# RAxML-NG:
# a fast, scalable and user-friendly tool for maximum likelihood phylogenetic inference

Alexandros Stamatakis and Alexey Kozlov

Computational Molecular Evolution Group,
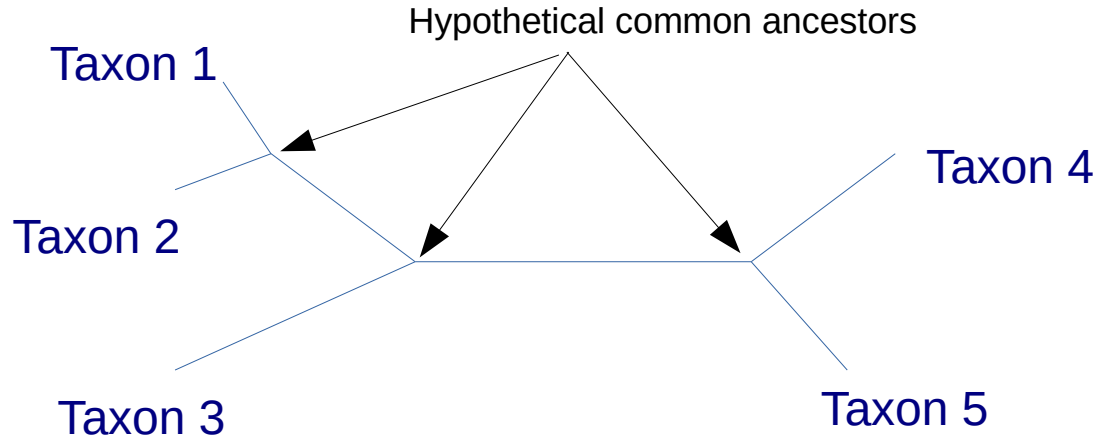Heidelberg Institute for Theoretical Studies

Institute for Theoretical Informatics,
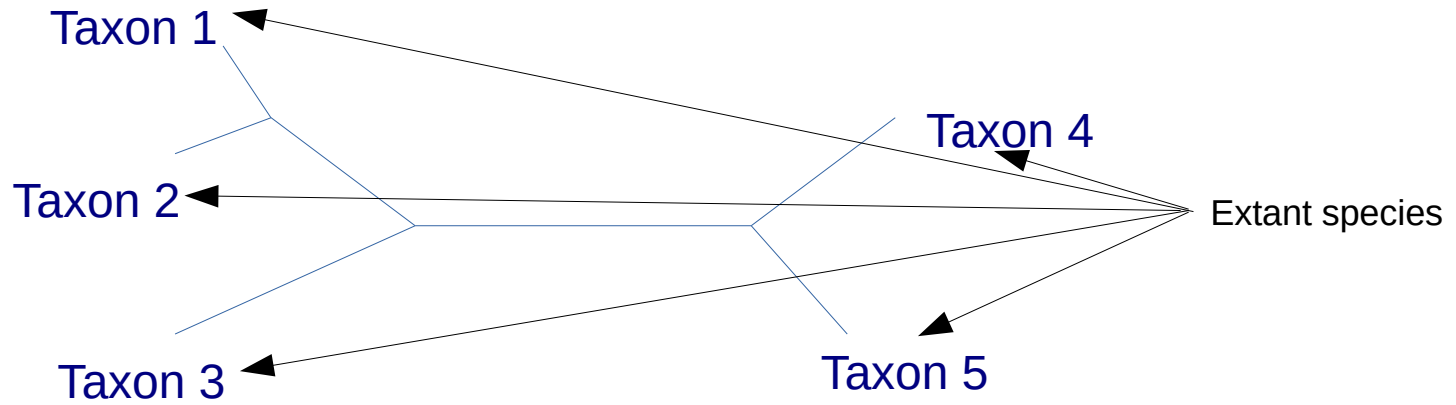Karlsruhe Institute of Technology

# Outline

- Introduction to Phylogenetic Inference - *Alexandros*

- The RAxML Search Algorithm - *Alexandros*

- Improvements in RAxML **N**ext **G**eneration - *Alexey*

- Tutorial - *Alexey*

# A phylogeny

Hypothetical common ancestors
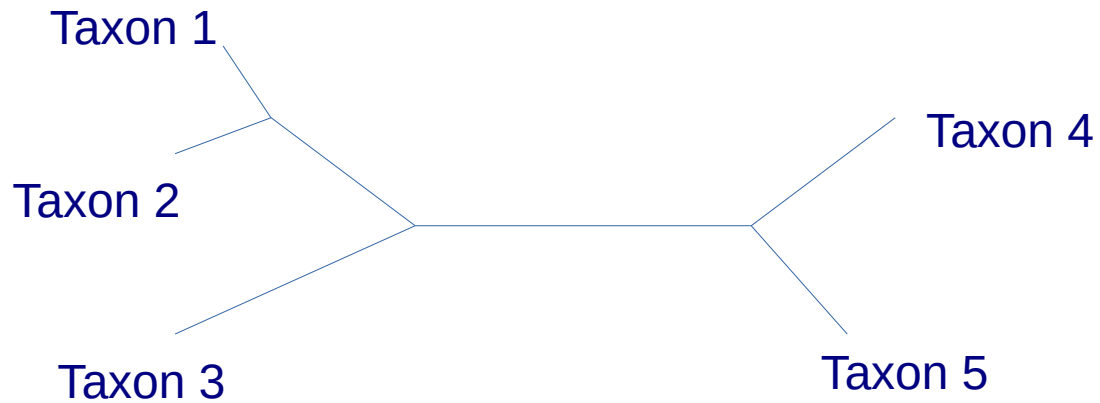
Taxon 1

Taxon 4

Taxon 2

Taxon 3

Taxon 5

Phylogenies describe evolutionary relationships among species

# A phylogeny
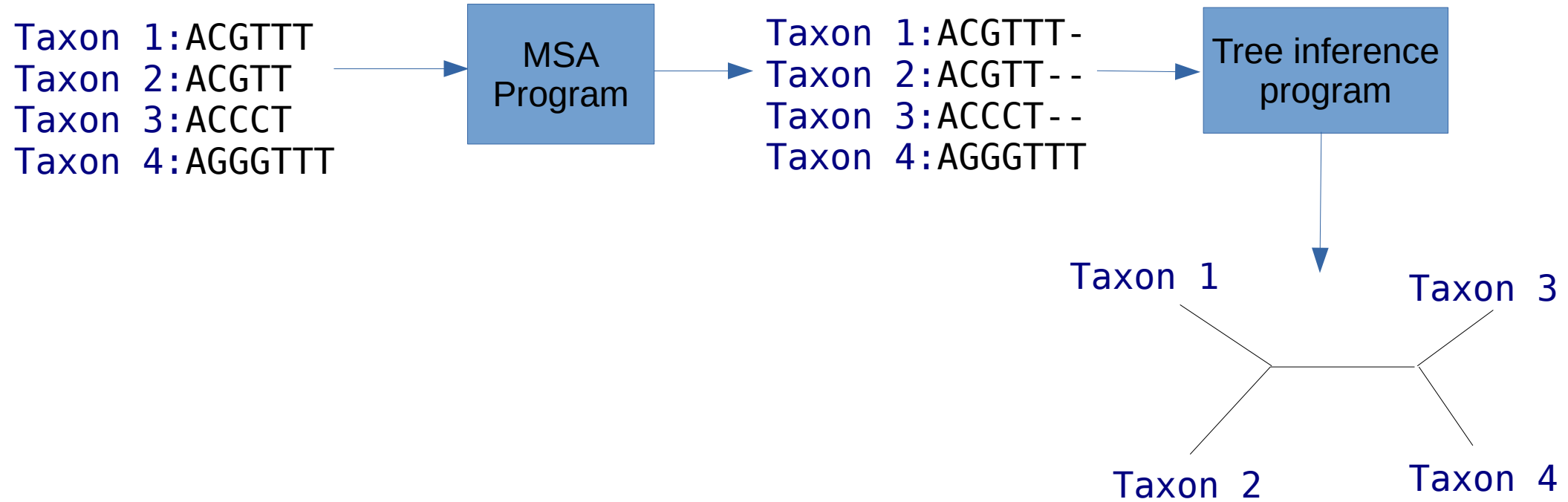
# A phylogeny



Taxon 1

Taxon 2

Taxon 3

Taxon 4

Taxon 5

**Phylogenetic trees are unrooted binary trees!!!**

# Tree Inference Pipeline

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

# Tree Inference Pipeline

# How many unrooted 4-taxon trees exist?

# How many unrooted 4-taxon trees exist?

# How do we select among them ?

# How do we select among them???



A   *Score:* 1.0   C
B                  D

A   *Score:* 3.0   B
C                  D

A   *Score:* 2.0   B
D                  C

We need scoring criteria!!!

# How do we select among them???



*Score:* 1.0

A — C
B — D

*Score:* 3.0

A — B
C — D

*Score:* 2.0

A — B
D — C

We need scoring criteria!!!
The currently most widely used criterion is maximum likelihood:
How likely is it that the tree, given a model of evolution, generated
the observed data?

# How do we select among them???

Score: 1.0

A            C

B            D

Maximum Likelihood tree

A    Score: 3.0    B

C            D

Score: 2.0

A            B

D            C

We need scoring criteria!!!
The currently most widely used criterion is maximum likelihood:
How likely is it that the tree, given a model of evolution, generated
the observed data?

# The number of trees

3 taxa → *1 tree*

# The number of trees



4 taxa → *3 trees*

# The number of trees

5 taxa → *15 trees*

# The number of trees

6 taxa → *105 trees*

# The number of trees explodes!



BANG !

# # possible trees with 2000 taxa

```
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: 30049638174211656151632910065681814981377232074237013089504954043012636525258308210827685996688247000464352735214265634288295
891502344600006314939691306329704360561848618774654822779912235368092334555631999108345976931267565250128998674331877528114019609916315223670306091217357097623798477054676677
77953247971826143852733382267277842507372528499166696875844035105795870206865058176870446663181237429010214385064324713609344916670211359697569403000666252646479269124551031
49423661955428241182776251148487582545812279142898011326489026740337612947127473676036267579086843169660718609847941818865957214557044744572288661729053583520744253688123124
0106613156948861960941195646736200342575241335277575085829161096422575727699767991408283343210161327401652830993803904592327690690035972919709940739349563486203899010742687
282297597465537710225767267684285801187722495010622181173405232082653973429622273525365905158656313832720311198419874675997386463182903203832523085979979922216101227215780805
248145831206844016760062393060097116167297155047284877996343375313489942303724373478791319890859537640701348494461138775725756952408702461720107874294978304622750525457066893 72
31941820644070689188400387059028977219751645449597582166213062050646177610994856637341681835489932907699338206780105243728461492403422961155182609778228619192627207129519895
893600995913097423307231638251842811033057101744115688430513186587754437630850031145111072383703970746518223204040615470827307862995754933103127520861670066079129801426 2230
0565123522718063811920933587265172862358902052001614436175607565428647142212661300443480708406750158924767316634153954057507447499490983149647303108041114018918497359128 11228
378770498848340542102420566424463860093899650857429619472690540152812375265109658152846997970379217112903556809818079169587951614159281049528179855847292534447846442 44359
9808531537204796814969465991768614533701051985928577157482455943377242369582576242663016946320482495182255939287403177623433881048604630975191556923871167513095213415098816
715464307862352606237786406838680424690252749113931927680261151599058260388673317293071367390340361863746398060576483647467027444672788088533707425442192272667774700332 9403
32010382880351126890262551830967919483586789293701637681753048206338943871497931152335698229625111163071482945992116208033026847620133356904410896681453 4615090515587758 1167
97700125639121511162374441704973717046040294811041148222864661319188219975713833683520725260552027698239746132184952492648970507903983602562556062898522888395613578874156576
6488999260873286612630642543260248979229113560071640573984516375245423376943755857384725455643975996042559146401122211447552355731762399730577471839565312741653229598 66575
9012941161239240722093250369673124884491553759210650560560154167207741592362408686676753482865129648887390597075788024733945340708481590116397727977474800417316 2687 00916728735
6121642268468160683198959801260376485615312781611689587215123123087600634733810972531184233396403909373783950668355787353078863586464005329949949060311874240290927792 72693
30032244537759572248734568915114585570783850541681667667425811301958063621907500790295031088209097271748136436898473971079932777700676301730617566538739726037771717300 8441
343940512366905554493248616508253995779503632670494784429349885317297348177797146567175511788763964340693324580763461107342143281950499009680874027397688914 7045174720555543
8969396668742601477241894693212902453317334188286773194653544113302100866575081713240342647580489218623663461607955943720516395154096949060486242309796947211169665961580041
78839922322646284984235227692639112436076750899767171896834593378907423463545571935306561537990027791626563636186197404859093823406223545976973313721365934371758559066 44439
6461328400113672302132448919921530438528154165963011854942836349090629162729564926584723003608555598989761916293248140094592489899468468623225866
01705514689056498372760039274870695504637888197416994290491028531804107765160072633847216389037822700138435995997302657271988043246643123597585552171969051392101022659 6324
78878309774053331315517629781528807186526031763272648280944993704562580999380530584976995700802893798080149029010052938984722799471678048216894241591182842576964647865 05731
5325178130233630729326516922103465842658944474649161238546897185079681729091390321828341111848213847677283165486532123173820041319905016780722201887049585687088890 3360
6930402937216038968917605587676955382318093705826257083898387409098468656634271397500013291835105943321729879825243707508272087959859437157667660155782699660343197 75623308
89899625878006280095609444169323779449554410336965862615562560106693903032038789709836737860870566414335851061116583145204245132085085899949323648316896711949516716195 67622
70709069738895888557597924666415365617235493018073940047605298017217713916876880002778519661730700612845173075825037335643102065112443730825229625040453160590741 34388187263
47791383066059093180025223100853401768402614015396169891920751471080337577088497401418345997539720598786820648791160649698581776011539720584982226989071813494326918018 21173
318806365391089368981171489135745668054280748517017585826663963357018935444983266976283509265792220174637219027311964175148994401007963687601782674710701994547321888 7832742
60889667243715747713420600093704251309893630537459784279980403132894172664922904257309583685344162156405572902820662240038632375263809102332698978388604237596231 5175262
6950798639868104294832333160267216555178120899264677804953741326387137408423885546538336158643451305439624281397279559725995110706314305992653202323270805768115 66690
4895866105220300573725298472118747827136713666058669271094875563974858489475910819727033878284439864486743456200958161930314727345961900499318424337975243662489363 321244850
59719925236625292493053462527641378534132089431289015237380925560459870909127666623296787033288820591349495800740744731433888007245323217473096597419671144414531271 32790205
1010047670101435063885795347844725538980154192331702751989618063515268254317319383292589193153016413054897231112866664654929719304792964328295567190928 8169209104233412 2007454
2420499078125489729634328295567190928816920910423341220074542420499078125489729634328295567190928816920910423341220074542420499078... (wrapped)
350341796875

Approximately 3.00 times 10^6328
```
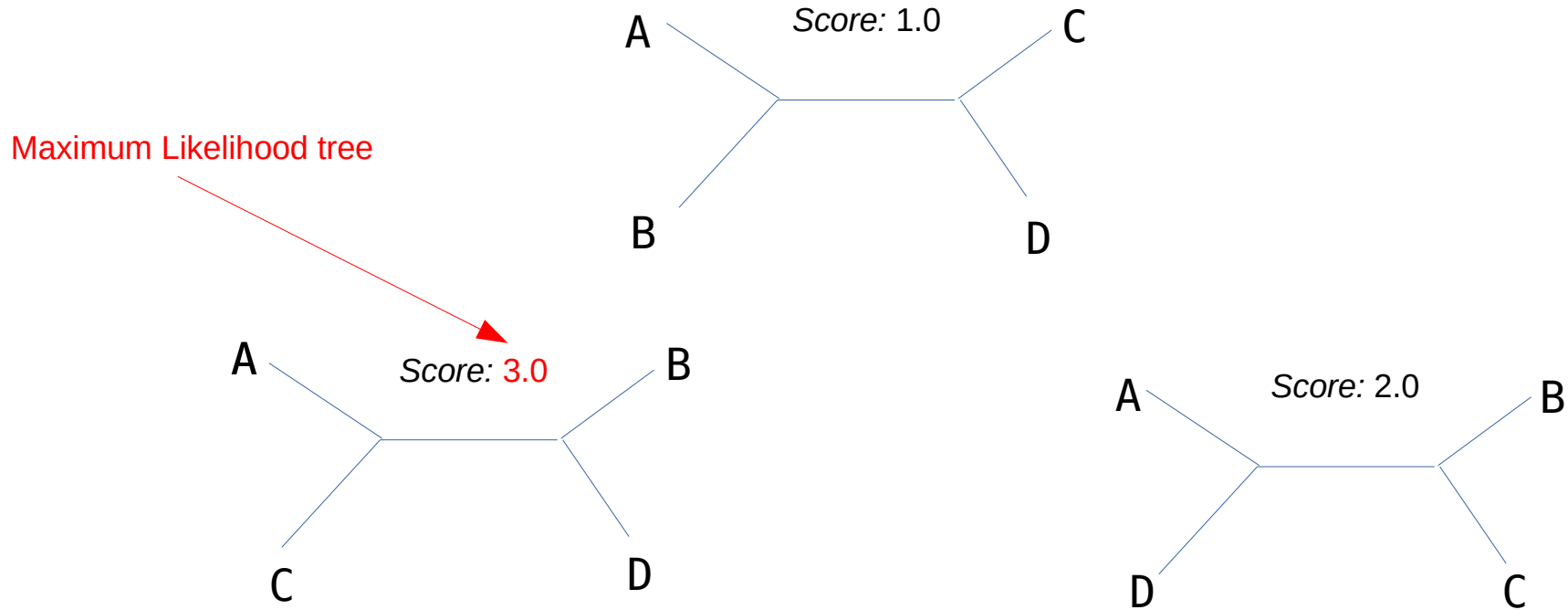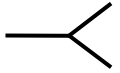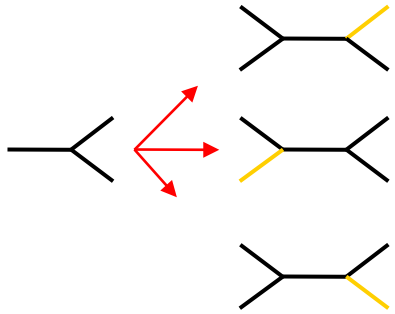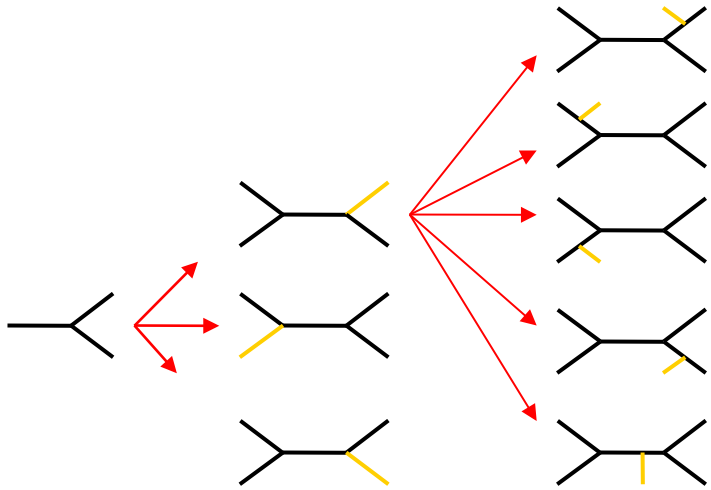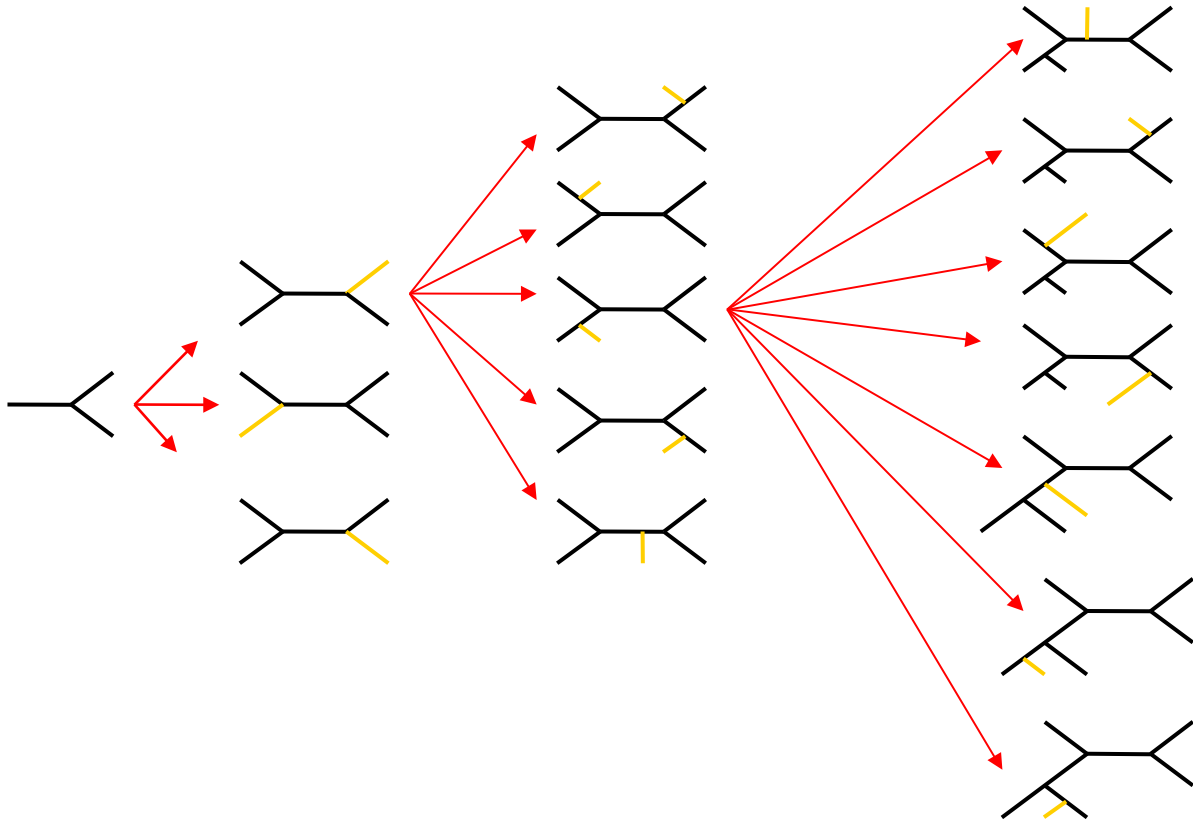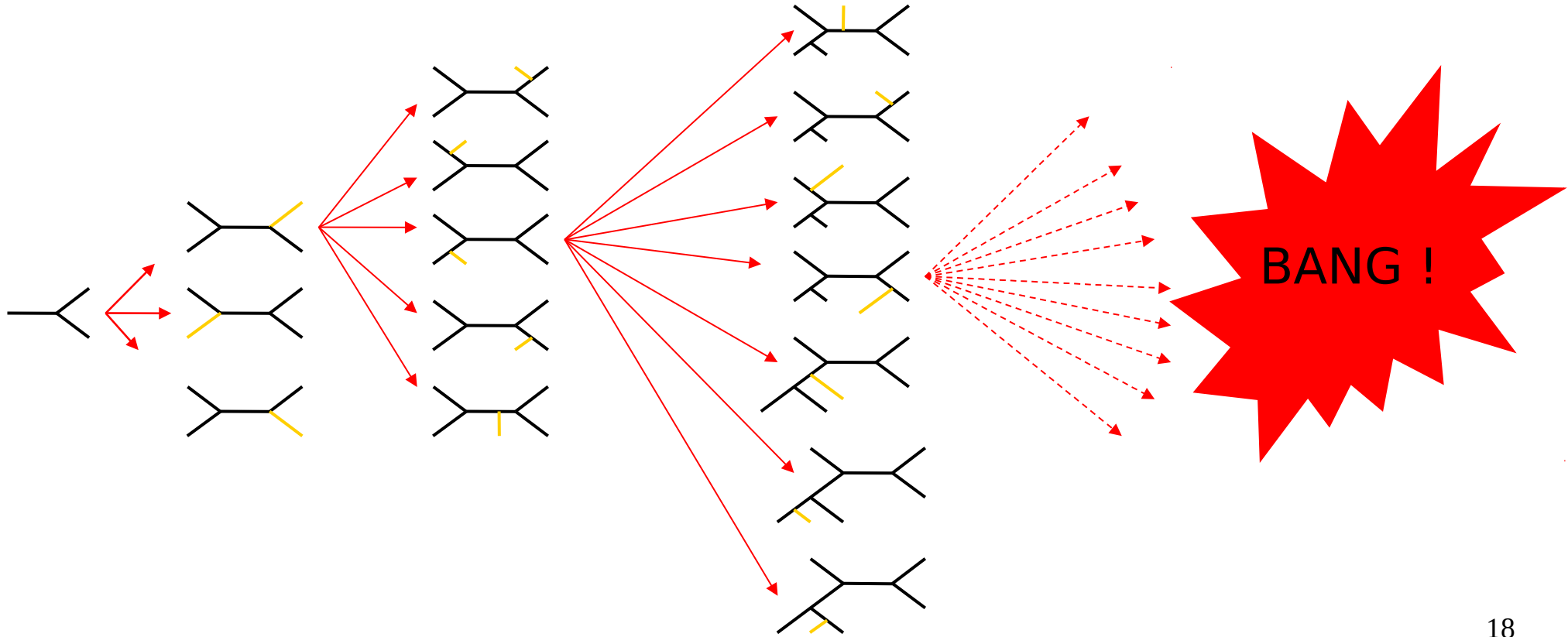
# Problem Complexity



Global maximum

search space

heuristic tree
search strategy

good scores

bad scores

# Problem Complexity

Global maximum
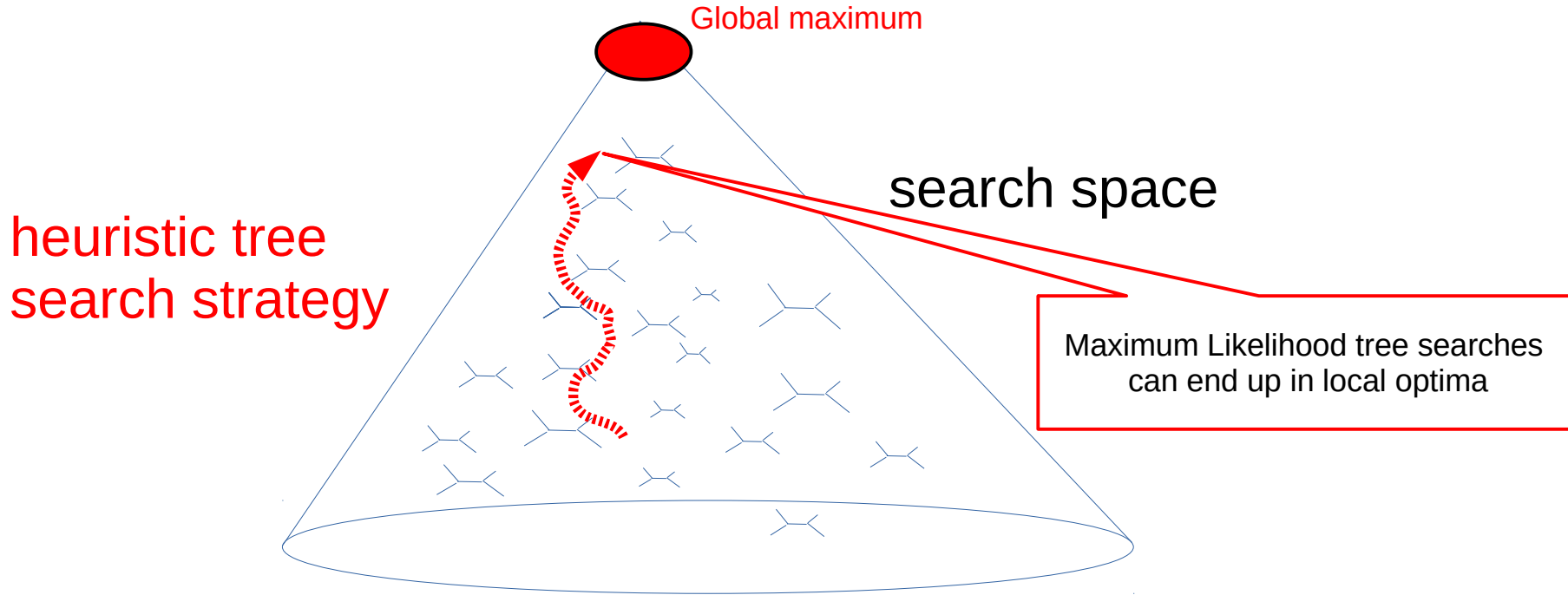
search space

good scores

heuristic tree
search strategy

bad scores

Finding the best tree under Maximum Likelihood is **NP-hard**!

# Problem Complexity

Global maximum

search space

heuristic tree
search strategy

Maximum Likelihood tree searches
can end up in local optima

# Maximum Likelihood

Length: m

Seq1
Seq2  Alignment
Seq3
Seq4

# Maximum Likelihood

Length: m

Seq1
Seq2
Seq3
Seq4

Alignment

|   | A | C | G | T |
|---|---|---|---|---|
| A |   |   |   |   |
| C |   |   |   |   |
| G |   |   |   |   |
| T |   |   |   |   |

Substitution
model

# Nucleotide Substitution Models



We model evolution as time-reversible Markov Process!

# Maximum Likelihood

Length: m

Seq1
Seq2    Alignment
Seq3
Seq4

A C G T

A
C    Substitution
G    model
T

Commonly denoted as $Q$ matrix: transition probs for time $dt$, for time $t$: $P(t)=e^{Qt}$

# Maximum Likelihood

Length: m

Seq1
Seq2      Alignment
Seq3
Seq4

A  C  G  T

A
C      Substitution
G      model
T

Prior probabilities,
Empirical base frequencies

$\pi_A \ \pi_C \ \pi_G \ \pi_T$

# Maximum Likelihood

Length: m

Seq1
Seq2
Seq3
Seq4

Alignment

$$\begin{array}{c|c|c|c|c} & A & C & G & T \\\hline A & & & & \\\hline C & & & & \\\hline G & & & & \\\hline T & & & & \end{array}$$

Substitution model

Prior probabilities,
Empirical base frequencies

$$\pi_A \ \pi_C \ \pi_G \ \pi_T$$

Seq 1   b1                                    b3   Seq 3
                        b5
b2                                                  b4
Seq 2                                              Seq 4

# Maximum Likelihood



Length: m

Seq1
Seq2
Seq3
Seq4

Alignment

A C G T

A
C
G
T

Substitution model

Prior probabilities,
Empirical base frequencies

$\pi_A \ \pi_C \ \pi_G \ \pi_T$

Seq 1

b1

Seq 3

b3

b5

b2

b4

Seq 2

Seq 4

virtual root: vr

# Maximum Likelihood

Length: m

Seq1
Seq2
Seq3
Seq4

Alignment

A C G T

A
C
G
T

Substitution
model

Prior probabilities,
Empirical base frequencies

$\pi_A$ $\pi_C$ $\pi_G$ $\pi_T$

Seq 1    b1                              b3    Seq 3

vr  b5

b2                                            b4

Seq 2                                              Seq 4

| P(A) | P(C) | P(G) | P(T) |
|------|------|------|------|
|      |      |      |      |

| P(A) | P(C) | P(G) | P(T) |
|------|------|------|------|
|      |      |      |      |

m

# Maximum Likelihood

# Maximum Likelihood

Length: m

Seq1
Seq2
Seq3
Seq4

Alignment

A C G T

A
C
G
T

Substitution
model

Prior probabilities,
Empirical base frequencies

$\pi_A$ $\pi_C$ $\pi_G$ $\pi_T$

Seq 1

3

**Floating-point & memory intensive**

b2

Seq 2

b4

Seq 4

| P(A) | P(C) |
|------|------|
|      |      |

| P(A) | P(C) | P(G) | P(T) |
|------|------|------|------|
|      |      |      |      |

m

# Maximum Likelihood

Length: m

A C G T

Prior probabilities,
Empirical base frequencies

Seq1
Seq2
Seq3
Seq4

Alignment

|   | A | C | G | T |
|---|---|---|---|---|
| A |   |   |   |   |
| C |   |   |   |   |
| G |   |   |   |   |
| T |   |   |   |   |

Substitution
model

$\pi_A \ \pi_C \ \pi_G \ \pi_T$

Seq 1

Sites evolve independently →
we can compute per-site likelihoods in
parallel :-)

b2

Seq 2

b4

Seq 4

| P(A) | P(C) | P(G) | P(T) |
|------|------|------|------|
|      |      |      |      |

| P(A) | P(C) | P(G) | P(T) |
|------|------|------|------|
|      |      |      |      |

m

# Outline

- Introduction to Phylogenetic Inference - *Alexandros*

- **The RAxML Search Algorithm - *Alexandros***

- Improvements in RAxML **N**ext **G**eneration - *Alexey*

- Tutorial - *Alexey*

# How does it work?

Compute randomized stepwise addition order
Maximum Parsimony tree

# How does it work?

Compute randomized stepwise addition order
Maximum Parsimony tree

Advantage of RAxML: search starts
from distinct point in search space
every time

# How does it work?

Compute randomized stepwise addition order
Maximum Parsimony tree

Advantage of RAxML: search starts
from distinct point in search space
every time

Alternatively, we can start from
a completely random tree

# Starting Trees

# Starting Trees

# How does it work?

Compute randomized stepwise addition order
Maximum Parsimony tree

Apply lazy subtree rearrangements

# How does it work?

Compute randomized stepwise addition order
Maximum Parsimony tree

Apply lazy subtree rearrangements

Iterate while tree improves

# Subtree Pruning & Re-Grafting

# Subtree Pruning & Re-Grafting

# Subtree Pruning & Re-Grafting

**ST1**

**ST2**

**+1**

re-graft

**ST6**

**ST3**

**ST5**

**ST4**

# Subtree Pruning & Re-Grafting

# Subtree Pruning & Re-Grafting

# Subtree Pruning & Re-Grafting

# Subtree Pruning & Re-Grafting

ST1

ST2

**+2**

ST3

ST5

**ST6**

ST4

# Subtree Pruning & Re-Grafting

**ST1**

**ST2**

**+2**

**ST3**

**ST5**

**ST6**

**ST4**

# Subtree Pruning & Re-Grafting



ST1

ST2

+2

ST5

ST6

ST4

This search parameter can be explicitly set in RAxML-NG via `--spr-radius`

# Subtree Pruning & Re-Grafting



ST1

ST2

+2

ST5

ST6

ST4

If you don't set it RAxML-NG will auto-tune it

# Subtree Pruning & Re-Grafting



**Optimize all branches!**

ST1

ST2

ST3

ST4

ST5

ST6

# Subtree Pruning & Re-Grafting



ST1

ST2

Do we need to optimize
all branches?

ST3

ST5

ST6

ST4

# *Lazy* Subtree Pruning & Re-Grafting

# *Lazy* Subtree Pruning & Re-Grafting



"FAST" SPR round:
no branches are optimized

55

# *Lazy* Subtree Pruning & Re-Grafting

ST1

ST2

ST3

ST5

ST4

ST6

"SLOW" SPR round:
three adjacent branches are optimized

# Outline

- Introduction to Phylogenetic Inference - *Alexandros*

- The RAxML Search Algorithm - *Alexandros*

- **Improvements in RAxML Next Generation - *Alexey***

- Tutorial - *Alexey*

# Evolution of RAxML(-NG)



58

# RAxML-NG design goals

- Full rewrite of RAxML

  – Search heuristic largely unchanged (as of v1.0)

- Improve maintainability & enable code reuse

- Eliminate known bugs & bottlenecks

- Improve user experience

  – by default, "do the right thing"

# RAxML-NG & family

# Improvements & new features

- Flexible evolutionary models
  - All "classical" + custom DNA models (eg. DNA010010 = HKY)
  - Per-partition rate heterogeneity (incl. FreeRate)
  - Proportional branch lengths
- Phylogenetic terrace detection (Biczok '17)
- Transfer bootstrap support metric (Lemoine '18)
- Energy monitoring

# Performance & scalability

- Checkpointing
- Advanced load balancing    from ExaML
- Binary alignment format
- "Site repeats" optimization (Kobert '17)
  - 10-60% speedup + memory savings
- Flexible and user-friendly parallelization

# Parallelization: hardware

raxmlHPC

   raxmlHPC-SSE

raxmlHPC-AVX

| CPU | Vectorization (SSE, AVX ...) |
|---|---|

raxml-ng

raxmlHPC-PTHREADS

raxmlHPC-PTHREADS-SSE

| Desktop | Multi-threading |
|---|---|

raxmlHPC-HYBRID

   raxmlHPC-MPI
raxmlHPC-MPI-AVX2

| Cluster | MPI |
|---|---|

raxml-ng-mpi

# Parallelization: software

Alignment



sites

4 threads

T1, T2, T3, T4

4 searches

e.g. from 4 starting trees

## Fine-grained

Search 1 | T1 | T2 | T3 | T4 |
Search 2 | T1 | T2 | T3 | T4 |
Search 3 | T1 | T2 | T3 | T4 |
Search 4 | T1 | T2 | T3 | T4 |

## Coarse-grained

Search 1 — T1
Search 2 — T2
Search 3 — T3
Search 4 — T4

## Mixed/hybrid

Search 1 — T1 | T2
Search 2 — T1 | T2
Search 3 — T3 | T4
Search 4 — T3 | T4

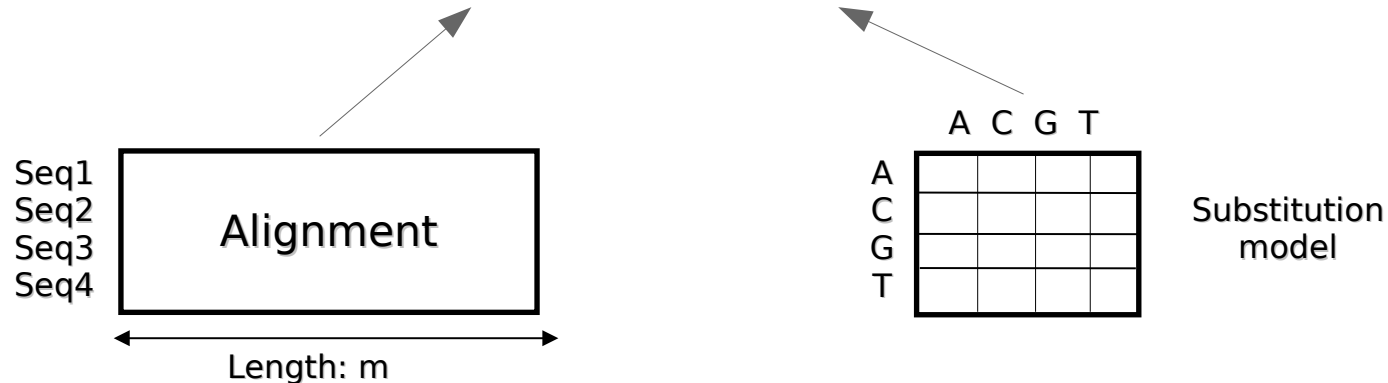**New in v1.0:** Full native support and automatic configuration!

# Outline

- Introduction to Phylogenetic Inference - *Alexandros*

- The RAxML Search Algorithm - *Alexandros*

- Improvements in RAxML Next Generation - *Alexey*

- **Tutorial - *Alexey***

# Quick start: ML tree search

- Default command: --search
  - 20 starting trees (10 random + 10 parsimony)
  - Pick the best-scoring one

```
$ raxml-ng --msa prim.phy --model GTR+G
```



Seq1
Seq2
Seq3
Seq4

Alignment

Length: m

A C G T

A
C
G
T

Substitution
model

# ML tree search: Output

```
Analysis options:
  run mode: ML tree search
  start tree(s): random (10) + parsimony (10)
...
Starting ML tree search with 20 distinct starting trees

[00:00:00 -7871.515760] Initial branch length optimization
...
[00:00:00 -5736.644605] FAST spr round 1 (radius: 5)
...
[00:00:00 -5709.394601] SLOW spr round 1 (radius: 5)
...
[00:00:00] ML tree search #1, logLikelihood: -5708.979717
...
[00:00:07] ML tree search #20, logLikelihood: -5709.014076


Final LogLikelihood: -5708.923977

Best ML tree saved to: /home/alexey/test/prim.phy.raxml.bestTree
Optimized model saved to: /home/alexey/test/prim.phy.raxml.bestModel
```

--log info to hide search progress

67

# Tree with support values

- All-in-one mode: <span style="color:red">--all</span>
  - ML tree search (as before)
  - Bootstrapping with autoMRE convergence test
  - Compute support values + map on ML tree

```
$ raxml-ng --all --msa prim.phy --model GTR+G --prefix A1 --seed 1
```

Output files: A1.raxml.*

Fixed RNG seed
for better reproducibility

# Tree with support values: Output

```
Analysis options:
  run mode: ML tree search + bootstrapping (Felsenstein Bootstrap)
  start tree(s): random (10) + parsimony (10)
  bootstrap replicates: max: 1000 + bootstopping (autoMRE, cutoff: 0.030000)

Starting ML tree search with 20 distinct starting trees
...
[00:00:02] ML tree search completed, best tree logLH: -5708.926130

[00:00:02] Starting bootstrapping analysis with 1000 replicates.
...
[00:00:14] Bootstrapping converged after 100 replicates.

Best ML tree with Felsenstein bootstrap (FBP) support values saved to:
/home/alexey/test/A1.raxml.support
...
Bootstrap trees saved to: /home/alexey/test/A1.raxml.bootstraps
```

# Customize analysis

- Starting trees: --tree

`--tree rand{1}`  `--tree pars{50},rand{50}`  `--tree user.nw`

- Bootstrap replicates: --bs-trees

`--bs-trees 100`  `--bs-trees autoMRE{500}`  `--bs-trees bs.bw`

- Branch length linkage mode: --brlen

`--brlen linked`  `--brlen scaled`  `--brlen unlinked`

- Branch length limits: --blmin / --blmax

`--blmin 1e-9`  `--blmax 10`

# Tree likelihood evaluation

- Optimize free model parameters and branch lengths on a fixed topology: --evaluate

```
$ raxml-ng --evaluate --msa prim.phy --tree A1.raxml.bestTree --model GTR+G --prefix E1
```

```
Evaluating 1 trees

[00:00:00] Tree #1, initial LogLikelihood: -6419.996676    ←

[00:00:00 -6419.996676] Initial branch length optimization
[00:00:00 -6276.983451] Model parameter optimization (eps = 0.100000)

[00:00:00] Tree #1, final logLikelihood: -5709.005148    ←
...
Best ML tree saved to: /home/alexey/test/E1.raxml.bestTree
Optimized model saved to: /home/alexey/test/E1.raxml.bestModel
```

# Tree likelihood evaluation (2)

- Compute and print tree log-likelihood: <span style="color:red">--loglh</span>
  - No branch length optimization
  - No model optimization
  - No output files created

```
$ raxml-ng --loglh --msa prim.phy --tree A1.raxml.bestTree --model GTR{1/2/3/4/5/6}+G{0.5}

Final LogLikelihood: -6334.023267
```

# Check & parse

- Check alignment for format errors: <span style="color:red">--check</span>

```
$ raxml-ng --check --msa prim.phy —model GTR+G
```

- Compress alignment into binary file: <span style="color:red">--parse</span>

```
$ raxml-ng --parse --msa prim.phy --model GTR+G --prefix prim
```

- ...which can be then used in parallel jobs

```
$ raxml-ng --search --msa prim.raxml.rba --prefix S1
```

# Parallelization tuning

- Fully automatic (default) → heuristic-based

```
$ raxml-ng --msa prim.rba

System: Intel(R) Xeon(R) CPU E5-2630 v3 @ 2.40GHz, 16 cores, 62 GB RAM
...
Analysis options:
  run mode: ML tree search
  start tree(s): random (10) + parsimony (10)
...
  parallelization: coarse-grained (auto), PTHREADS (auto)
...
[00:00:00] Alignment comprises 12 taxa, 1 partitions and 413 patterns
...
Parallelization scheme autoconfig: 16 worker(s) x 1 thread(s)
...
[00:00:00] Data distribution: max. partitions/sites/weight per thread: 1 / 413 / 6608
[00:00:00] Data distribution: max. searches per worker: 2
```

74

# Parallelization tuning

- Automatic with upper limits

```
$ raxml-ng --msa prim.rba --threads auto{16} --workers auto{2}
```

- Manual

```
$ raxml-ng --msa prim.rba --threads 16 --workers 2
```

- Also works with MPI

```
$ mpirun -n 4 raxml-ng-mpi --msa prim.rba --threads 16 --workers 8
```

4 ranks * 16 threads = 64 = 8 workers * 8 threads

# Energy monitoring

- New in RAxML-NG v1.0: energy usage report
  - Measured with Intel RAPL → CPU+DRAM only
  - Supported on Linux systems only
  - To disable, add: --extra energy-off
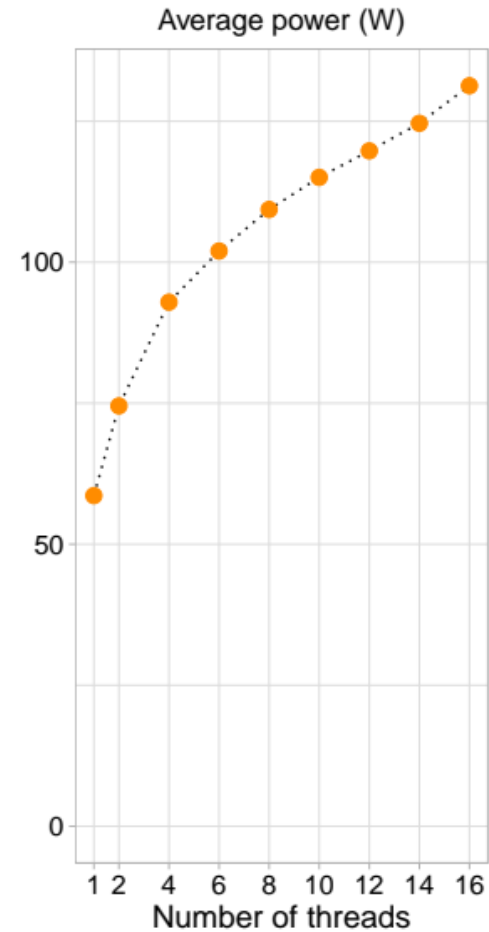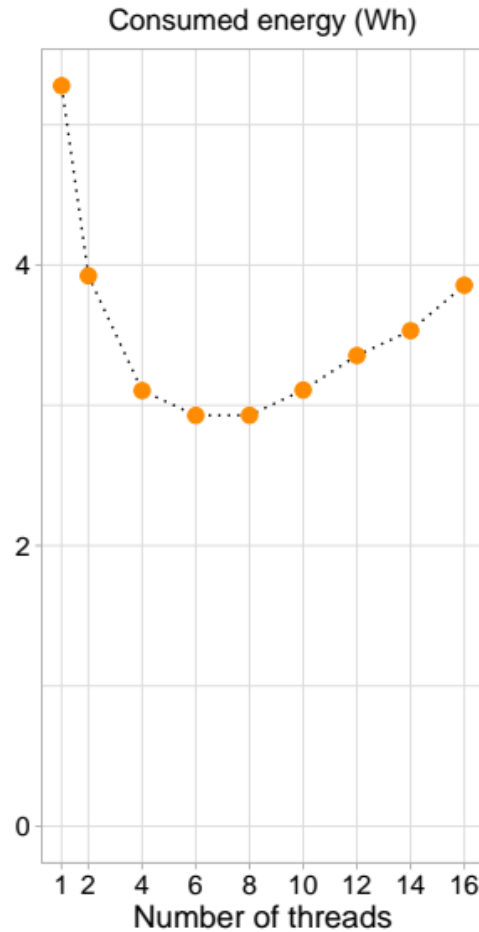
# Energy monitoring

- New in RAxML-NG v1.0: energy usage report
  - Measured with Intel RAPL → CPU+DRAM only
  - Supported on Linux systems only
  - To disable, add: --extra energy-off
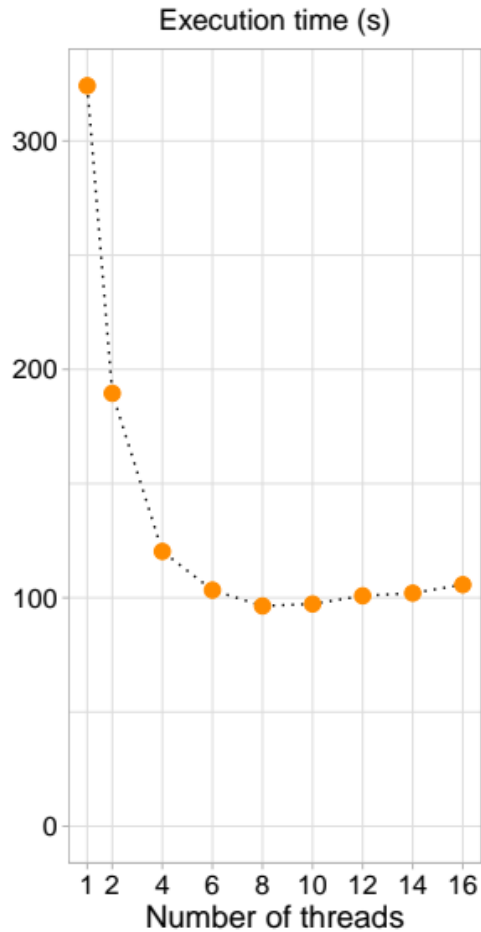
```
Elapsed time: 42846.287 seconds

Consumed energy: 162370.469 Wh (= 812 km in an electric car, or 4059 km with an e-scooter!)
```

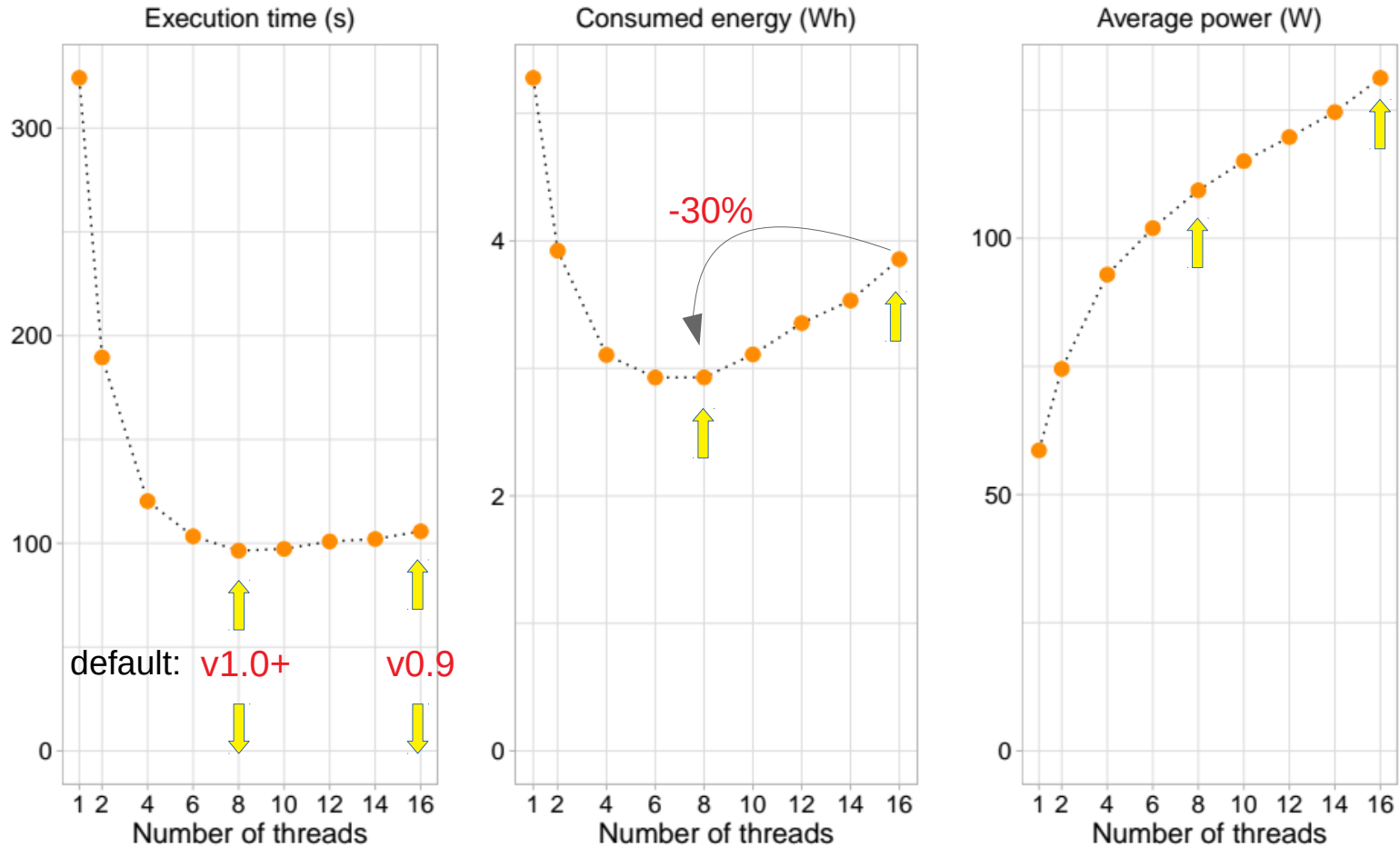Single tree search (96 nodes x 12h): >160 kWh
My apartment per month: ~100 kWh

# You can't improve what you can't measure!

# You can't improve what you can't measure!

# RAxML-NG availability

- Web services
  - Vital-IT: https://raxml-ng.vital-it.ch/#/ → free, for small datasets
  - CIPRES: http://www.phylo.org/ → registration required
- Source+binaries for Linux & macOS
  - GitHub: https://github.com/amkozlov/raxml-ng
- Conda: https://anaconda.org/bioconda/raxml-ng
- GUI: https://github.com/AntonelliLab/raxmlGUI

# Where to get help?

- Documentation

  https://github.com/amkozlov/raxml-ng/wiki

- Tutorial

  https://github.com/amkozlov/raxml-ng/wiki/Tutorial

- User support group

  https://groups.google.com/forum/#!forum/raxml

# Q & A

# References

- Barbera et al. (2018) **EPA-ng: Massively Parallel Evolutionary Placement of Genetic Sequences.** *Systematic Biology*, syy054, https://doi.org/10.1093/sysbio/syy054
- Biczok et al. (2017) **Two C++ libraries for counting trees on a phylogenetic terrace.** Bioinformatics. https://doi.org/10.1093/bioinformatics/bty384
- Kobert et al. (2017) Efficient detection of repeating sites to accelerate phylogenetic likelihood calculations. *Syst. Biol https://doi.org/10.1093/sysbio/syw075*
- Kozlov et al. (2019) **RAxML-NG: A fast, scalable, and user-friendly tool for maximum likelihood phylogenetic inference**. Bioinformatics*. https://doi.org/10.1093/bioinformatics/btz305*
- Kozlov, Aberer, Stamatakis (2015) **ExaML version 3: a tool for phylogenomic analyses on supercomputers**. *Bioinformatics*. https://doi.org/10.1093/bioinformatics/btv184
- Lemoine et al. (2018) **Renewing Felsensteinen phylogenetic bootstrap in the era of big data.** *Nature. https://doi.org/10.1038/s41586-018-0043-0*
- Morel, Kozlov, Stamatakis (2019) **ParGenes: a tool for massively parallel model selection and phylogenetic tree inference on thousands of genes.** *Bioinformatics.* https://doi.org/10.1093/bioinformatics/bty839
- Morel et al (2020). **GeneRax: A tool for species tree-aware maximum likelihood based gene tree inference under gene duplication, transfer, and loss**. *MBE*. https://doi.org/10.1093/molbev/msaa141
- Stamatakis (2006) **RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models**. *Bioinformatics*, https://doi.org/10.1093/bioinformatics/btl446
- Stamatakis (2014) **RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies.** *Bioinformatics*, https://doi.org/10.1093/bioinformatics/btu033
- Stamatakis and Aberer (2013) **Novel parallelization schemes for large-scale likelihood-based phylogenetic inference.** *In Parallel Distributed Processing (IPDPS)* https://doi.org/10.1109/IPDPS.2013.70