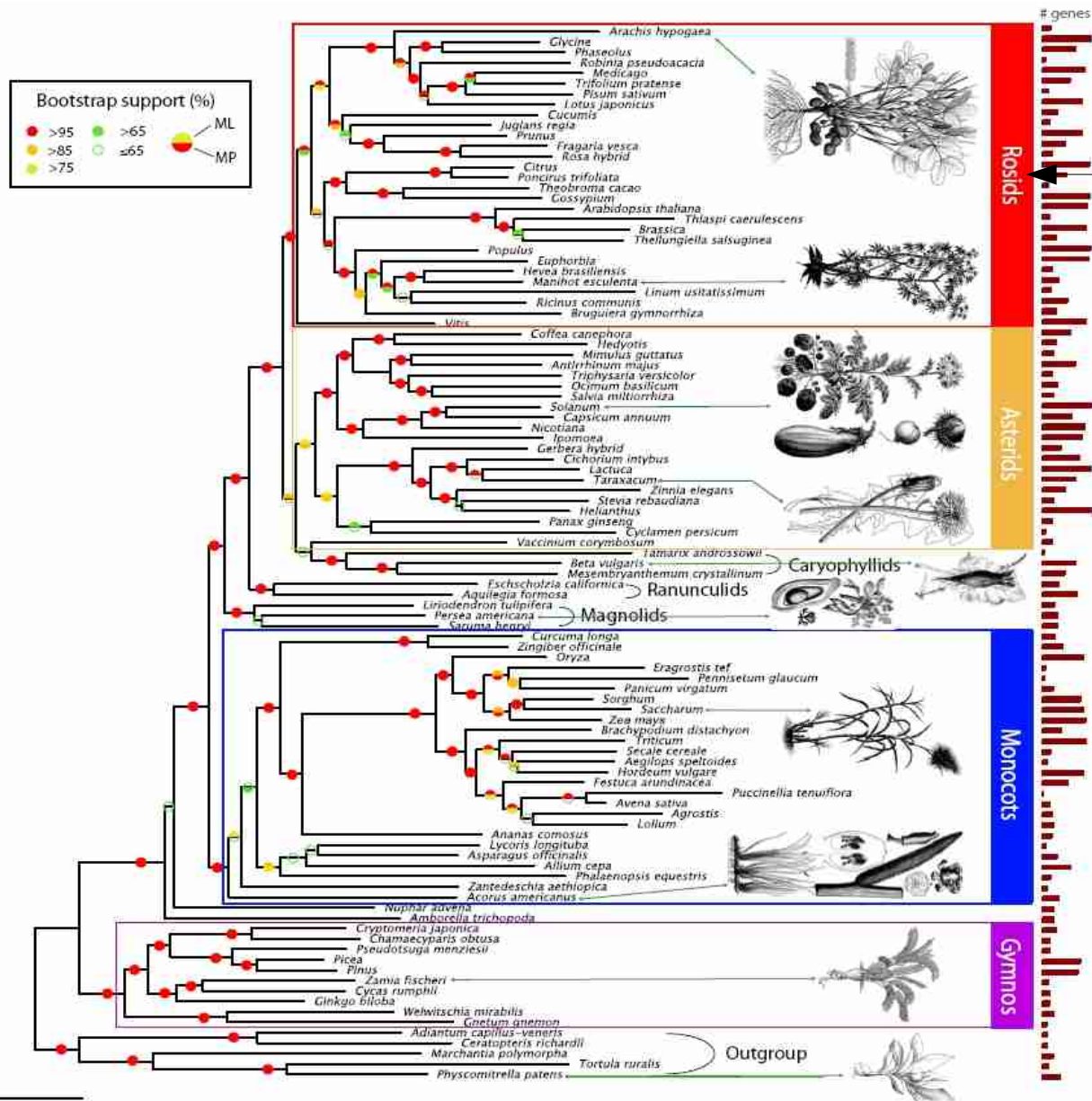*Vertebrates* versus *Invertebrates*

Wirbeltiere

# Taxonomy

- Group biological organisms (species) into groups with similar characteristics
- Define characteristics of groups at different hierarchy levels, e.g., animals > mammals > great apes
- Taxonomic ranks

  - Domain → three domains of life
  - Kingdom
  - Phylum
  - Class
  - Order
  - Family
  - Genus
  - Species

# A Phylogeny or Phylogenetic Tree



A taxonomic subclass

The ingroup

This tree is unrooted

The outgroup

9

# A Phylogeny or Phylogenetic Tree



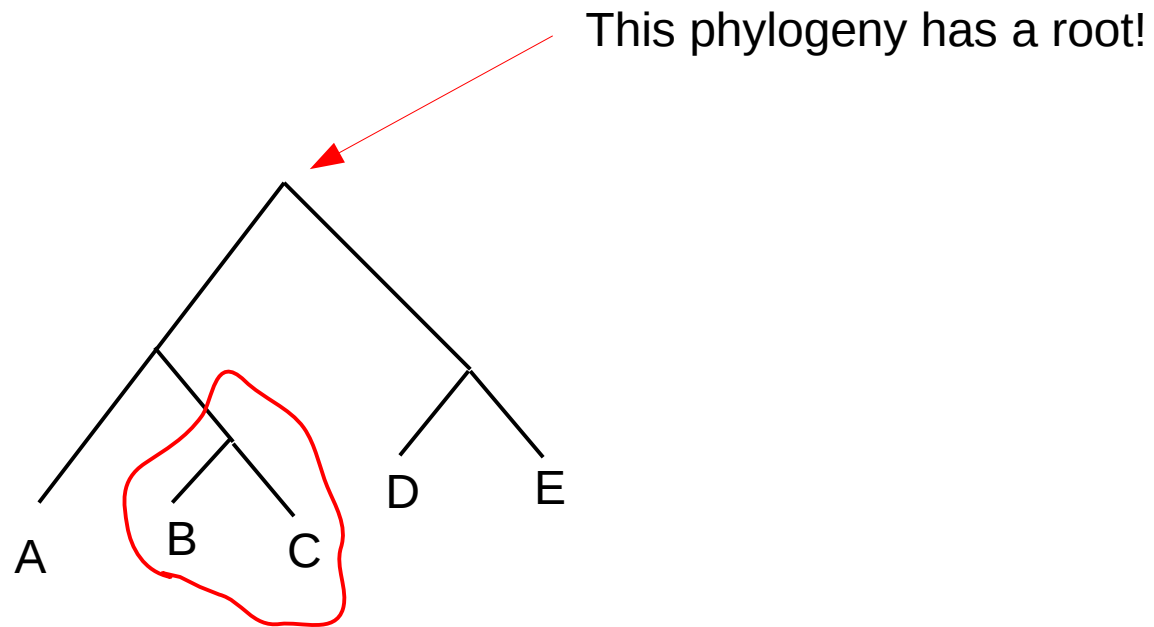In Phylogenetics such a subtree is often also called *Lineage*!

10

# Phylogeny

- An unrooted strictly binary tree

- Leafs are labeled by *extant* "übrig geblieben" (currently living) organisms represented by their DNA/Protein sequences

  → we can also sequence ancient DNA, see, for instance, the neandertal genome: "The complete genome sequence of a Neanderthal from the Altai Mountains", *Nature* 2013

  → depends on temperature, time, and other environmental conditions

  → up to 300,000 years back, see

  http://www.pnas.org/content/110/39/15758.abstract

- Inner nodes represent *hypothetical common ancestors*

- *Outgroup*: one or more closely related, but different species → allows to root the tree
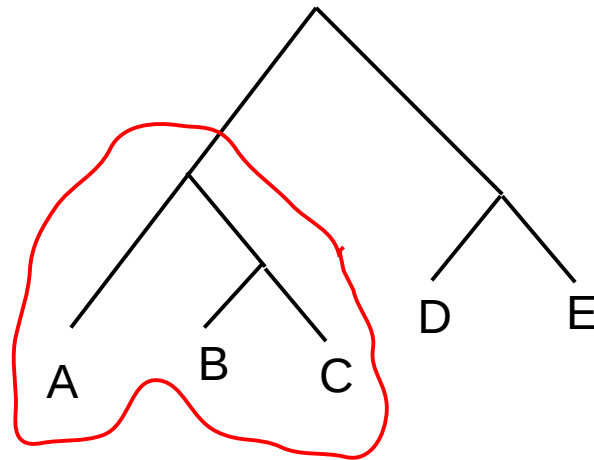
# Taxon

- Used to denote clades/subtrees in phylogenies or taxonomies

- A group of one or more species that form a biological unit

- As defined by taxonomists

    → subject of controversial debates

    → part of the culture/fuzziness of Biology

- In phylogenetics we often refer to a single leaf as taxon

    → the plural of taxon is *taxa*

    → we often say that a tree with *n* leaves (sequences) has *n* taxa

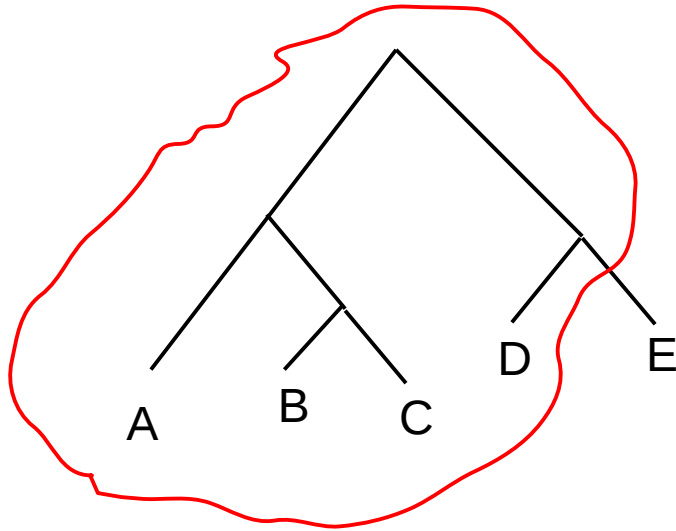# Some more terminology

This phylogeny has a root!



B and C are a *monophyletic* group; they are sister species

# Some more terminology



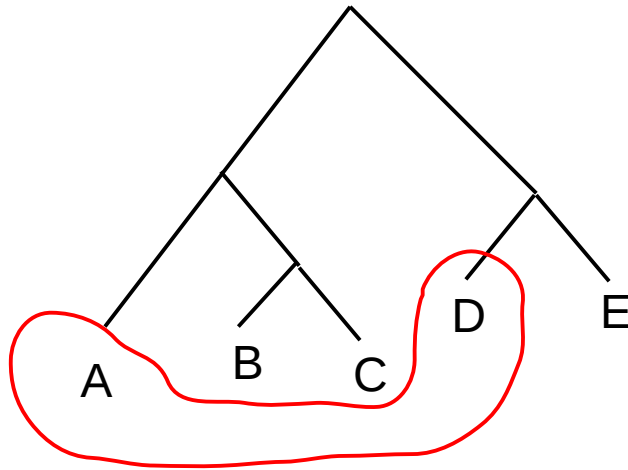**(A,B,C)** is a *monophyletic* group; it is sister to **(D, E)**

# Some more terminology



**(A,B,C,D)** is *paraphyletic* → **E** is excluded

# Some more terminology



**(A,D)** is a *polyphyletic* group → their most recent common ancestor (MRCA) is excluded

# Some more terminology



*Tree-based* or *patristic distance* between two taxa:
Sum over branch lengths along the path in the tree, e.g.:

# Some more terminology



*Tree-based* or *patristic distance* between two taxa:
Sum over branch lengths along the path in the tree, e.g.:
**A** ↔ **B**: 0.2

# Some more terminology



*Tree-based* or *patristic distance* between two taxa:
Sum over branch lengths along the path in the tree, e.g.:
**A** ↔ **B**: 0.2
**A** ↔ **D**: 0.35

# Tree Rooting

# Tree Rooting

# Tree Rooting



root

Outgroup species 1

Outgroup species 2

Ingroup species 2

Ingroup species 1

Ingroup species 3

# Tree Rooting

# Outgroup Choice

Ingroup species 4

Ingroup species 3

Ingroup species 2

Ingroup species 1

**?**

**Fuzzy signal**

Distant Outgroup

Ingroup species 4

Ingroup species 3

Ingroup species 2

Ingroup species 1

**Clear signal**

Close Outgroup

24

# Tree Inference

Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT

$\longrightarrow$

MSA
Program

$\longrightarrow$

Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT

$\longrightarrow$

Tree inference
program

Obtain *homologous* sequences from the same gene
(e.g., 16S RNA) of different species from a sequence database
(e.g., GenBank)

Taxon 1        Taxon 3

Taxon 2        Taxon 4

# Tree Inference

Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT

→

**MSA Program**

→

Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT

→

**Tree inference program**

↓

Most widely-used alignment formats:
- PHYLIP
- NEXUS
- FASTA

Taxon 1          Taxon 3

Taxon 2          Taxon 4

# Tree Inference

Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT

MSA Program

Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT

Tree inference program

Taxon 1        Taxon 3

Taxon 2        Taxon 4

Most widely-used tree formats:
•NEWICK
•NEXUS

# Tree Inference

Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT

MSA Program

Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT

Tree inference program

Taxon 1          Taxon 3

Taxon 2          Taxon 4

Newick example: Remember that this is an unrooted tree!
 (Taxon1, Taxon2, (Taxon3,Taxon4));
or
 ((Taxon1, Taxon2), Taxon3,Taxon4);

*Top level trifurcation*

# Tree Inference

Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT

MSA Program

Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT

Tree inference program

((Taxon1, Taxon2), (Taxon3,Taxon4));

Taxon 1

Taxon 3

root

Taxon 2

Taxon 4

# Tree Inference

Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT

→ MSA Program →

Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT

→ Tree inference program →

Taxon 1          Taxon 3
       0.1          0.15
          0.3
0.2              0.15

Taxon 2          Taxon 4

Trees may have *relative* branch lengths, depending on the tree inference method that was used

# Tree Inference

Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT

→

MSA
Program

→

Taxon 1:ACGTTT-
Taxon 2:ACGTT--
Taxon 3:ACCCT--
Taxon 4:AGGGTTT

→

Tree inference
program

Trees may have *relative*
branch lengths, depending
on the tree inference method
that was used

Taxon 1        Taxon 3
        0.1              0.15
            0.3
    0.2              0.15

Taxon 2        Taxon 4
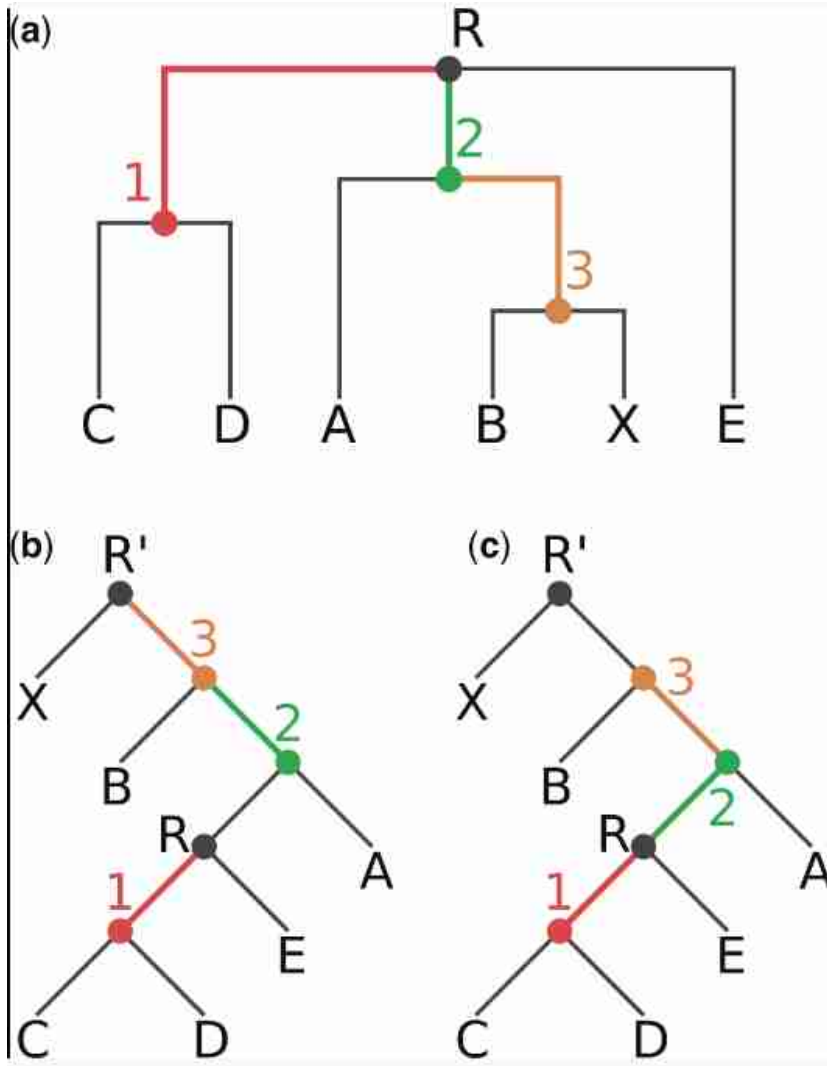
Newick format with branch lengths:
(Taxon1:0.1,Taxon2:0.2,(Taxon3:0.15,Taxon4:0.15):0.3);

# Problems with Newick tree format

- Except for branch length values: no way to associate meta-data to branch lengths

- However, there is important meta-data, e.g., branch support: how well is a branch in the tree supported?

  → ad hoc solution: represent branch support values as node meta-data!

  → this causes problems

# Problems with Newick tree format



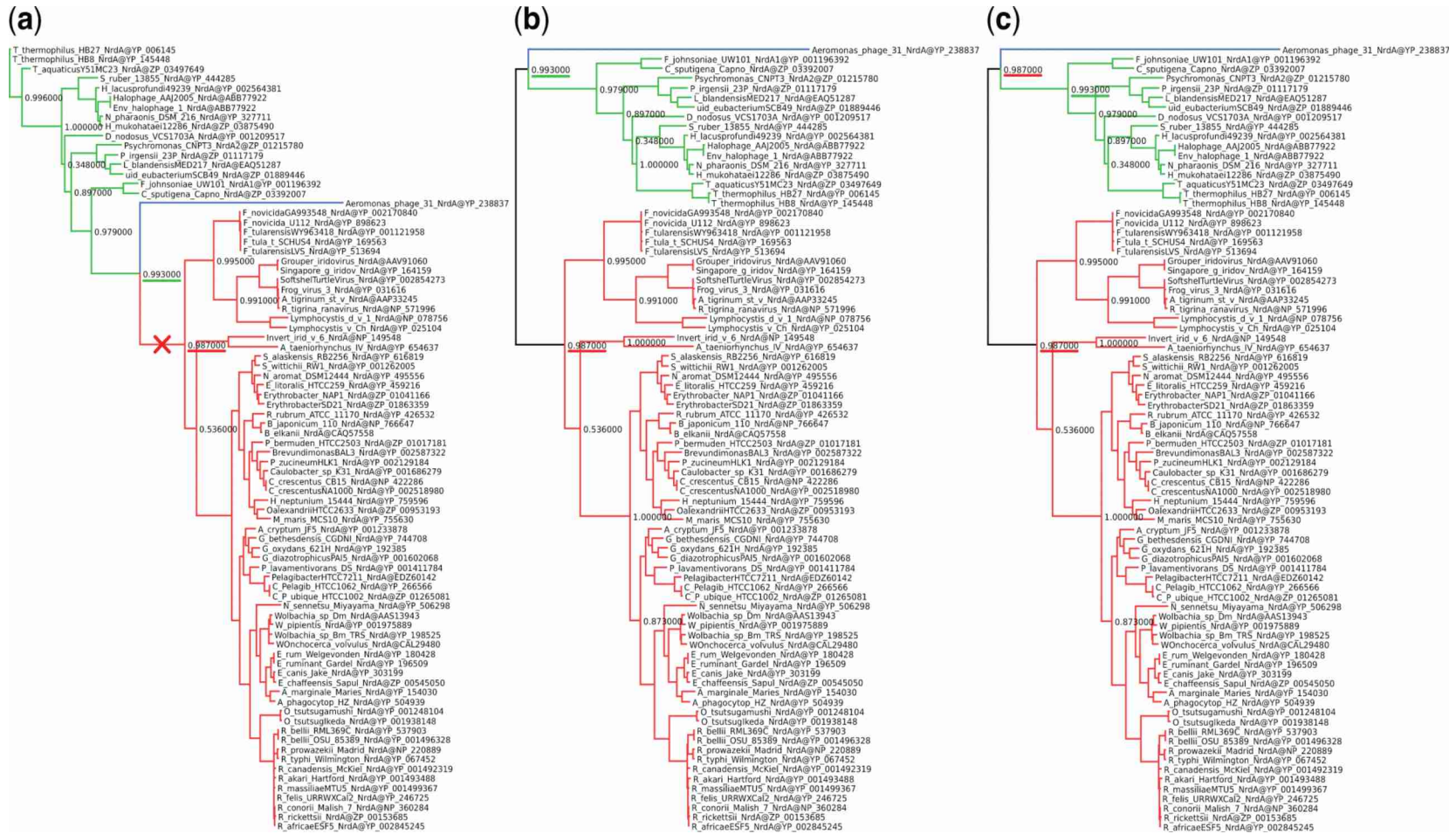Branch support values represented as node meta-data can be assigned incorrectly to branches after re-rooting.

About 50% of the tools we checked had this Problem. For details see:
https://academic.oup.com/mbe/article/34/6/1535/3077051
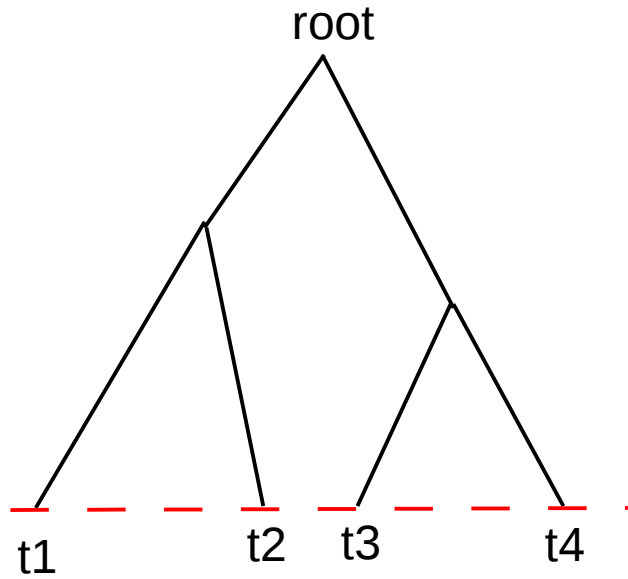
Which representation is correct?

# A real example



a) original tree
b) re-rooted tree with shifted support values
c) re-rooted tree with correct support values

# Tree Shapes



Evolutionary time →

root

t1    t2  t3    t4

**Ultrametric tree**

root

t2

t3

t1

t4

**Non-ultrametric tree**

# Tree Shapes

Evolutionary time

root

t1    t2   t3        t4

**Ultrametric tree**

root

t2

t3

t4

t1

**Non-ultrametric tree**

Most tree inference
models/algorithms/programs
produce non-ultrametric trees

# Tree Shapes

# Tree Shapes

# Dating Trees



**Ultrametric tree**

dated fossil

root

t1    t2    t3    t4

Evolutionary time

# Dating Trees

# Dating Trees



Root: 3 million years

1 million years

dated fossil 2 million years

2 million years

Evolutionary time

t1    t2    t3    t4

**Ultrametric tree**

# Dating Trees



Evolutionary time

Root: 3 million years

1 million years

dated fossil 2 million years

2 million years

**Ultrametric tree**

t1    t2   t3    t4

We need a rooted &
ultrametric tree!
→ rooting with outgroups
→ ultrametricity with programs
for *divergence time estimation*
→ active research area
→ most codes rely on the phylogenetic
likelihood function and Bayesian
Statistics (MCMC methods)

# Dating Trees

Evolutionary time →

Root: 3 million years

1 million years

dated fossil 2 million years

2 million years

t1   t2   t3   t4

**Ultrametric tree**

But how do we place the fossil?
→ typically no DNA data available

Fossil placement:
→ ad hoc using empirical knowledge
→ computationally using morphological data

**The input for a phylogenetic analysis need not be molecular data!**

**We can also use sequences of morphological traits ("Merkmale")!**

# Remember that we deal with extant species!



Evolutionary time

t1   t2   t3   t4   2018

**Ultrametric tree**

# Morphological Traits

```
t1: 1000
t2: 0100
t3: 0010
T4: 0001

or:

t1: 0
t2: 1
t3: 2
t4: 3
```

What image best matches the extent of your natural brow line (without hair removal)?



52%

37%

7%

1%

I'm not sure    1%

None of the above    1%

# Morphological Traits

```
t1: 1000
t2: 0100
t3: 0010
T4: 0001

or:

t1: 0
t2: 1
t3: 2
t4: 3
```



What image best matches the extent of your natural brow line (without hair removal)?

52%

37%

7%

1%

I'm not sure — 1%

None of the above — 1%

**Traits need not be discrete,
they can also be continuous, e.g., bone ratios**

# Alignment-Free Tree Inference

Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT

Pair-wise distances
e.g., pair-wise sequence
alignment scores

Tree inference
program

Taxon 1                    Taxon 3

Taxon 2                    Taxon 4

# Alignment-Free Tree Inference

Taxon 1:ACGTTT
Taxon 2:ACGTT
Taxon 3:ACCCT
Taxon 4:AGGGTTT

Pair-wise distances
e.g., pair-wise sequence
alignment scores

Tree inference
program

Taxon 1

Taxon 3

Alignment-free
tree inference
is typically less
accurate → we have
not established homology
via a MSA

Taxon 2

Taxon 4

# How many unrooted 4-taxon trees exist?

# How many rooted 4-taxon trees exist?

# Tree Counts

- Unrooted binary trees

  - *4* taxa → *3* distinct trees
  - A tree with *n* taxa has *n-2* inner nodes
  - And *2n-3* branches

- Rooted binary trees

  - *4* taxa → *3* unrooted trees * *5* branches each (rooting points) = *15* trees
  - *n-1* inner nodes
  - *2n-2* branches

# The number of trees

3 taxa = 1 tree

# The number of trees



4 taxa: 3 trees
u: # trees of size 4-1 := 1
v: # branches in a tree of size 4-1 := 3
Number of unrooted binary trees with 4 taxa: u * v = 3

# The number of trees



5 taxa: 15 trees
u = 3
v = 5
Number of unrooted trees with 5 taxa: 3 * 5 = 15

# The number of trees



6 taxa: 105 trees
u = 15
v = 7
u * v = 105

# The number of trees explodes!



BANG !

# Some Numbers

| Number of Organisms | Number of alternative Trees |
|:---:|:---:|
| 3 | 1 |
| 4 | 3 |
| 5 | 15 |
| 6 | 105 |
| 7 | 945 |
| 10 | 2.027.025 |
| 15 | 7.905.853.580.625 |
| 20 | $2.21 * 10^{20}$ |
| 50 | $2.84 * 10^{76}$ |

Table 2.1: Number of possible trees for phylogenies with 3–50 organisms

# Equation for the number of unrooted trees

- Simple proof via induction

$$\prod_{i=3}^{n}(2i-5)$$

- The number of rooted trees for *n* taxa simply is the number of unrooted trees for *n+1* taxa

- The additional (*n+1*[th]) taxon represents all possible rootings for all unrooted trees with *n* taxa

58

# # trees with 2000 tips

```
stamatak@exelixis:~/Desktop/GIT/TreeCounter$ ./treeCounter -n 2000

GNU GPL tree number calculator released June 2011 by Alexandros Stamatakis

Number of unrooted binary trees for 2000 taxa: 300496381742116561516329100656818149813772320742370130895049540430126365252583082108276859966882470004643527352142656342882958915023446000631493969130632970436056184861877465482277991223536809233455563199910834597693126756525012899867433187752811401960991631522367030609121735709762379847705467667779532479718261438527333822672778425073725284991666968758440351057958702068650581768704466631812374290102143850643247136093449166702113596975694030066625264647926912455103149423661955428241182776251148487582545812279142898011326489026740337612947127457670362675790868431696607186098479418188659572145570447445722886617290535835207442536881231240106613156948861960941195646736200342575241335277575085829161096422575727699767991408283343210161327401652830993803904592327690690035972919709940739349563486203899010742687282297597465537710225767267684285801187722495010621811734052320826539734296222735253659051586563138327203111984198746759973864631829032038325230859799799221610122721578080524814583120684401676062393060097116167297155047284877996343375313489942303724373478791319890859537640701348494461138775725769524087024617201078742973804622750525457066689372319418206440706891884003870590289772197516454495975821662130620506461776109948566373416818358498932907699338206780105243728461492403422961155182609778228619192672071295189589360099591309742330723163825184281103305710174411568843051318658775443763085003114511107238370397074651822320404061547082730786299575493310312752086167006607912980142622300565123522718063819509335872651728623589020520016144361756075654286471422126613004434807084067501589247673166341539540575074474994909831496473031080411401891849735912811228378774049884836445210242056642446386009389965085742961947269054301528123752651096581528469979703679217112903556809818079169587951614159281049528179855847292534447864442443599808531537204796814969465991768614533701051985928577157482455943377242369582576242663016946320482495182255939287403177623433881048604630975191556923871167513095213415098816715464307862352606237786406838680424690252749113931927680261151599058260388673317293071367390340361863746398060576483647467027444672788088533707425442192272667774700332940332010382880351126890262551830967919483586789293701637681753048206338943871497931152353698229625111116307148294599211620803302684762013335690441089668145436150905155877581167977001256391215111623744417049737170460402948110411482228646613191882199757138336835207252605520276982397461321849524926489705079039836025625560628985228883956135787415657564889992608732866126306425432602489792291135600716405739845163752452433769437558573847255455643975996042559146401122211447552355731762399730577471839565312174165322959866759012941161239240722093250369673124884491553759210650656015416720774159236240868667675342865129648887390597075788024733934634708481590116397727977474804173162687009167287356121642268468160683198959801260376485615312781611689587215123123308760063473381097253118423339640390937378395066835578735307886358646400563299499490631187424029092779272693300322445377597972248734568915114585570783850541681667667425811301958063621907500790295031088209097271748136436989473971079932777700676301730617566538739726037771730084413439405123669055544932486165082539957795036326704947844293498853172797348177797146567175151178876396434069332458076346110734214328195049909680874027397688914704517472055543896939666874260147724189469321290245331733418828677319465354411330210086657508171324034264758048921862366346160795594372051639515409694980648624230979694721116966596158004178839922322646284984423522769263911243607675089767171896834593378907423463545571935306656153799002779162656363618619740485909382340622354596797331372136593437175855906644393646132830011367260193406870644233948919921530435281541659630118549942363486352458574664283460906291627956492658472300360855559889761916129324814009459248989946846886232225860170551468905649837276003927487069550463788819741699429049102853180410776516007263384721638903782270013843599599730265727198804324664312359758555217196905139210102265963247887830977405333131551762978152880718652603176327264828094499370456258099380530584976995700802893798080149029010052938984722799471678048216894241591182842576964647865057315325178302336307298251692210346584265894447464916123854689718507968172901390321828341111848213847677283165486532123173820041319905105189670222018870495856871805095907303606930402937216038968917605587676955382318093705826257083898387409098468656634271397500013291835105943321729879825243707508272087959859437157667660155782699660343197752623308898996258780062800956094441693237794495544103369658626155625601066939030320387897098367378608705664143358510611165831452042451320850858999493236483168967119495167161956762270709069738895888555795624666415365617235493018073940047605298017217713916867880002778519661730700612845173075825037356431020651124437308252296250404531605907413438818725634779138306605909311880252231008534017684026140153961698919207514710803375770884974014183459975397205987868206487911606496985817760115397205849822269890718134943269180182117331880636539108936898117148913574566805428074851701758582666396335701893544498326697628350926579222017463721902731196417514899440100796368760178267471070199454732188878327426088966724371574713420600009370425130989363053745978427998040313298941726649229042573095836853441621564055729028206622400386323752638091023326989783886042375962560156797526269507986398668104294832333160267216555178120899264677804935741326387137408423885546538336158643451305439624281397279559725995110706314305992615495622958320232708057681156690489586610522030057372529847211187478271367136660586692710948755639748584947591081972703387828443986448674345620095816193031472734596190049931842433797524366248936332124485059719952523668529249305346252764137853413208943128901523738092556045987090912766662329678703328882059134949580074074473143388800724532321747309659741967114444531271327902051010047671014350638857953478447255389801541923317027519896180635152682543173193832925891931530164130548972311128666465492971930479296432829556719092881692091042334122007454242049900872585046208051104875883059495990311188736668509414882172573457635523396403848131821316740835900691640005326225818478376506780445117771732865818989215358309447765350341796875

Approximately 3.00 times 10^6328
```

# A side-note
# The `treeCounter` tool

- Evidently, the tree count can not be computed using normal integers

  → we need an arbitrary precision library

  → I used the GNU GMP (Multiple Precision Arithmetic) library

  → `treeCounter` available as open-source code at

    https://github.com/stamatak

  → Has anybody already used GNU GMP?

# Scoring Trees

- Now we know how many **unrooted** candidate trees there exist for *n* taxa

- How do we chose among them?

  → we need some scoring criterion f() to evaluate them

  → finding the optimal tree under most criteria is NP-Hard



```
A:  ACGG
B:  AGGG
C:  GA-A
D:  AAGG
```
A, B, C, D tree — f() — 1.0

```
A:  ACGG
B:  AGGG
C:  GA-A
D:  AAGG
```
A, C, B, D tree — f() — 2.0

```
A:  ACGG
B:  AGGG
C:  GA-A
D:  AAGG
```
A, B, D, C tree — f() — 3.0

# What can we do with Phylogenies?

# What can we do with Phylogenies?



Phylogenetic placement for identifying anonymous sequences
Examples:
· Bird strike
· Bacteria
· Viral strains

Unknown/anonymous sequence/species

?

Known Species 1

Known Species 2

Known Species 4

reference phylogeny

# What can we do with Phylogenies?



reference phylogeny

# Diversification Rates



**From:** Charles C. Davis, Hanno Schaefer: "Plant Evolution: Pulses of Extinction and Speciation in Gymnosperm Diversity", *Current Biology*, 2011.

# Diversification Rates

- With former PostDoc Stephen Smith: "Understanding angiosperm diversification using small and large phylogenetic trees", *American Journal of Botany* 98 (3), 404-414, 2011.

- Largest tree of angiosperms computed to date

- 55,000 taxa

# Diversification Rates

- With former PostDoc Stephen Smith: "Understanding angiosperm diversification using small and large phylogenetic trees", *American Journal of Botany* 98 (3), 404-414, 2011.

- Largest tree of angiosperms computed to date

- 55,000 taxa

Visualizing big trees also represents a challenge → graph drawing & layout algorithms.

# Influenza Outbreaks

# And of course SARS-CoV-2



**Phylogenetic analysis of SARS-CoV-2 data is difficult**

Benoit Morel[*,1], Pierre Barbera[*,1], Lucas Czech[3], Ben Bettisworth[1], Lukas Hübner[1,2], Sarah Lutteropp[1], Dora Serdari[1], Evangelia-Georgia Kostaki[5], Ioannis Mamais[6], Alexey M Kozlov[1], Pavlos Pavlidis[4], Dimitrios Paraskevis[5], and Alexandros Stamatakis[1,2]

# Snakebites

Australia has more poisonous snakes than any other continent, and many people die from snakebites each year. Developing effective antivenins is thus a high priority, but little is known about the venins of most species.

Phylogenetic analysis is helping with this task because venin properties correlate strongly with evolutionary relationships.

Although the red-bellied black snake looks very different from the king brown, it is actually closely related and can be treated with the same antivenin.

Conversely, the western brown looks very similar to the king brown, but it is only distantly related and thus responds best to different antivenin.

The phylogeny is also predictive: the recent demonstration that the poorly-known barclick is closely related to the death adder (orange lineage) predicts that the former is also highly dangerous and might respond to widely-available death adder antivenin.

70



Black whip snake
Talpan
Fierce snake
Common brown
Western brown
Dugte
Collatt's snake
Spotted black
Butler's snake
King brown
Red-bellied black
Death adder
Barclick
Small-eyed snake
Australian copperhead
Tiger snake
Rough-scaled snake
Broad-headed snake

**Recommended Antivenene**

Taipan
Brownsnake
Blacksnake
Death adder
Tiger snake

# Snakebites
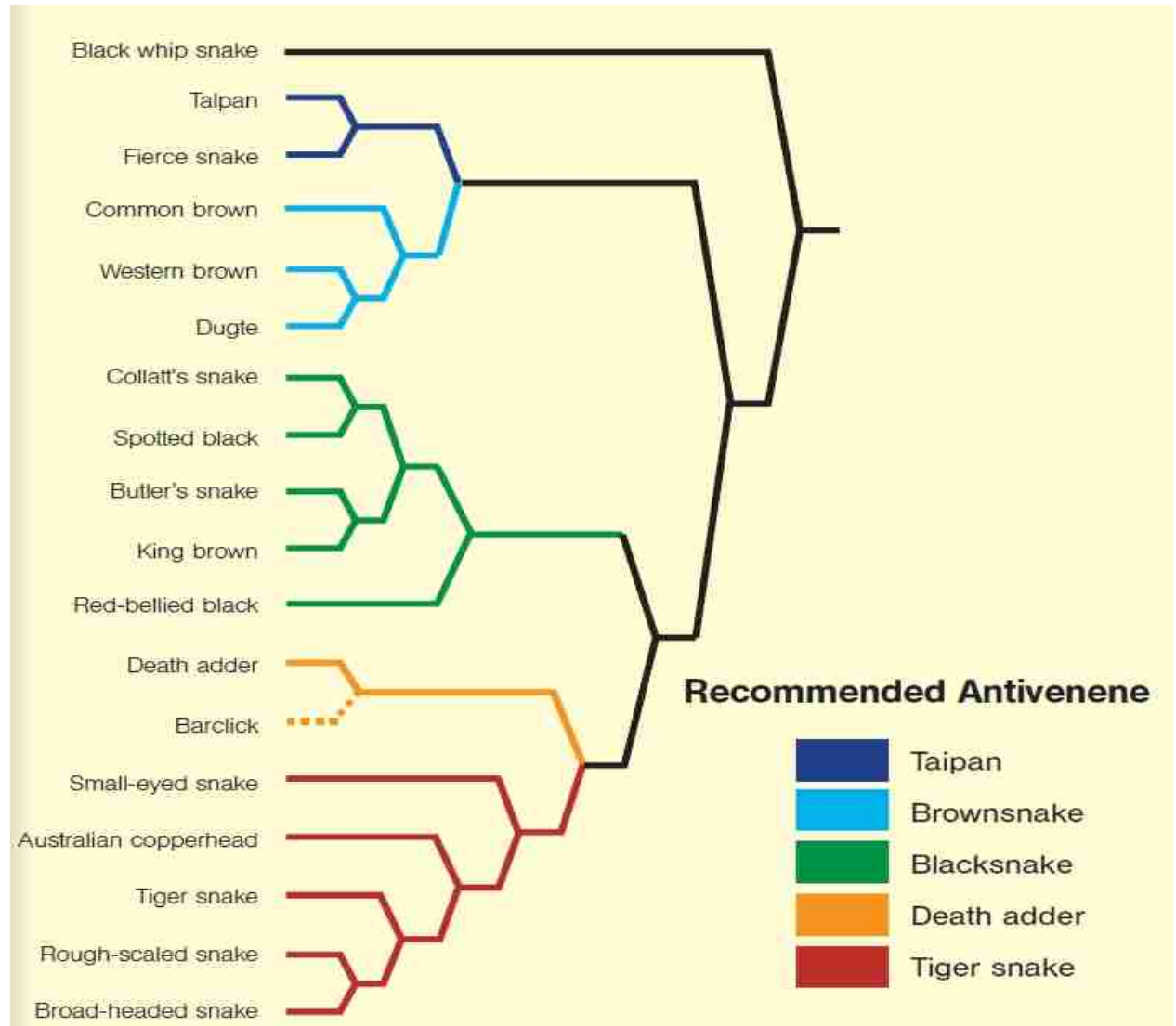
Australia has more poisonous snakes than any other continent, and many people die from snakebites each year. Developing effective antivenins is thus a high priority, but little is known about the venins of most species.

Phylogenetic analysis is helping with this task because venin properties correlate strongly with evolutionary relationships.

Although the red-bellied black snake looks very different from the king brown, it is actually closely related and can be treated with the same antivenin.

Conversely, the western brown looks very similar to the king brown, but it is only distantly related and thus responds best to different antivenin.

The phylogeny is also predictive: the recent demonstration that the poorly-known barclick is closely related to the death adder (orange lineage) predicts that the former is also highly dangerous and might respond to widely-available death adder antivenin.



Potentially: convergent evolution

Black whip snake
Taipan
Fierce snake
Common brown
Western brown
Dugite
own
bellied black
Death adder
Barclick
Small-eyed snake
Australian copperhead
Tiger snake
Rough-scaled snake
Broad-headed snake

**Recommended Antivenene**

Taipan
Brownsnake
Blacksnake
Death adder
Tiger snake

71

# What can we do with phylogenetic trees?

- identifying unknown species

- divergence time estimates

- diversification rates

- viral outbreaks

- forensics → M.L. Metzker, D.P. Mindell, X.M. Liu, R.G. Ptak, R.A. Gibbs, D.M. Hillis: "Molecular evidence of HIV-1 transmission in a criminal case" PNAS: 99(22):14292-7, 2002.

# *"Nothing in Biology makes sense, except in the light of evolution"*

Why this increase in
Phylogenetics papers?
Advances in:
•Sequencing technology
•Hardware
•Methods & Tools



Number of 'Phylogen*' Publications Per Year

# Building Trees

- We distinguish between
  - *Distance-based methods*
    - → use MSA to compute a matrix of pair-wise distances
    - → build a tree using these distances
    - → Heuristics (essentially hierarchical clustering methods)
      - → *Neighbor Joining:* NJ
      - → *Unweighted Pair Group Method with Arithmetic Mean:* UPGMA
    - → least-squares method: explicit optimality criterion
  - *Character-based methods*
    - → optimality criteria *f()* operate directly on the MSA & tree
      - → parsimony
      - → maximum likelihood
      - → Bayesian inference
    - → take the current tree topology & MSA to calculate a score
    - → the score tells us how well the MSA data fits the tree

# Building Trees

- We distinguish between

  - *Distance-based methods*
    - → use MSA to compute a matrix of pair-wise distances
    - → build a tree using these distances
    - → Heuristics (essentially hierarchical clustering methods)
      - → *Neighbor Joining:* NJ
      - → *Unweighted Pair Group Method with Arithmetic Mean:* UPGMA
    - → least-squares method: explicit optimality criterion
  - *Character-based methods*
    - → optimality criteria *f()* operate directly on the MSA & tree
      - → parsimony
      - → maximum likelihood
      - → Bayesian inference
    - → take the current tree topology & MSA to calculate a score
    - → the score tells us how well the MSA data fits the tree

Less accurate, but faster

Slow, but more accurate

# Building Trees

- We distinguish between

  - *Distance-based methods*

    - → use MSA to compute a matrix of pair-wise distances
    - → build a tree using these distances
    - → Heuristics (essentially hierarchical clustering methods)
      - → *Neighbor Joining:* NJ
      - → *Unweighted Pair Group Method with Arithmetic Mean:* UPGMA
    - → least-squares method: explicit optimality criterion

  - *Character-based methods*

    - → optimality criteria *f()* operate directly on the MSA
      - → parsimony
      - → maximum likelihood
      - → Bayesian inference
    - → take the current tree topology & MSA to calculate a score
    - → the score tells us how well the MSA data fits the tree

Less accurate, but faster

Memory-intensive!

Slow, but more accurate

# Building Trees

- We distinguish between
  - *Distance-based methods*
    - → use MSA to compute a matrix of pair-wise distances
    - → build a tree using these distances
    - → Heuristics (essentially hierarchical clustering methods)
      - → *Neighbor Joining:* NJ
      - → *Unweighted Pair Group Method with Arithmetic Mean:* UPGMA
    - → least-squares method: explicit optimality
  - *Character-based methods*
    - → optimality criteria *f()* operate directly
      - → parsimony
      - → maximum likelihood
      - → Bayesian inference
    - → take the current tree topology & MSA to calculate a score
    - → the score tells us how well the MSA data fits the tree

Less accurate, but faster

Slow, but more accurate

What could be the computational limitation here?

Memory-intensive!

# Building Trees

- We distinguish between

  - *Distance-based methods*

    - → use MSA to compute a matrix of pair-wise distances
    - → build a tree using these distances
    - → Heuristics (essentially hierarchical clustering methods)
      - → *Neighbor Joining:* NJ
      - → *Unweighted Pair Group Method with Arithmetic* ~~~~ UPGMA
    - → least-squares method: explicit optimality criterion

  - *Character-based methods*

    - → optimality criteria *f()* operate directly on
      - → parsimony
      - → maximum likelihood
      - → Bayesian inference
    - → take the current tree topology & MSA to calculate a score
    - → the score tells us how well the MSA data fits the tree

Less accurate, but faster

Storing this matrix can become problematic memory-wise
→ out-of-core/external memory algorithms
→ e.g.: NINJA tool for Neighbor joining
"Large-scale neighbor-joining with ninja"
T Wheeler,
*Algorithms in Bioinformatics*, 2009

# Out-of-core Algorithms

- Definition from Wikipedia:

  *Out-of-core or External memory algorithms are algorithms that are designed to process data that is too large to fit into a computer's main memory at one time. Such algorithms must be optimized to efficiently fetch and access data stored in slow bulk memory such as hard drive or tape drives.*
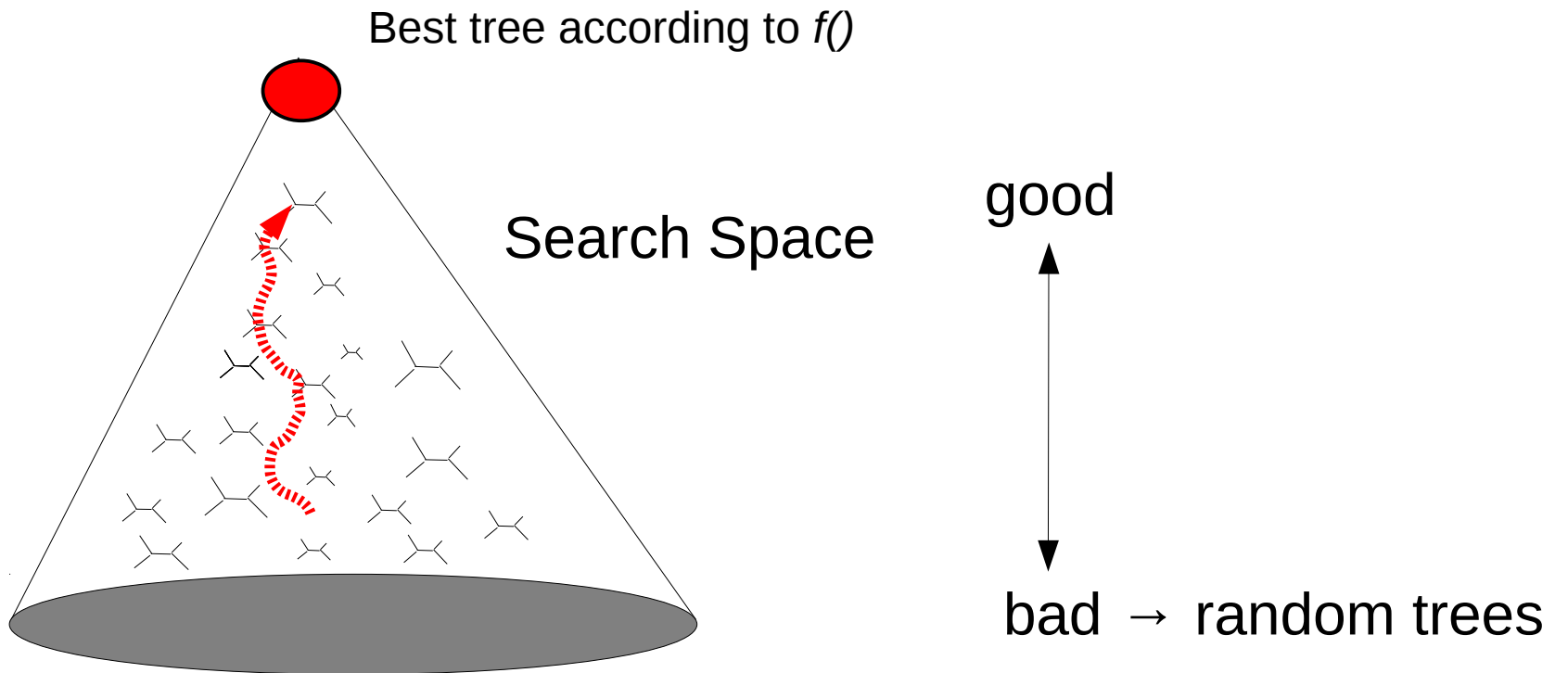
- We do the data transfer RAM ↔ disk explicitly from within the application code by using application-specific knowledge (e.g., about the data access patterns)

- This is to circumvent the paging procedure that would normally be initiated by the OS

- Out-of-core algorithms are typically much faster than the *application-agnostic* paging procedure carried out by the OS

- For an example from phylogenetics see:

  Fernando Izquierdo-Carrasco, Alexandros Stamatakis: "Computing the Phylogenetic Likelihood Function Out-of-Core", *IEEE HICOMB 2011 workshop*, Anchorage, USA, May 2011.
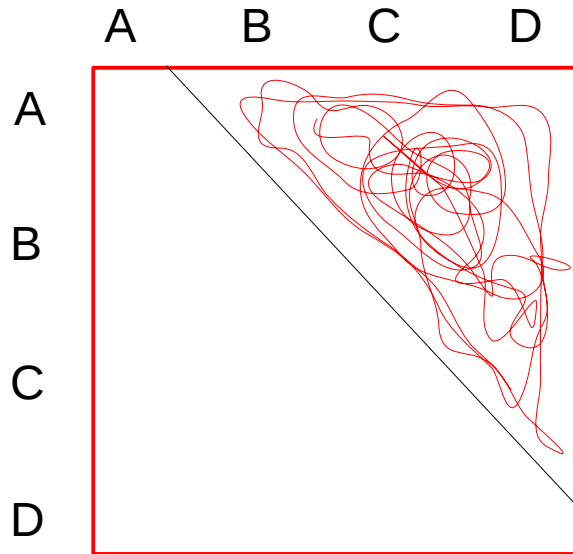
# NP-Hardness

- Because of the super-exponential increase in the number of possible trees for *n* taxa ...

- all interesting criteria on trees are NP-hard:

  - Least squares

  - Parsimony → discrete criterion

  - Likelihood → statistical criterion

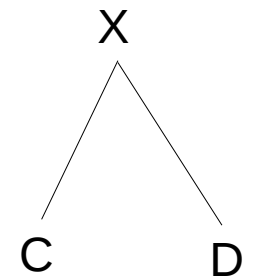  - Bayesian → integrate likelihood over entire tree space

# Search Space

Best tree according to *f()*
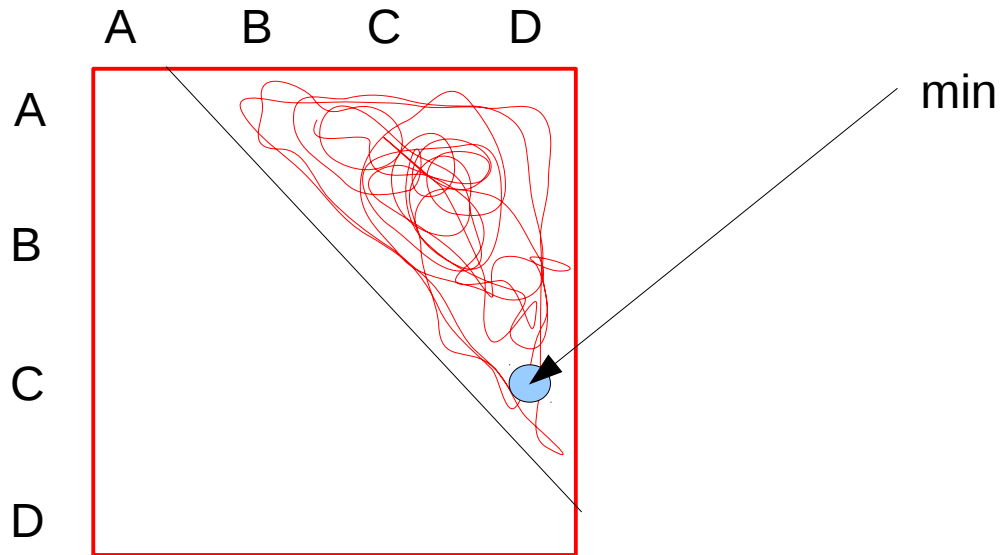


Search Space

good

bad → random trees

# Neighbor Joining → Principle


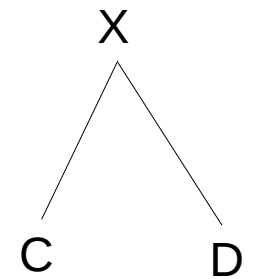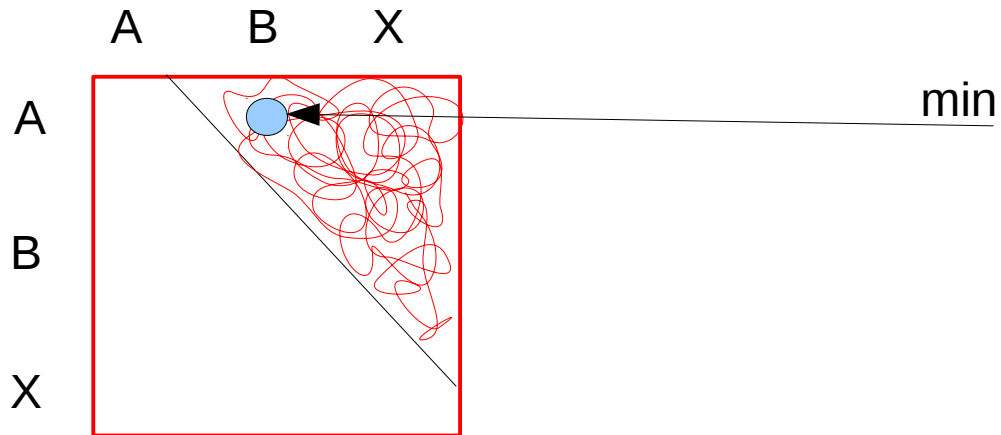
Given a kind of distance matrix $D_{i,j}$ where *i,j=1...4*

# Neighbor Joining → Principle



Given a kind of distance matrix $D_{i,j}$ where $i,j=1...4$
Find minimum and merge taxa

# Neighbor Joining → Principle



Given a kind of distance matrix $D_{i,j}$ where *i,j=1...4*
Find minimum and merge taxa
Compute a new distance matrix of size *n-1 = 3*
Find minimum

# Neighbor Joining → Principle



Given a kind of distance matrix $D_{i,j}$ where $i,j=1...4$
Find minimum and merge taxa
Compute a new distance matrix of size $n-1 = 3$
Find minimum and merge taxa

# Neighbor Joining → Principle



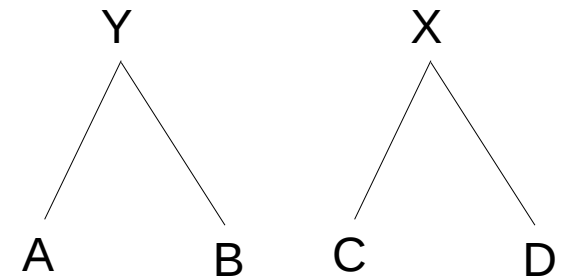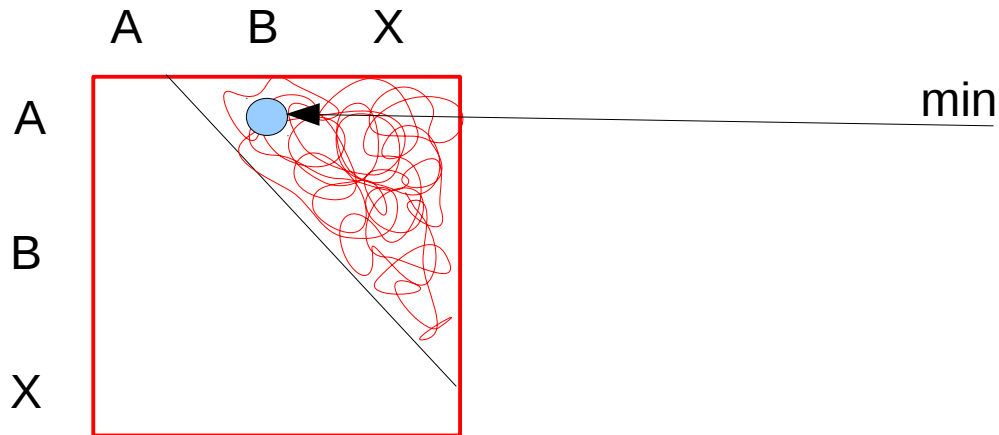Given a kind of distance matrix $D_{i,j}$ where *i,j=1...4*
Find minimum and merge taxa
Compute a new distance matrix of size *n-1 = 3*
Find minimum and merge taxa
Etc.
Space complexity: *O(n²)*
Time complexity: *O(n³)*
Key question: how do we compute distance between *X* and *A* or *X* and *B* respectively
→ for progressive alignment we may align the profile of *X* with all remaining sequences

# Neighbor Joining Algorithm

- For each tip compute

  $$u_i = \Sigma_j\, D_{ij} / (n-2)$$

    → this is in principle the average distance to all other tips

    → the denominator is *n-2* instead of *n,* see below why

- Find the pair of tips, *(i, j)* for which $D_{ij}-u_i-u_j$ is minimal
- Connect the tips *(i,j)* to build a new ancestral node *X*
- The branch lengths from the ancestral node *X* to *i* and *j* are:

    $b_i = 0.5\ D_{ij} + 0.5\ (u_i-u_j)$

    $b_j = 0.5\ D_{ij} + 0.5\ (u_j-u_i)$

- Update the distance matrix:
    → Compute distance between the new node *X* and each remaining tip as follows:

    $D_{ij,k} = (D_{ik}+D_{jk}-D_{ij})/2$

- Replace tips *i* and *j* by the new node *X* which is now treated as a tip
- Repeat until only two nodes remain
    → connect the remaining two nodes with each other

# Neighbor Joining Algorithm

|   | A | B | C | D |
|---|---|---|---|---|
| A | - | 17 | 21 | 27 |
| B |   | - | 12 | 18 |
| C |   |   | - | 14 |
| D |   |   |   | - |

# Neighbor Joining Algorithm

```
      A    B    C    D
A  -   17   21   27
B        -   12   18
C             -   14
D                  -
```

<span style="color:red">Distance matrix, usually denoted as *D*</span>

```
i            uᵢ
A    (17+21+27)/2=32.5
B    (17+12+18)/2=23.5
C    (21+12+14)/2=23.5
D    (27+18+14)/2=29.5
```

<span style="color:red">Average distance</span>

# Neighbor Joining Algorithm

|   | A | B | C | D |
|---|---|---|---|---|
| A | - | 17 | 21 | 27 |
| B |   | - | 12 | 18 |
| C |   |   | - | 14 |
| D |   |   |   | - |

| i | $u_i$ |
|---|---|
| A | (17+21+27)/2=32.5 |
| B | (17+12+18)/2=23.5 |
| C | (21+12+14)/2=23.5 |
| D | (27+18+14)/2=29.5 |

|   | A | B | C | D |
|---|---|---|---|---|
| A | - | -39 | -35 | -35 |
| B |   | - | -35 | -35 |
| C |   |   | - | -39 |
| D |   |   |   | - |

$$D_{ij}-u_i-u_j$$

Usually denoted as *Q* matrix

# Neighbor Joining Algorithm

|   | A | B | C | D |
|---|---|---|---|---|
| A | - | 17 | 21 | 27 |
| B |   | - | 12 | 18 |
| C |   |   | - | 14 |
| D |   |   |   | - |

| i | $u_i$ |
|---|-------|
| A | (17+21+27)/2=32.5 |
| B | (17+12+18)/2=23.5 |
| C | (21+12+14)/2=23.5 |
| D | (27+18+14)/2=29.5 |

|   | A | B | C | D |
|---|---|---|---|---|
| A | - | -39 | -35 | -35 |
| B |   | - | -35 | -35 |
| C |   |   | - | **-39** |
| D |   |   |   | - |

$$D_{ij} - u_i - u_j$$

# Neighbor Joining Algorithm

|   | A | B | C | D |
|---|---|---|---|---|
| A | - | 17 | 21 | 27 |
| B |   | - | 12 | 18 |
| C |   |   | - | 14 |
| D |   |   |   | - |

| i | $u_i$ |
|---|-------|
| A | (17+21+27)/2=32.5 |
| B | (17+12+18)/2=23.5 |
| C | (21+12+14)/2=23.5 |
| D | (27+18+14)/2=29.5 |

|   | A | B | C | D |
|---|---|-----|-----|-----|
| A | - | -39 | -35 | -35 |
| B |   | - | -35 | -35 |
| C |   |   | - | -39 |
| D |   |   |   | - |

$$D_{ij}-u_i-u_j$$

# Neighbor Joining Algorithm

|   | A | B  | C  | D  |
|---|---|----|----|----|
| A | - | 17 | 21 | 27 |
| B |   | -  | 12 | 18 |
| C |   |    | -  | 14 |
| D |   |    |    | -  |

| i | $u_i$ |
|---|-------|
| A | (17+21+27)/2=32.5 |
| B | (17+12+18)/2=23.5 |
| C | (21+12+14)/2=23.5 |
| D | (27+18+14)/2=29.5 |

|   | A | B   | C   | D   |
|---|---|-----|-----|-----|
| A | - | -39 | -35 | -35 |
| B |   | -   | -35 | -35 |
| C |   |     | -   | -39 |
| D |   |     |     | -   |

$D_{ij}-u_i-u_j$

$b_C$ = 0.5 x 14 + 0.5 x (23.5-29.5) = 4
$b_D$ = 0.5 x 14 + 0.5 x (29.5-23.5) = 10

# Neighbor Joining Algorithm

|   | A | B | C | D | X |
|---|---|---|---|---|---|
| A | - | 17 | 21 | 27 | |
| B | | - | 12 | 18 | |
| C | | | - | 14 | |
| D | | | | - | |
| X | | | | | - |

C     D
 \   /
 4\ /10
   X

# Neighbor Joining Algorithm

|   | A | B | C | D | X |
|---|---|---|---|---|---|
| A | - | 17 | 21 | 27 |   |
| B |   | - | 12 | 18 |   |
| C |   |   | - | 14 |   |
| D |   |   |   | - |   |
| X |   |   |   |   | - |

$D_{XA} = (D_{CA} + D_{DA} - D_{CD})/2$
$\quad\quad = (21 + 27 - 14)/2$
$\quad\quad = 17$

$D_{XB} = (D_{CB} + D_{DB} - D_{CD})/2$
$\quad\quad = (12 + 18 - 14)/2$
$\quad\quad = 8$

C      D
\      /
4 \  / 10
   X

# Neighbor Joining Algorithm

|   | A | B | C | D | X |
|---|---|---|---|---|---|
| A | - | 17 | 21 | 27 | 17 |
| B |   | - | 12 | 18 | 8 |
| C |   |   | - | 14 |   |
| D |   |   |   | - |   |
| X |   |   |   |   | - |

$$D_{XA} = (D_{CA} + D_{DA} - D_{CD})/2$$
$$= (21 + 27 - 14)/2$$
$$= 17$$

$$D_{XB} = (D_{CB} + D_{DB} - D_{CD})/2$$
$$= (12 + 18 - 14)/2$$
$$= 8$$

C        D

4    10

X

# Neighbor Joining Algorithm

|   | A | B | X |
|---|---|---|---|
| A | - | 17 | 17 |
| B |   | - | 8 |
| X |   |   | - |

$$D_{XA} = (D_{CA} + D_{DA} - D_{CD})/2$$
$$= (21 + 27 - 14)/2$$
$$= 17$$

$$D_{XB} = (D_{CB} + D_{DB} - D_{CD})/2$$
$$= (12 + 18 - 14)/2$$
$$= 8$$

C       D

4    10

X

# Neighbor Joining Algorithm

|   | A | B | X |
|---|---|---|---|
| A | - | 17 | 17 |
| B |   | - | 8 |
| X |   |   | - |

| i | $u_i$ |
|---|---|
| A | (17+17)/1 = 34 |
| B | (17+8)/1 = 25 |
| X | (17+8)/1 = 25 |

C       D

4    10

X

# Neighbor Joining Algorithm

```
    A   B   X                   i          u_i

A   -   17  17              A   (17+17)/1 = 34

B       -   8               B   (17+8)/1 = 25

X           -              X   (17+8)/1 = 25
```
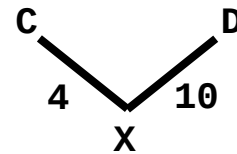
```
    A    B     X

A   -    -42   -28

B        -     -28

X              -


    D_ij - u_i - u_j
```
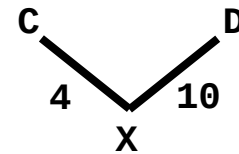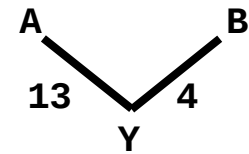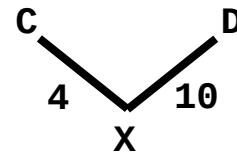
# Neighbor Joining Algorithm

|   | A | B | X |
|---|---|---|---|
| A | - | 17 | 17 |
| B |   | - | 8 |
| X |   |   | - |

| i | $u_i$ |
|---|---|
| A | (17+17)/1 = 34 |
| B | (17+8)/1 = 25 |
| X | (17+8)/1 = 25 |

|   | A | B | X |
|---|---|---|---|
| A | - | **-42** | -28 |
| B |   | - | -28 |
| X |   |   | - |

$D_{ij} - u_i - u_j$

# Neighbor Joining Algorithm

|   | A | B | X |
|---|---|---|---|
| A | - | 17 | 17 |
| B |   | - | 8 |
| X |   |   | - |

| i |  | $u_i$ |
|---|---|---|
| A | (17+17)/1 | = 34 |
| B | (17+8)/1 | = 25 |
| X | (17+8)/1 | = 25 |

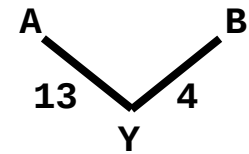|   | A | B | X |
|---|---|---|---|
| A | - | **-42** | -28 |
| B |   | - | -28 |
| X |   |   | - |

$$D_{ij} - u_i - u_j$$
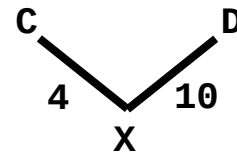


$b_A = 0.5 \times 17 + 0.5 \times (34-25) = 13$
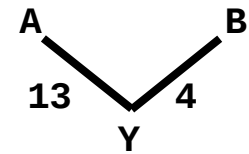$b_D = 0.5 \times 17 + 0.5 \times (25-34) = 4$

# Neighbor Joining Algorithm

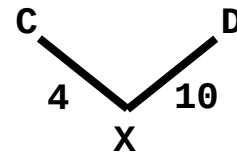|   | A | B | X | Y |
|---|---|---|---|---|
| A | - | 17 | 17 |   |
| B |   | - | 8 |   |
| X |   |   | - |   |
| Y |   |   |   |   |

```
C        D          A        B
 \      /            \      /
  4    10            13    4
   \  /                \  /
    X                   Y
```

# Neighbor Joining Algorithm

|   | A | B | X | Y |
|---|---|---|---|---|
| A | - | 17 | 17 | |
| B | | - | 8 | |
| X | | | - | 4 |
| Y | | | | |

$$D_{YX} = (D_{AX} + D_{BX} - D_{AB})/2$$
$$= (17 + 8 - 17)/2$$
$$= 4$$

C          D          A          B

4      10        13      4

X                      Y

# Neighbor Joining Algorithm

```
    X   Y

X   -   4

Y       -
```

$$D_{YX} = (D_{AX} + D_{BX} - D_{AB})/2$$
$$= (17 + 8 - 17)/2$$
$$= 4$$

```
C            D        A            B
 \          /          \          /
  4       10            13       4
   \      /              \      /
      X                     Y
```

# Neighbor Joining Algorithm

|   | X | Y |
|---|---|---|
| X | - | 4 |
| Y |   | - |

$$D_{YX} = (D_{AX} + D_{BX} - D_{AB})/2$$
$$= (17 + 8 - 17)/2$$
$$= 4$$

# Neighbor Joining Algorithm

|   | A | B | C | D |
|---|---|---|---|---|
| A | - | 17 | 21 | 27 |
| B |   | - | 12 | 18 |
| C |   |   | - | 14 |
| D |   |   |   | - |