



Let's have a drink now!



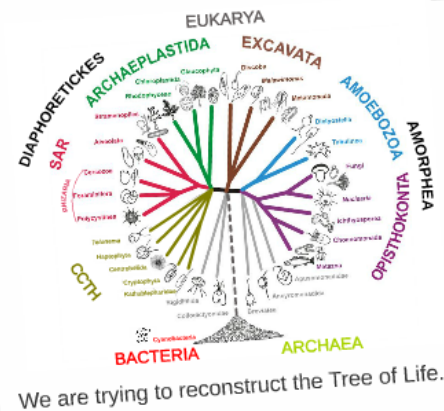
***HTC Bayesian Phylogenomics
on the Zenobe Tier-1 supercomputer***

Denis BAURAIN

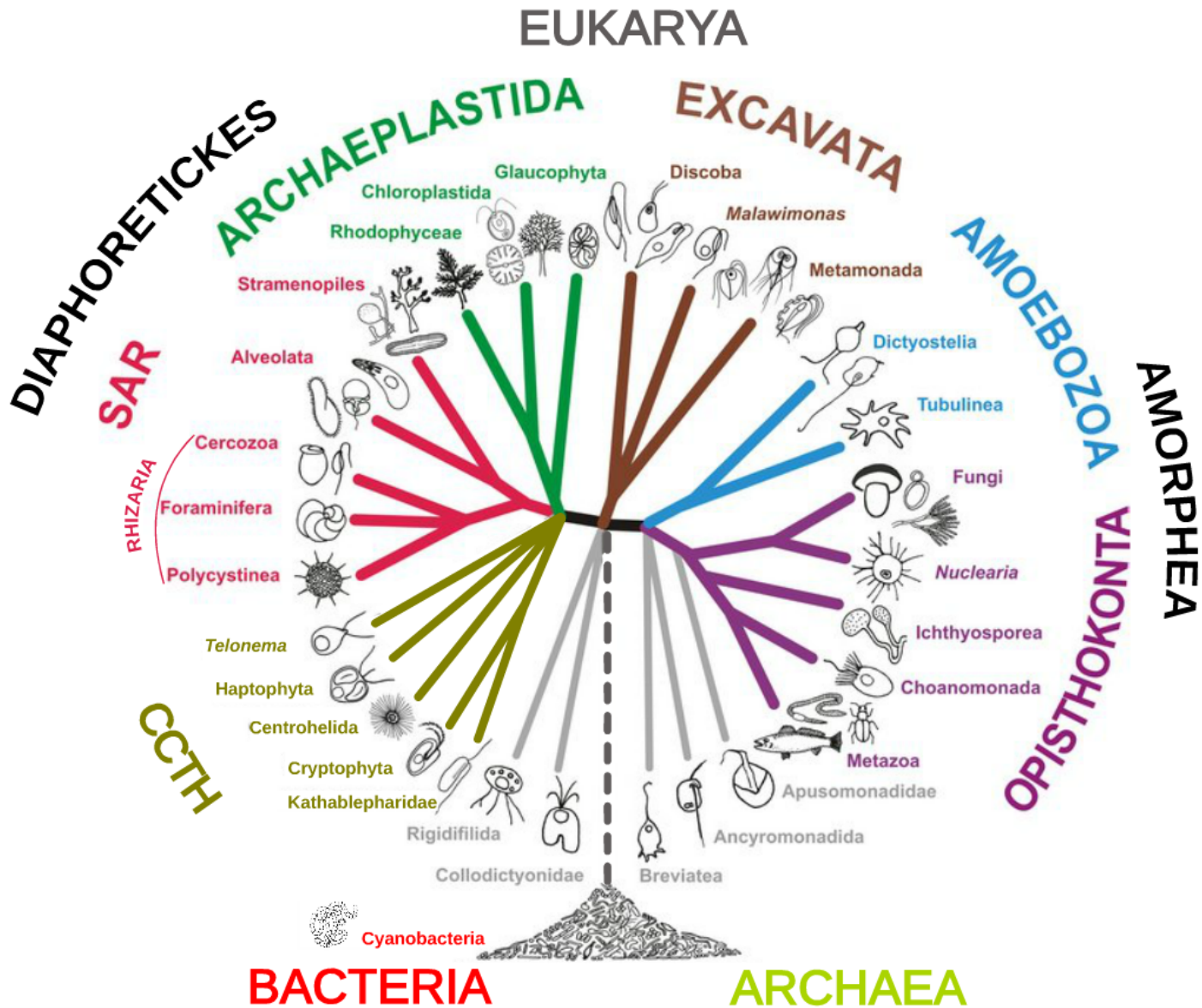
InBioS / Dept of Life Sciences
University of Liège

CÉCI scientific day / Apr 22, 2016

Background

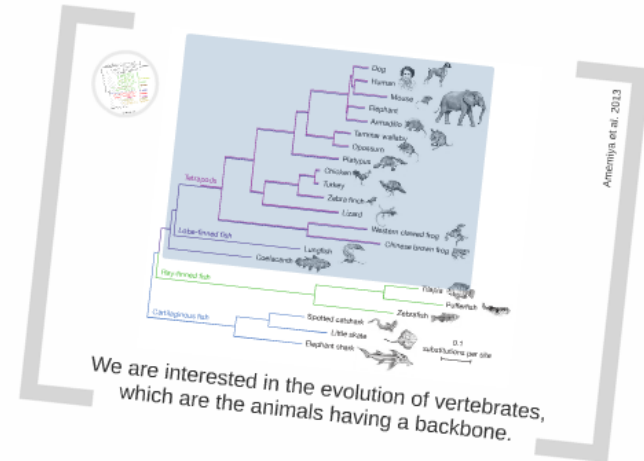


What is the general problem people in your field of research are trying to solve?

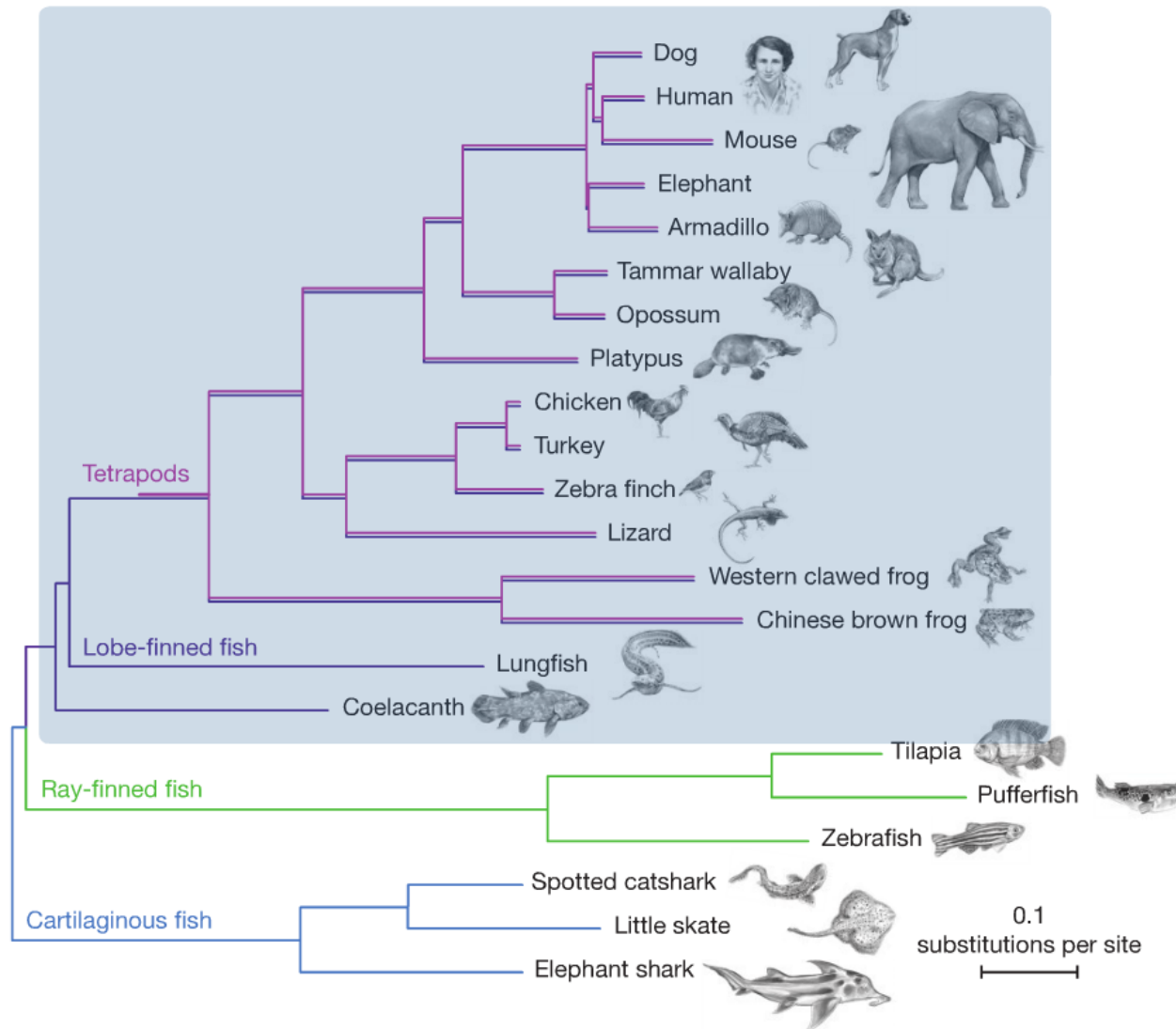
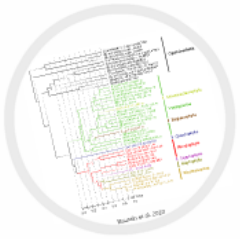


We are trying to reconstruct the Tree of Life.

Objectives

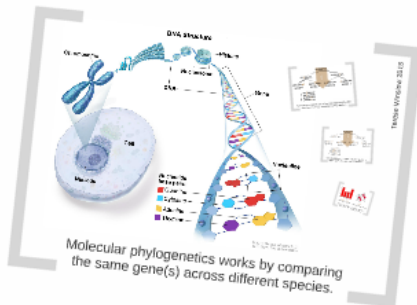
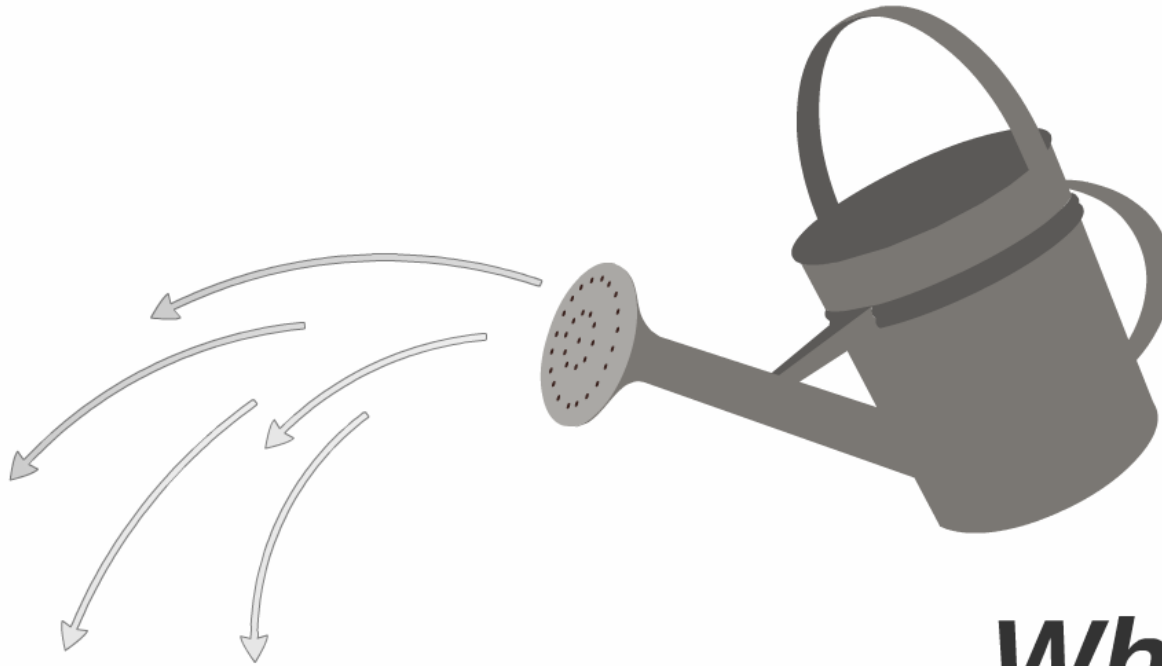
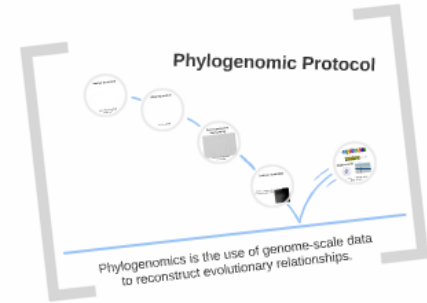


What is the specific problem in that context that you are trying to solve and why does that matter?

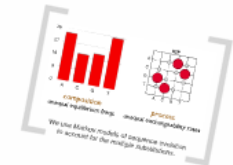
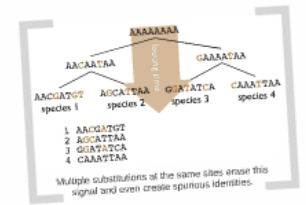
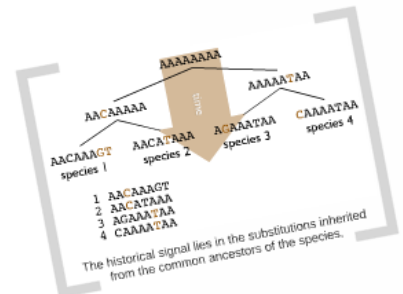
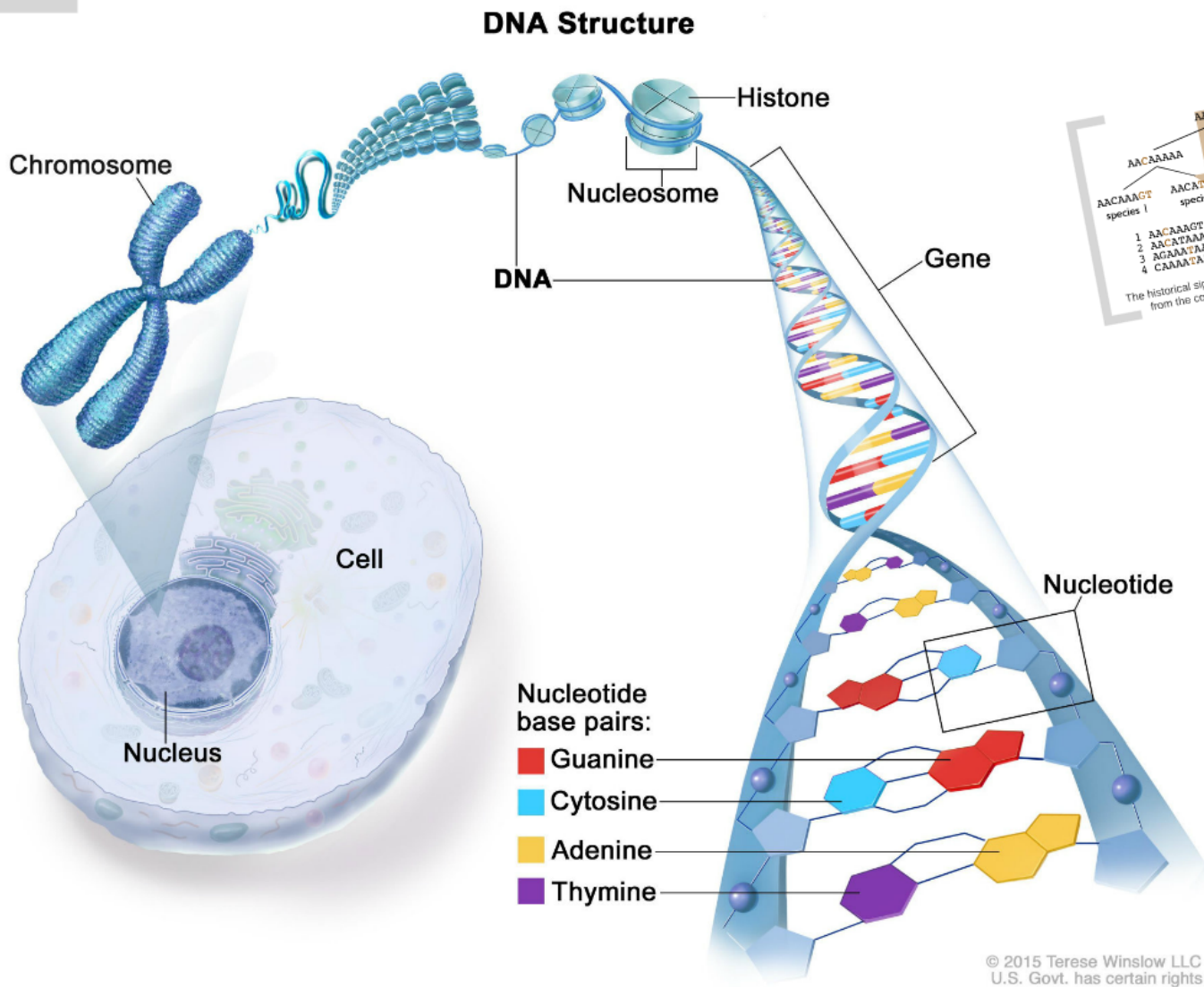


We are interested in the evolution of vertebrates, which are the animals having a backbone.

Methods

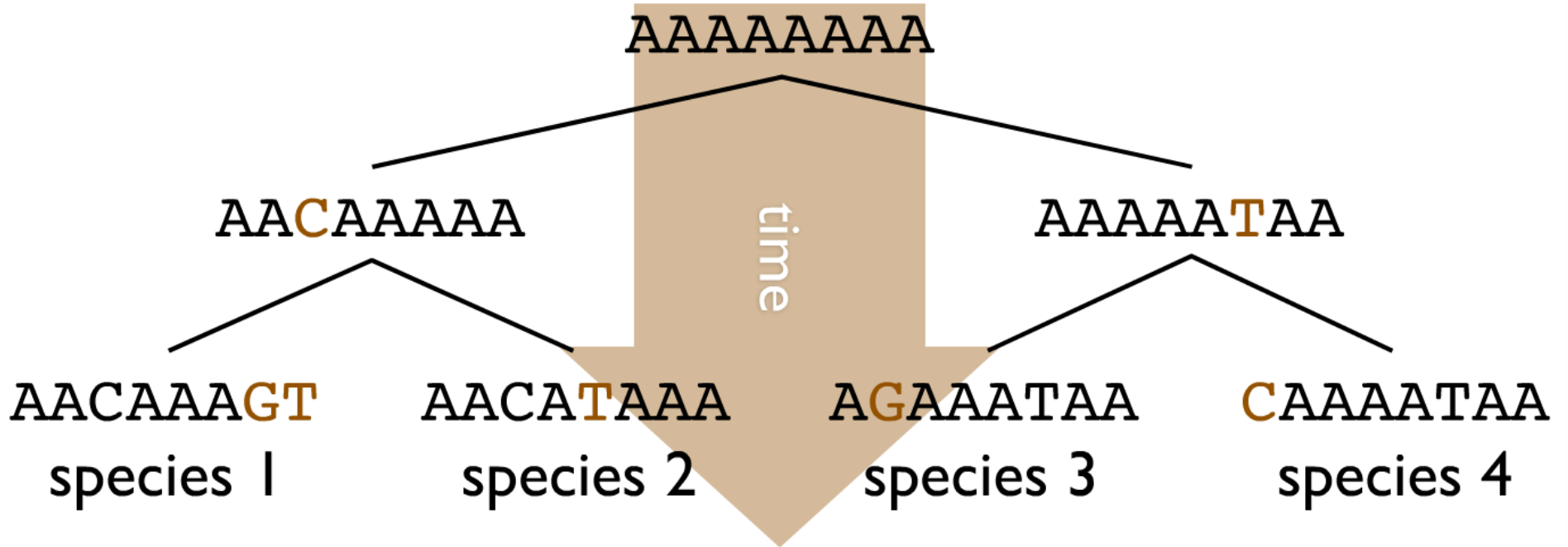


What strategy are you implementing to solve that problem?



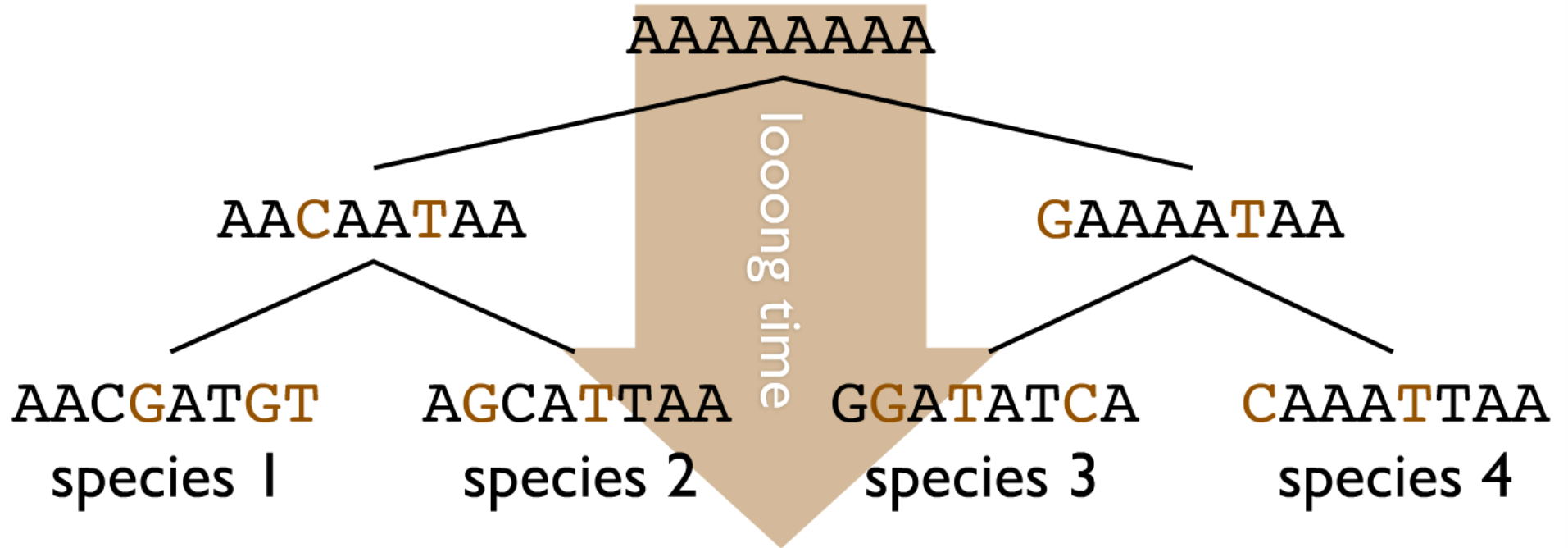
© 2015 Terese Winslow LLC
U.S. Govt. has certain rights

Molecular phylogenetics works by comparing the same gene(s) across different species.



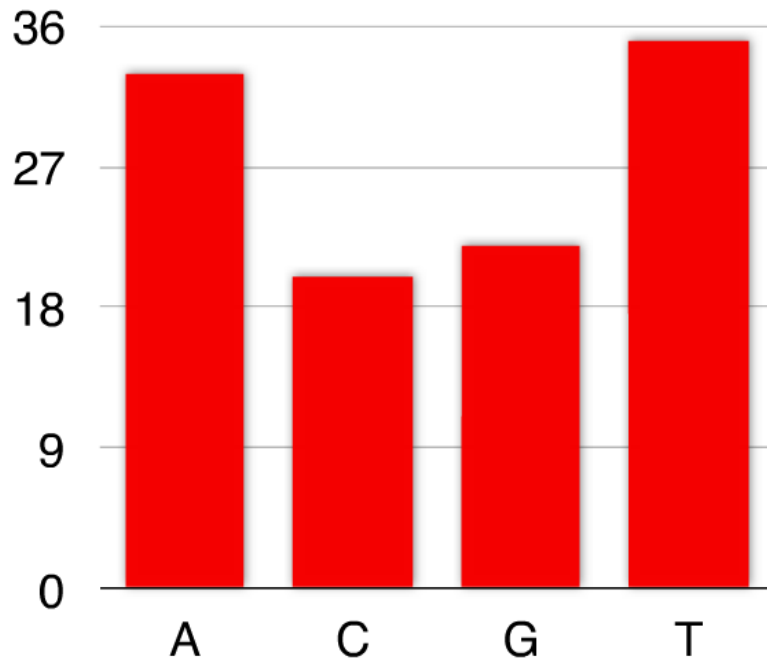
- 1 AACAAAGT
- 2 AACATAAA
- 3 AGAAATAA
- 4 CAAAATAA

The historical signal lies in the substitutions inherited from the common ancestors of the species.



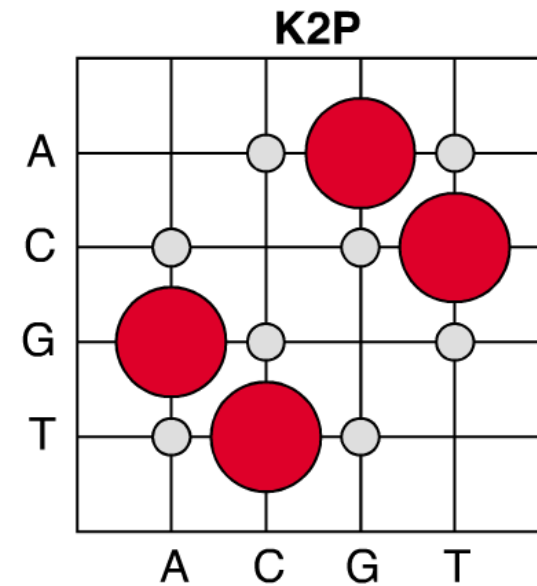
- 1 AACGATGT
- 2 AGCATTAA
- 3 GGATATCA
- 4 CAAATTAA

Multiple substitutions at the same sites erase this signal and even create spurious identities.



composition

unequal equilibrium freqs

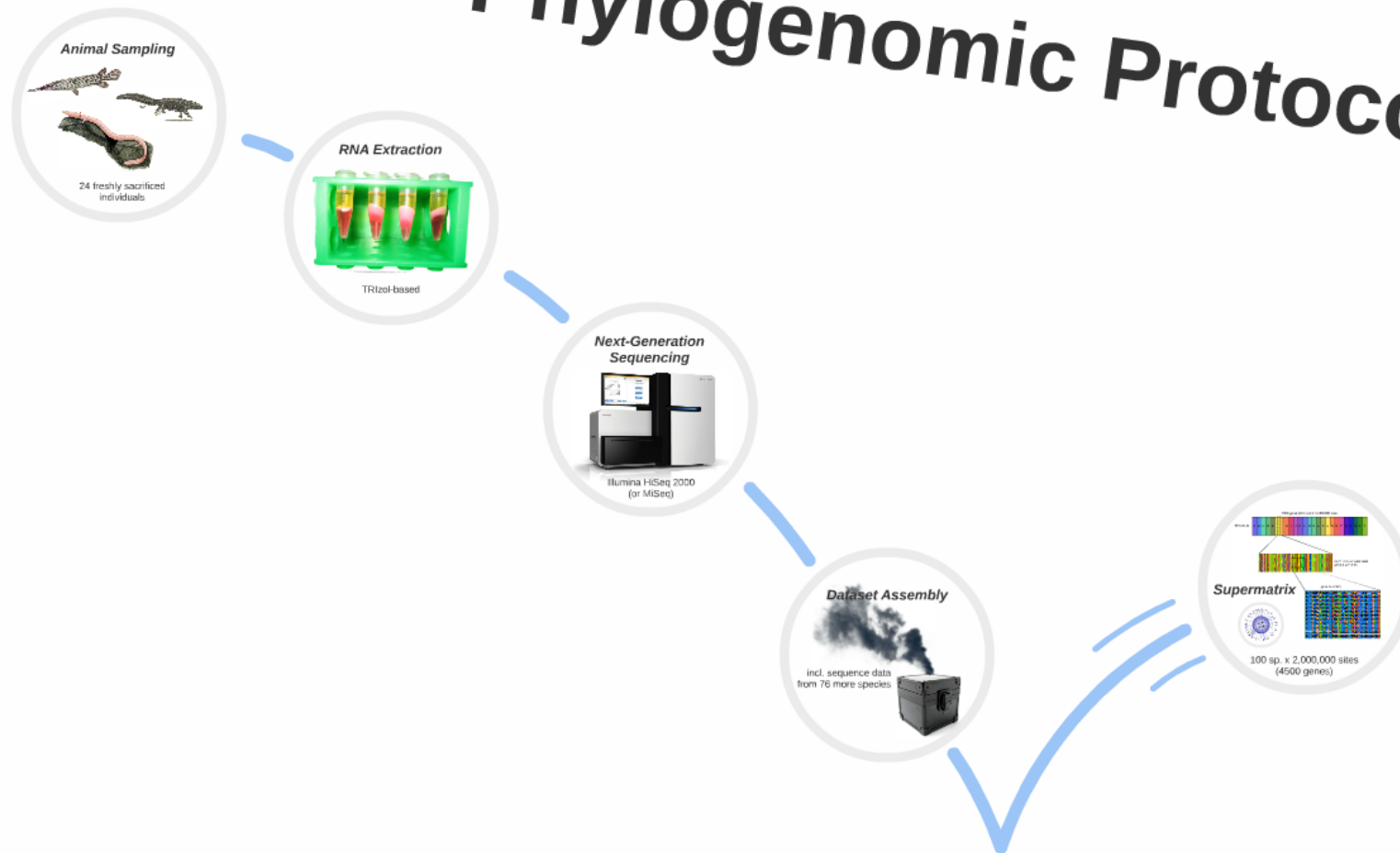


process

unequal exchangeability rates

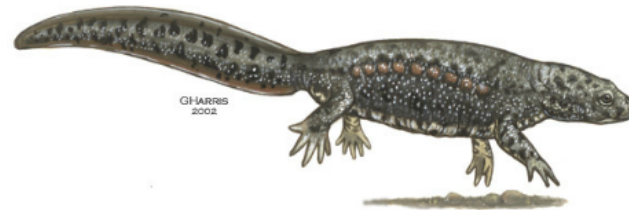
We use Markov models of sequence evolution to account for the multiple substitutions.

Phylogenomic Protocol



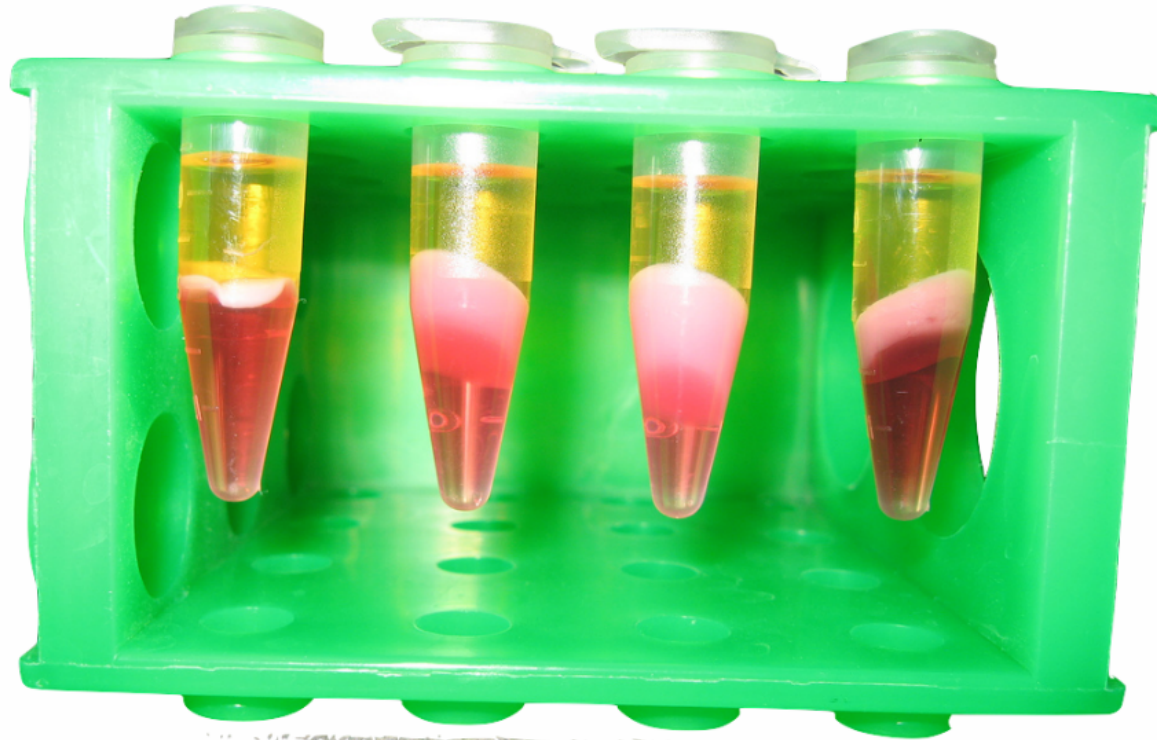
Phylogenomics is the use of genome-scale data to reconstruct evolutionary relationships.

Animal Sampling



24 freshly sacrificed
individuals

RNA Extraction



TRIZOL-based

Next-Generation Sequencing



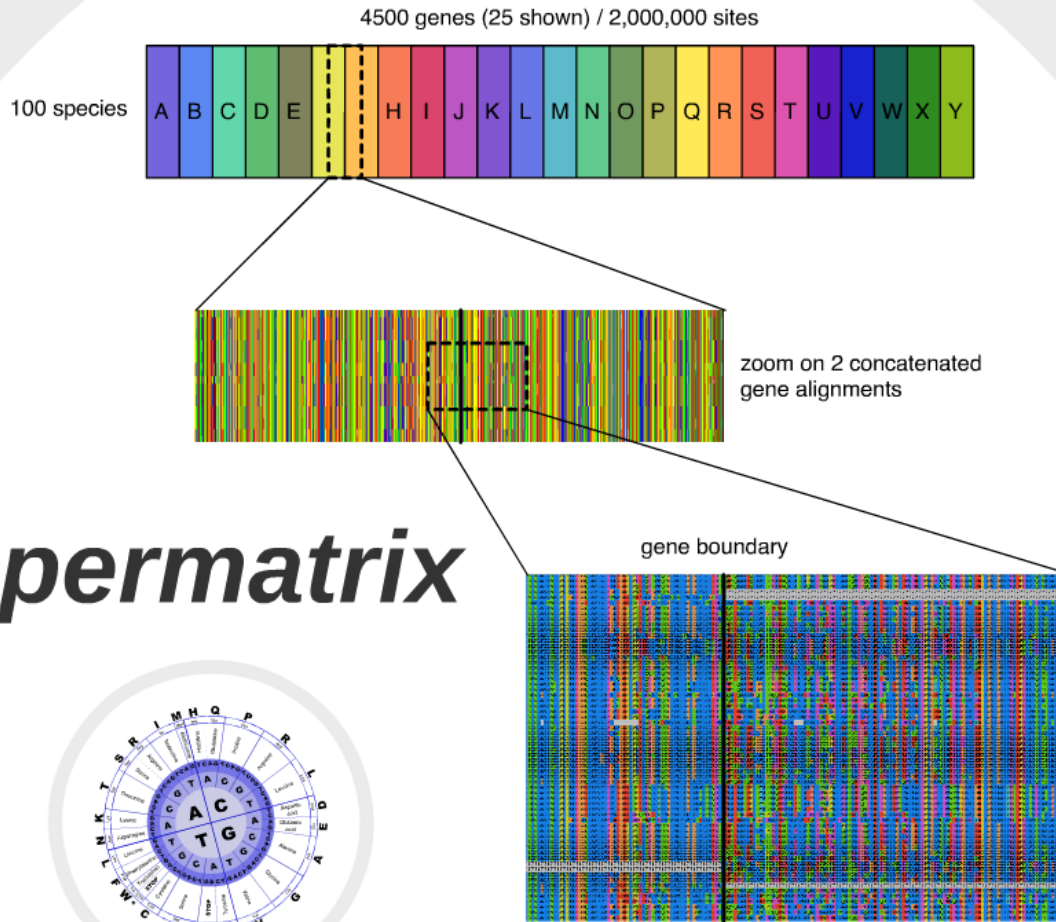
Illumina HiSeq 2000
(or MiSeq)

Dataset Assembly

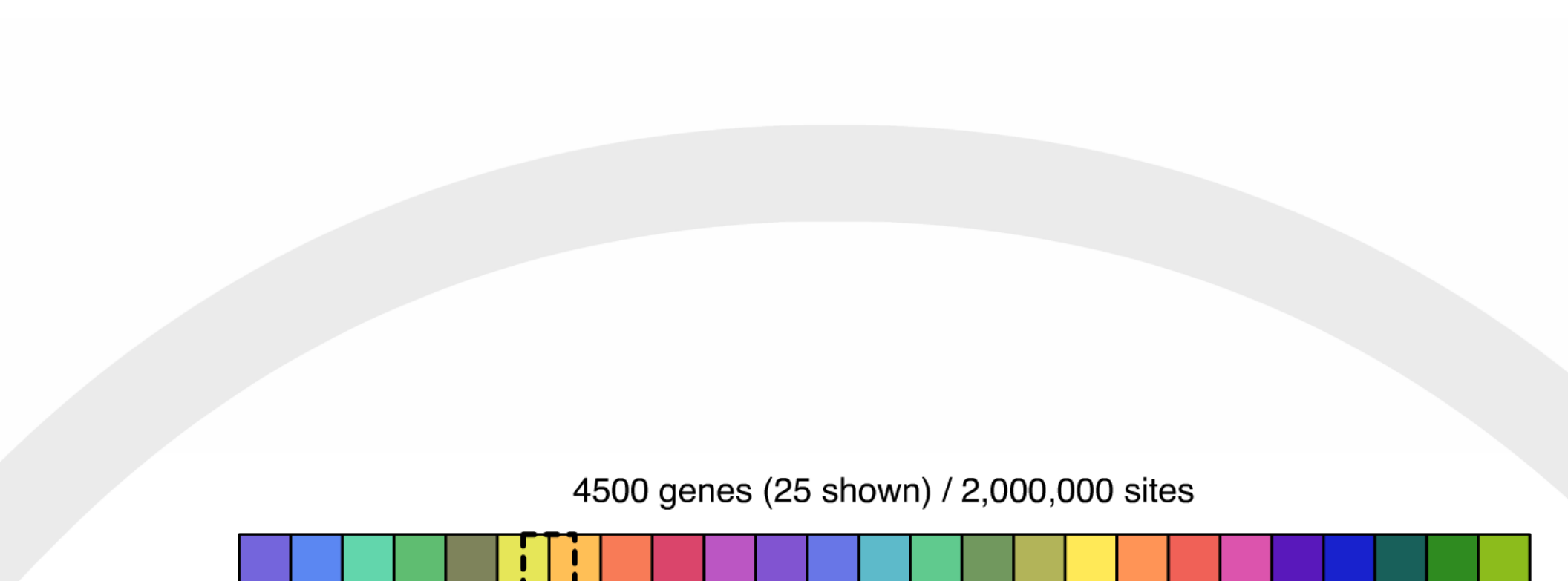
incl. sequence data
from 76 more species



Supermatrix

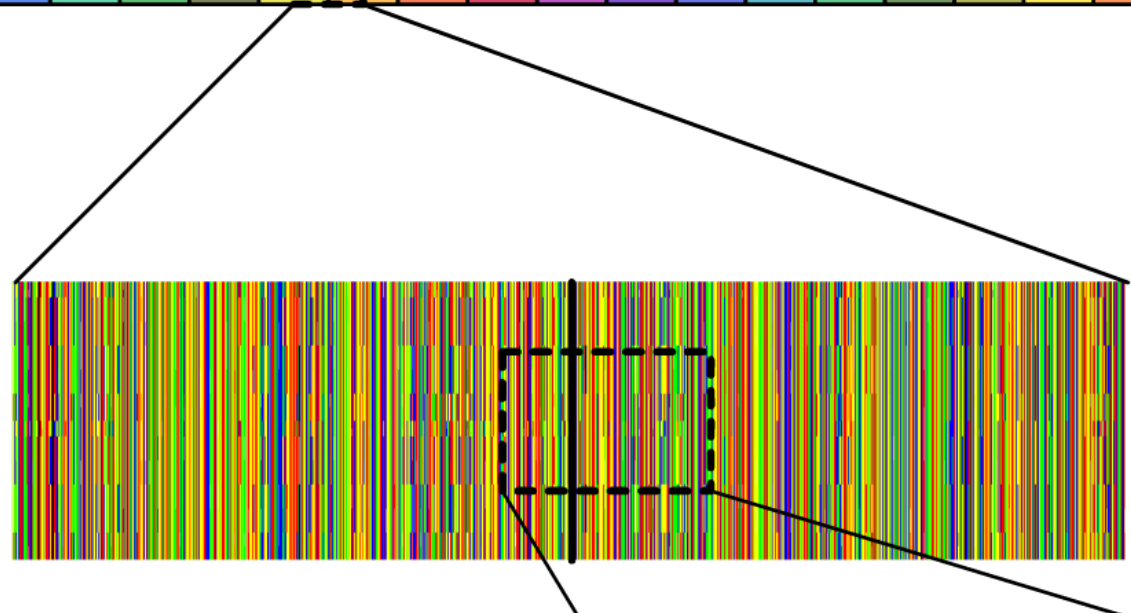
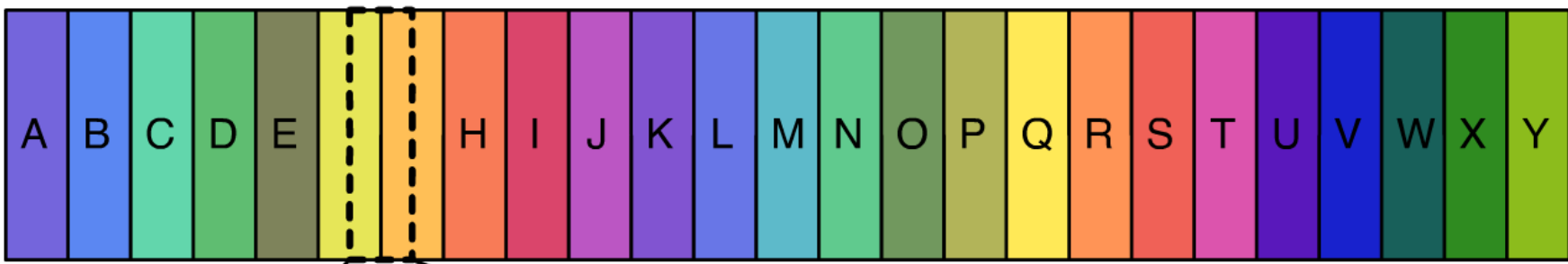


100 sp. x 2,000,000 sites
(4500 genes)

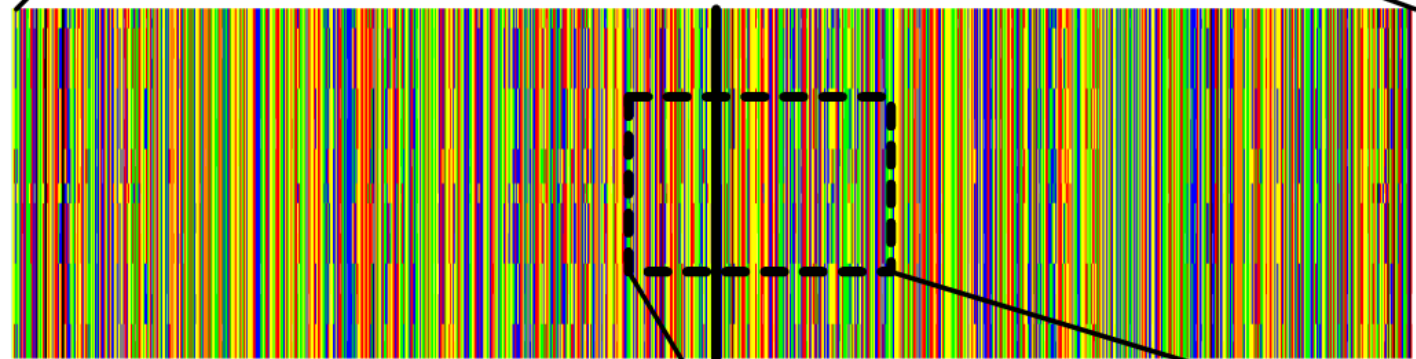


4500 genes (25 shown) / 2,000,000 sites

100 species



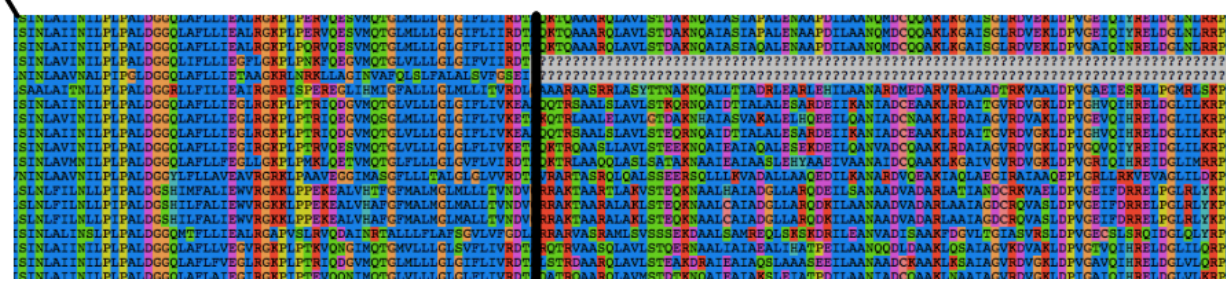
zoom on 2 concatenated gene alignments

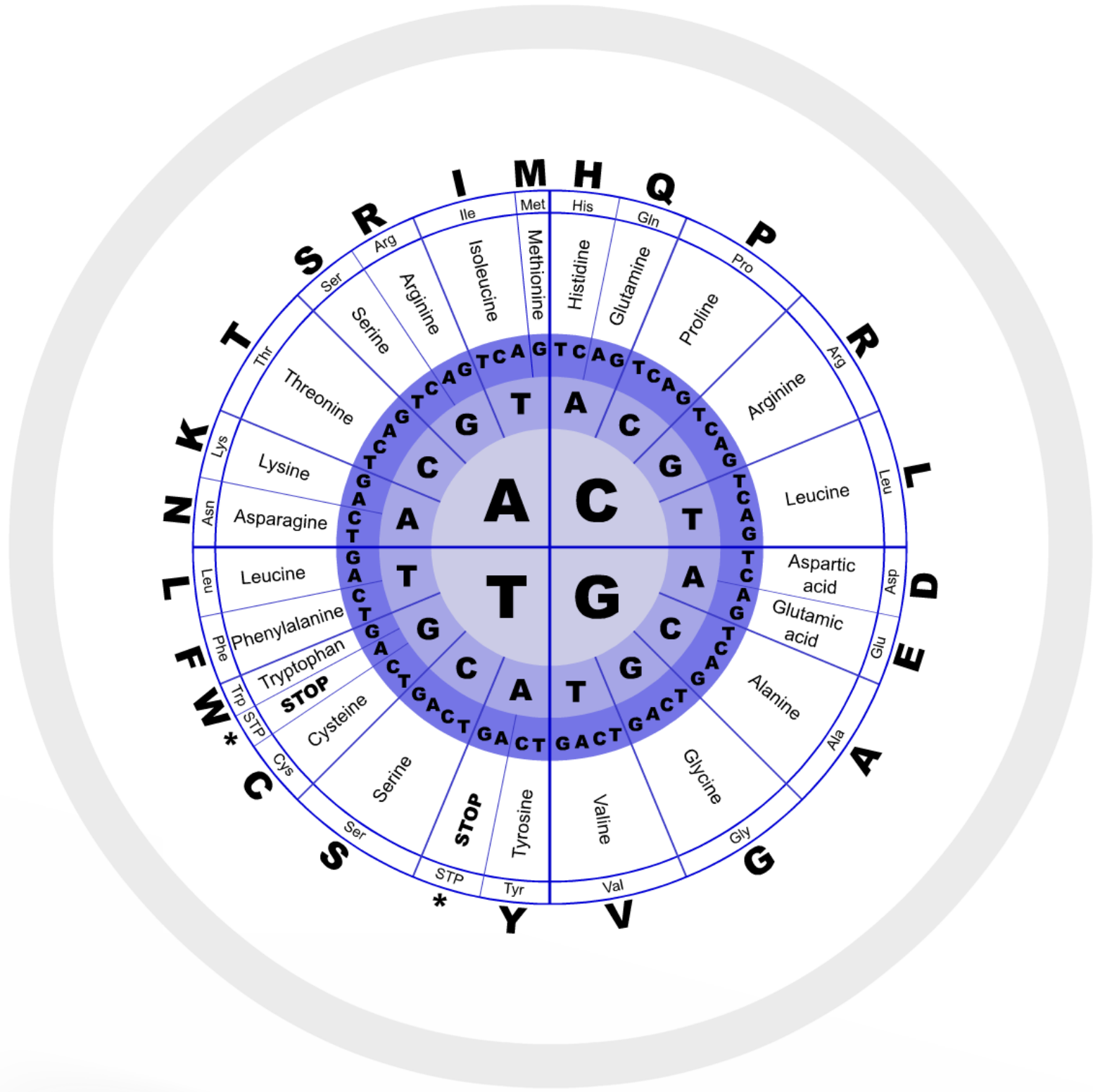


zoom on 2 concatenated gene alignments

matrix

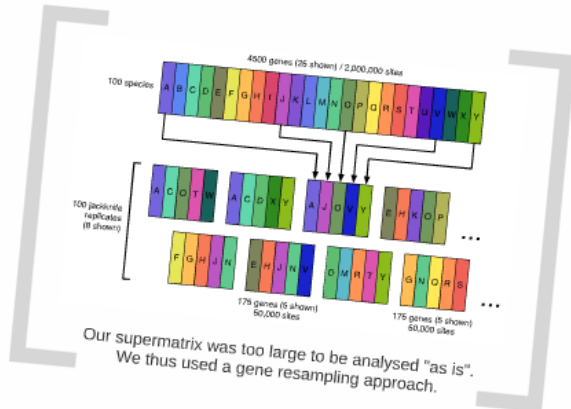
gene boundary







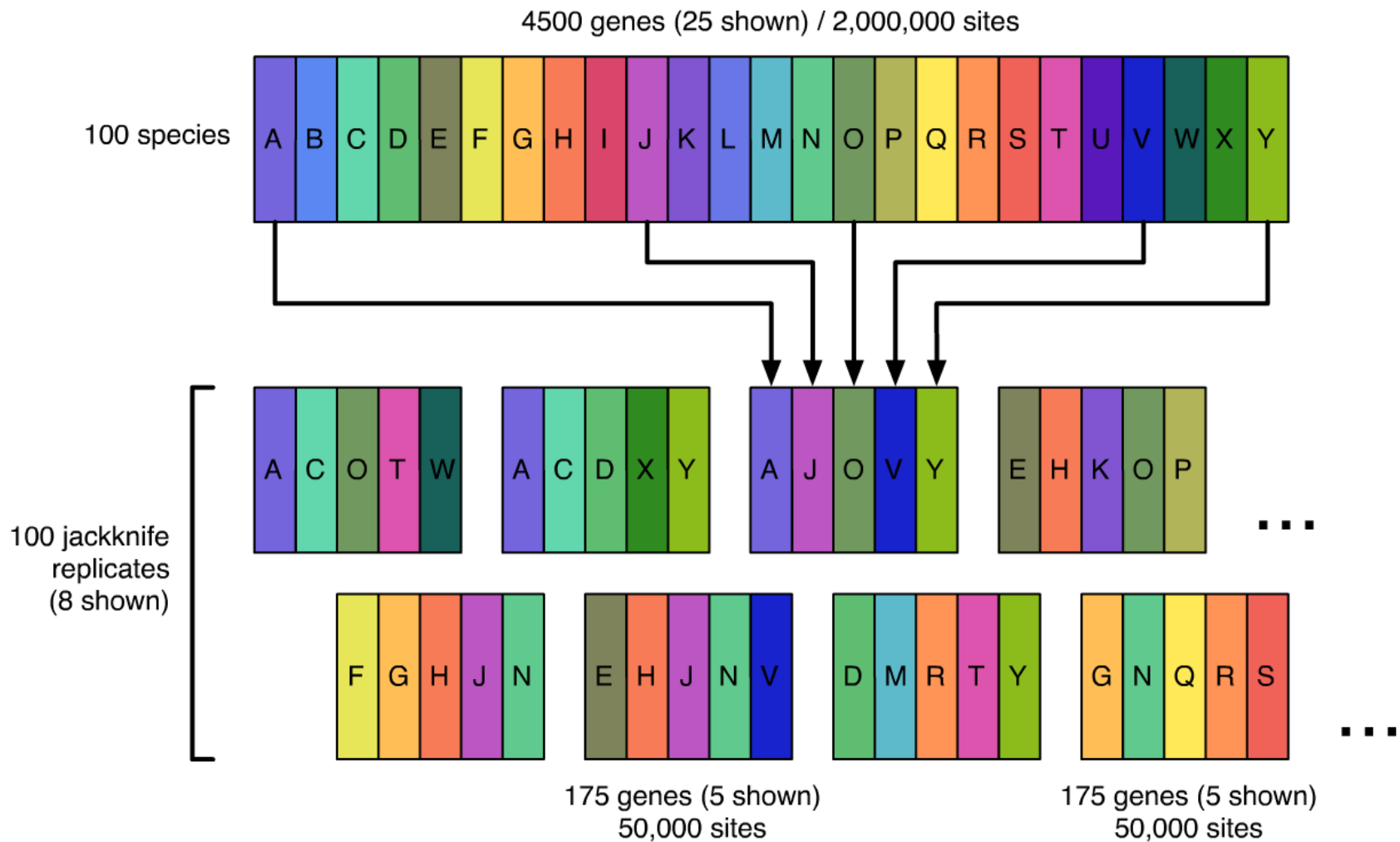
Methods



Our supermatrix was too large to be analysed "as is".
We thus used a gene resampling approach.

Each jackknife replicate was then analysed using a powerful Bayesian phylogenetic software: PhyloBayes MPI.

***Where do computers fit
in that strategy?
How do you (ab)use them?***



Our supermatrix was too large to be analysed "as is".
We thus used a gene resampling approach.



$$f(\theta|X) = \frac{f(\theta)f(X|\theta)}{f(X)} \text{ with } \theta = (\tau, v)$$

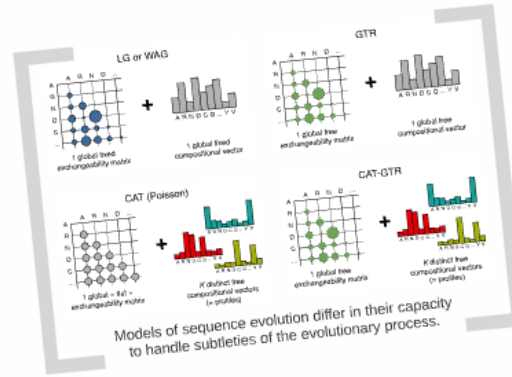
$$f(X) = \int f(\theta)f(X|\theta)d\theta$$

$$= \sum_{\tau} \int v f(v)f(X|\tau, v)dv$$

Bayes' theorem applied to phylogenetics

The tree is the interesting part of the model while the model of sequence evolution is a necessary "nuisance".

We want to compute the (posterior) probability distribution of the model (tree and sequence model) given the data (supermatrix).



N. Lartillot

Posterior Levels

MCMC random walk

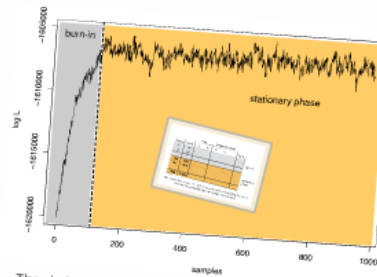


sample	tree					sequence model		log L
	cycle	τ	v	R	π			
1	10							
2	10							
...	...							
100	1000							
101	1010							
...	...							
1000	10,000							

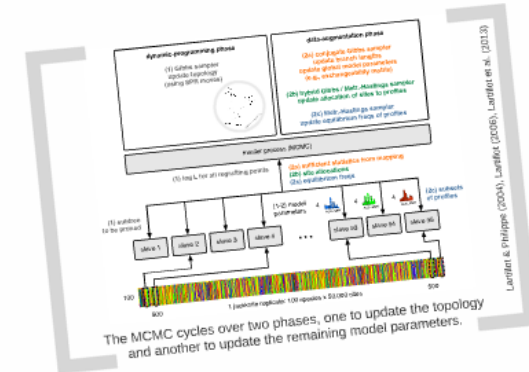
MCMC samples of the parameter space

The best sequence models use complex mixtures that can only be implemented in a Bayesian framework.

The principle is to sample the posterior distribution using numerical simulation (Markov Chain Monte Carlo).



The chain is left running for days until it reaches convergence. Being memoryless, it is naturally restartable at will.



Each jackknife replicate was then analysed using a powerful Bayesian phylogenetic software: PhyloBayes MPI.

$$f(\theta|X) = \frac{f(\theta)f(X|\theta)}{f(X)} \quad \text{with } \theta = (\tau, v)$$

$$f(X) = \int f(\theta)f(X|\theta)d\theta$$

$$= \sum_{\tau} \int_v f(v)f(X|\tau, v)dv$$

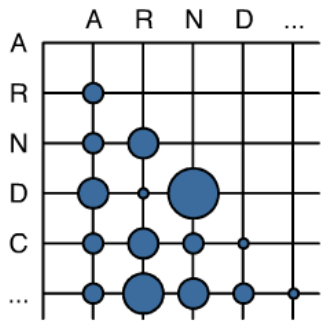
Bayes' theorem applied to phylogenetics



The tree *is* the interesting part of the model while the model of sequence evolution is a necessary "nuisance".

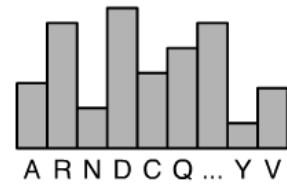
We want to compute the (posterior) probability distribution of the model (tree and sequence model) given the data (supermatrix).

LG or WAG



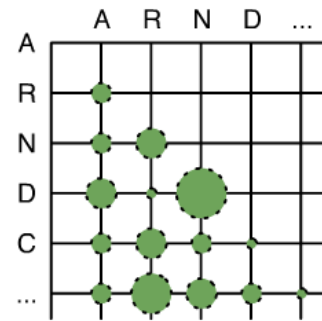
1 global fixed exchangeability matrix

+



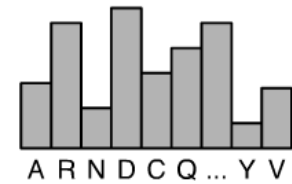
1 global fixed compositional vector

GTR



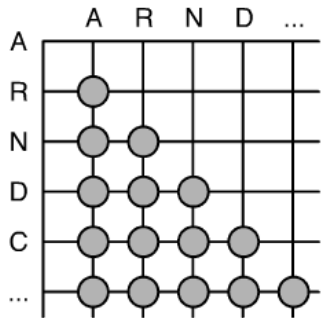
1 global free exchangeability matrix

+



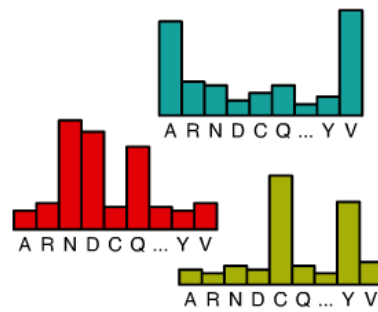
1 global free compositional vector

CAT (Poisson)



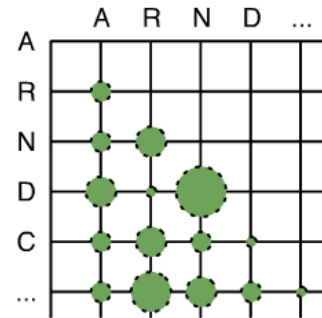
1 global « flat » exchangeability matrix

+



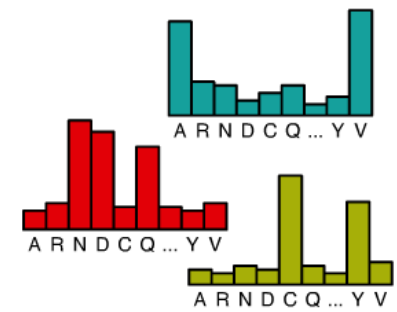
K distinct free compositional vectors (= profiles)

CAT-GTR



1 global free exchangeability matrix

+



K distinct free compositional vectors (= profiles)

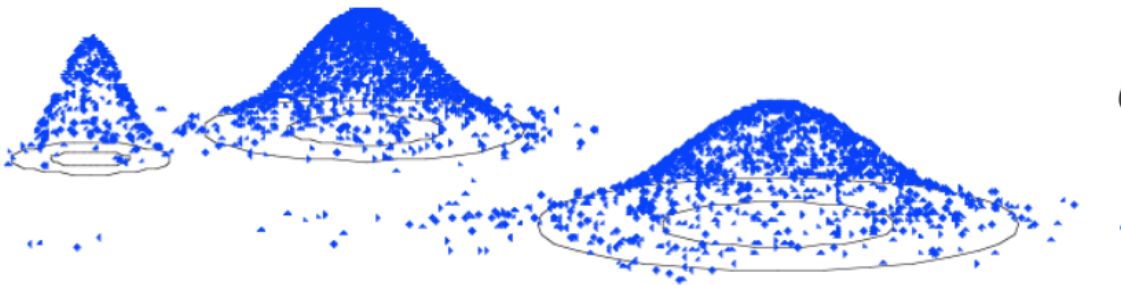
Models of sequence evolution differ in their capacity to handle subtleties of the evolutionary process.



MCMC random walk

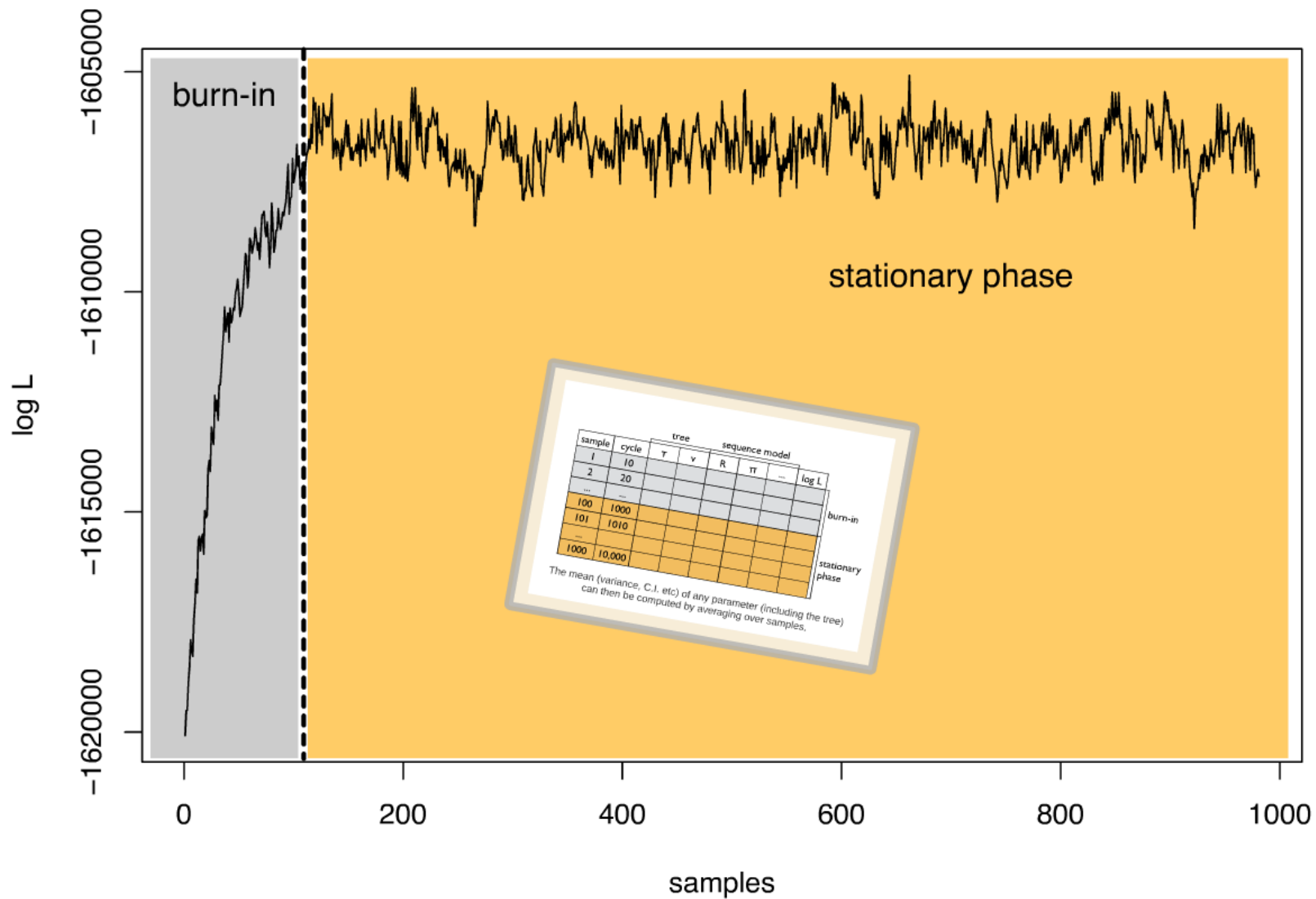
sample	cycle	tree		sequence model			
		τ	ν	R	π	...	$\log L$
1	10						
2	20						
...	...						
100	1000						
101	1010						
...							
1000	10,000						

***MCMC samples
of the parameter space***



The best sequence models use complex mixtures that can only be implemented in a Bayesian framework.

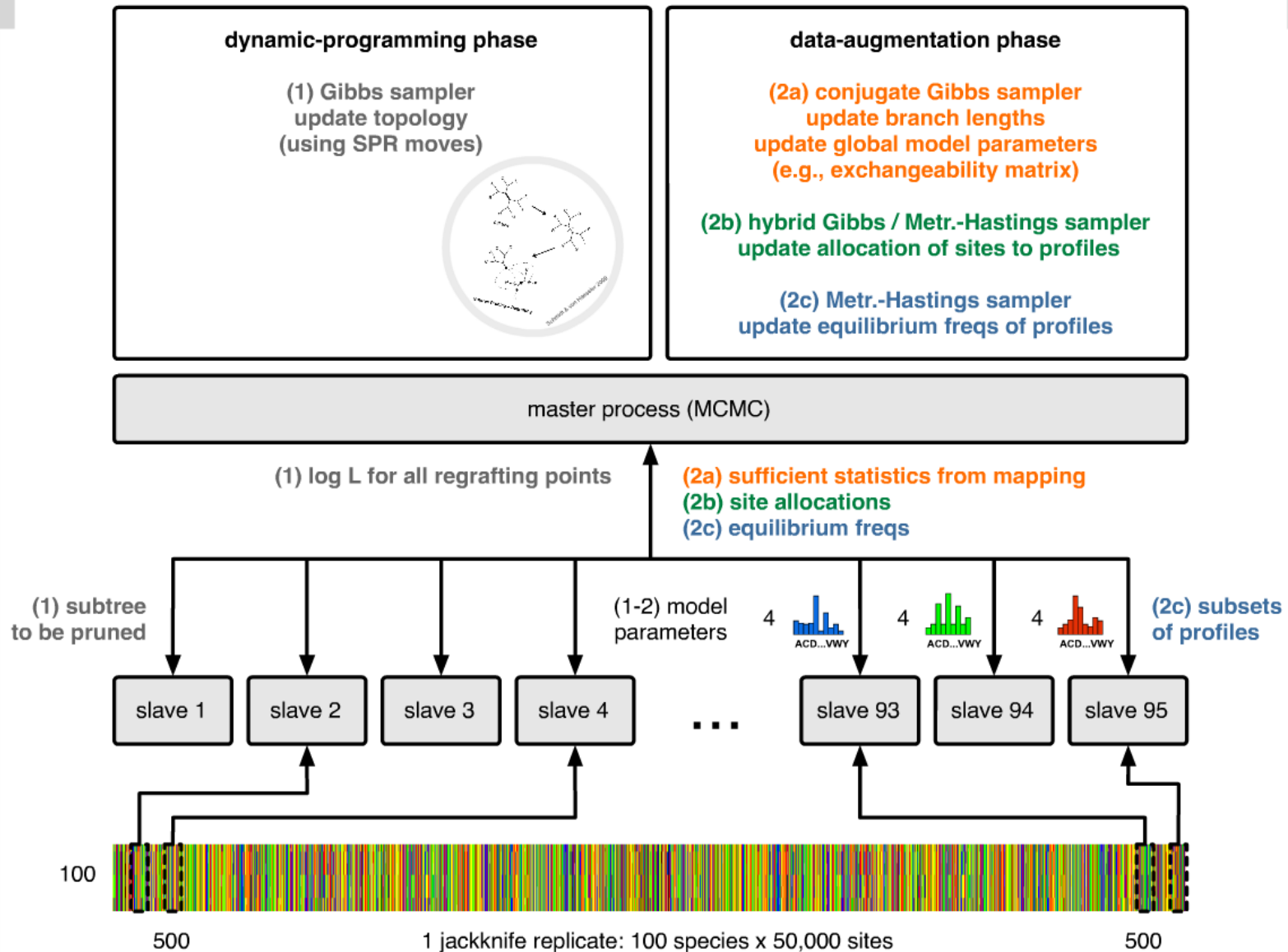
The principle is to sample the posterior distribution using numerical simulation (Markov Chain Monte Carlo).



The chain is left running for days until it reaches convergence.
 Being memoryless, it is naturally restartable at will.

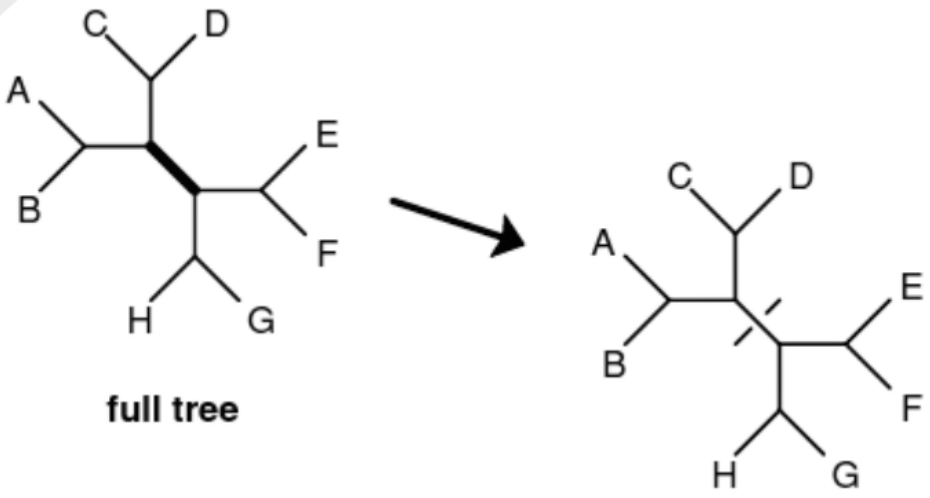
sample	cycle	tree		sequence model			log L	
		τ	v	R	π	...		
1	10							burn-in
2	20							
...	...							
100	1000							stationary phase
101	1010							
...								
1000	10,000							

The mean (variance, C.I. etc) of any parameter (including the tree) can then be computed by averaging over samples.

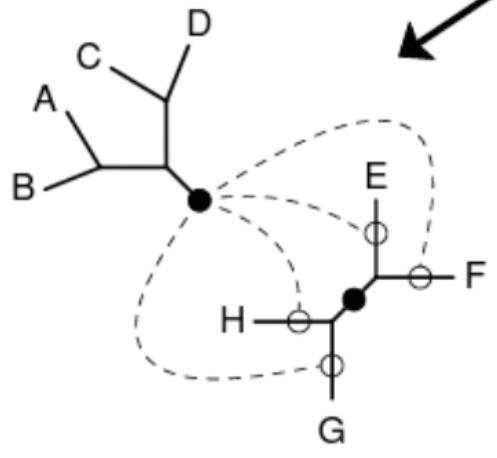


The MCMC cycles over two phases, one to update the topology and another to update the remaining model parameters.

oves)



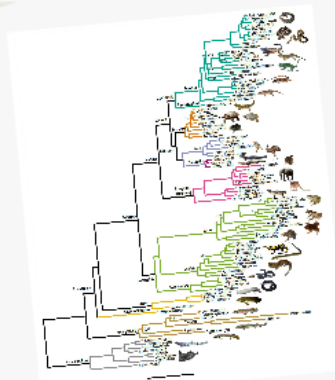
full tree



Subtree Pruning + Regrafting

Schmidt & von Haeseler 2009

Results

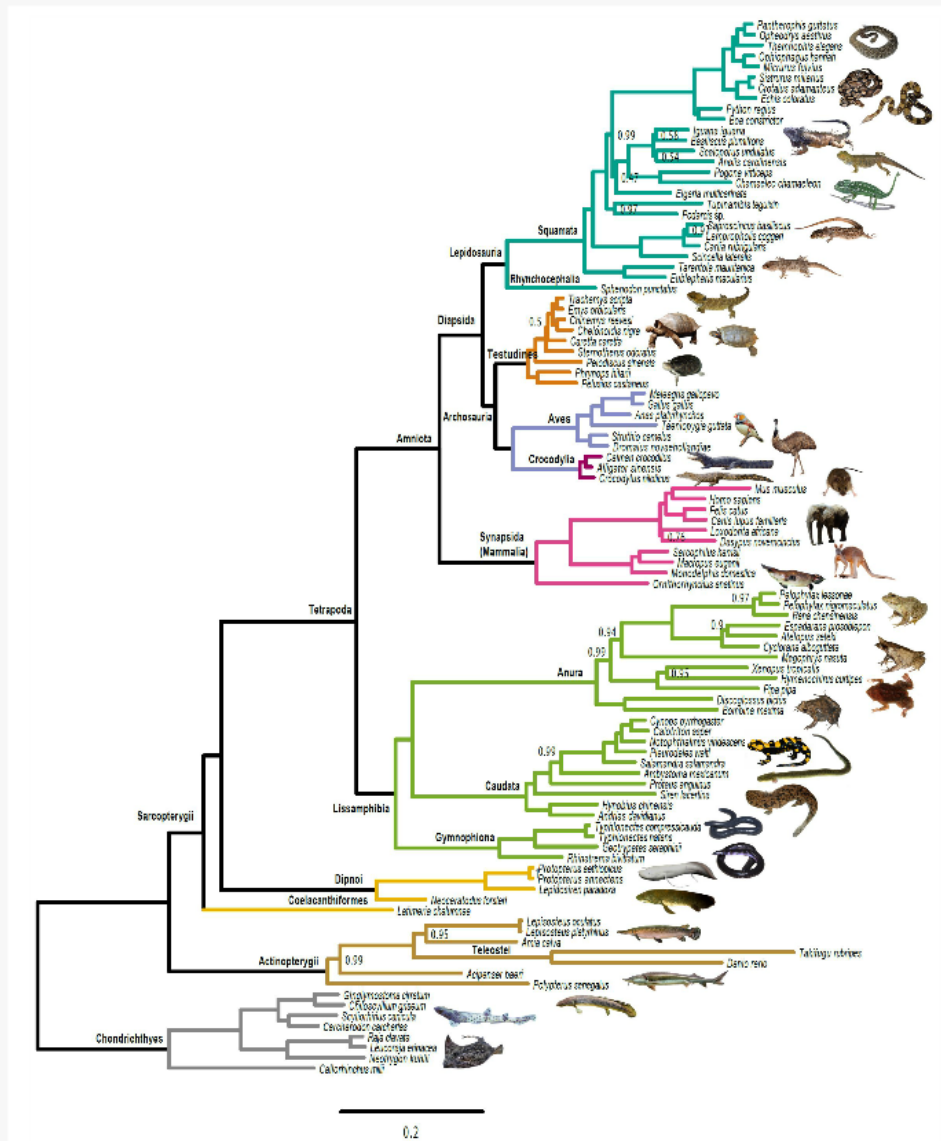


Irisani, Baurain et al. (in prep.)

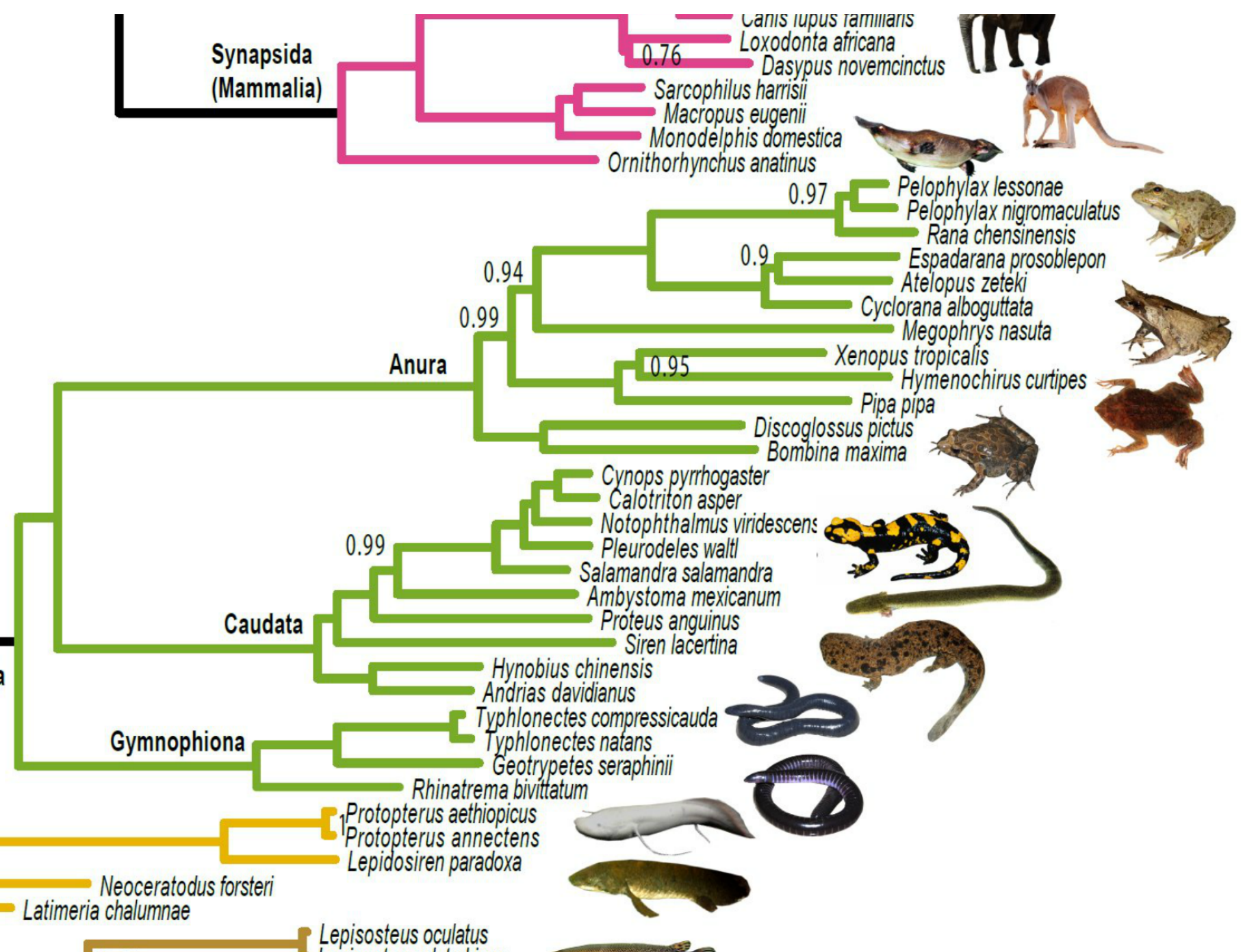
We obtained a new reference phylogenetic framework for the evolution vertebrates.

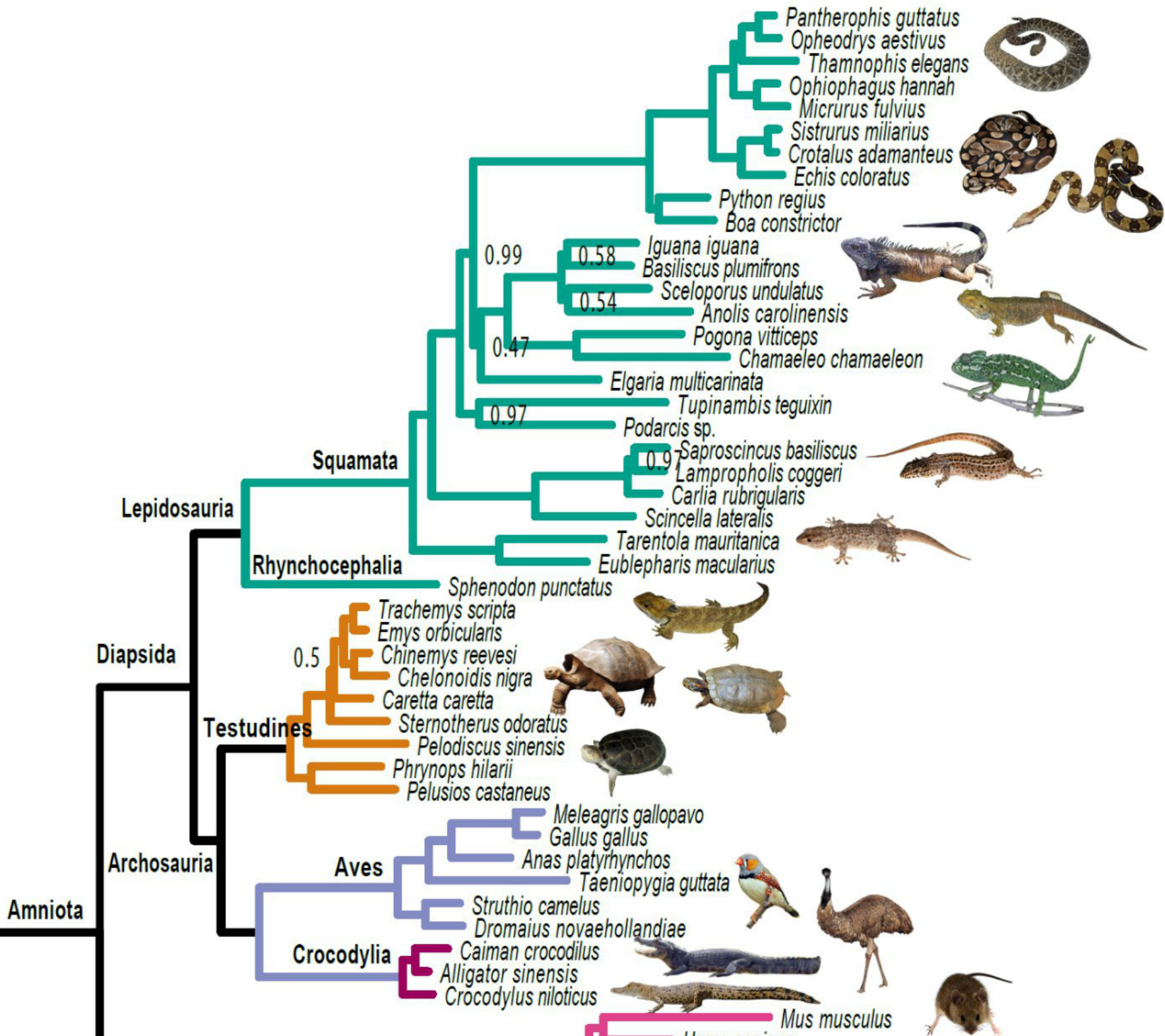
***What results
have computers enabled
you to achieve?***





We obtained a new reference phylogenetic framework for the evolution vertebrates.





Perspectives

```

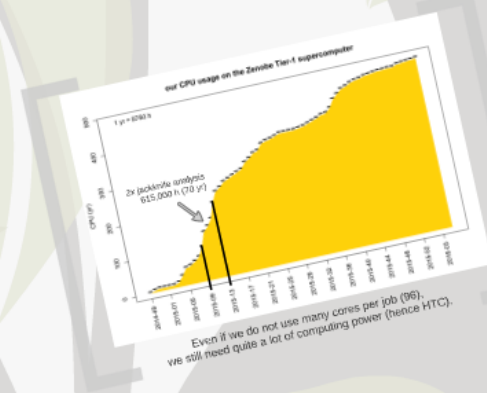
# script job template
cat > ./templates/jb-01
#!/bin/bash

source /usr/share/modules/init/bash
module load compiler/intel/composer/2011_gpl_1.2.0x_1p46
module load intelmpi/4.2.1-020764

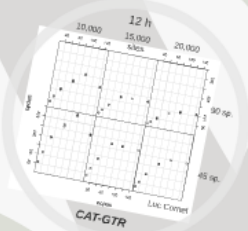
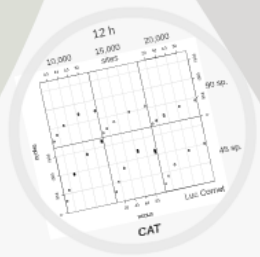
mpirun -n 18 (rank0 1) \
  /projects/cond/p44k44/software/jb_mpi_1.3e/data/jb_mpi_1
  { 18 hostname |} psu -de -out \
  {psu 1 -> 10 1000} |} gnu1 8256 |} polarsort 1 stop |}
  {psu 1 -> 10 1000} { hostname |} { psu1 |}

# prepare job name
for f in $(ls 2 > /dev/null); do
  do qsub --define thread=${f} |}
  --define hostname=jack2000-100-01 --define node=c200 \
  --define rank ./templates/jb_01 > jack2000-100-01.sh
done
  
```

We use templating and job arrays whenever possible to ease up submission and restart.



What experience do you feel can be transposed to another research field? Can you share any tips and tricks?



12 h

10,000

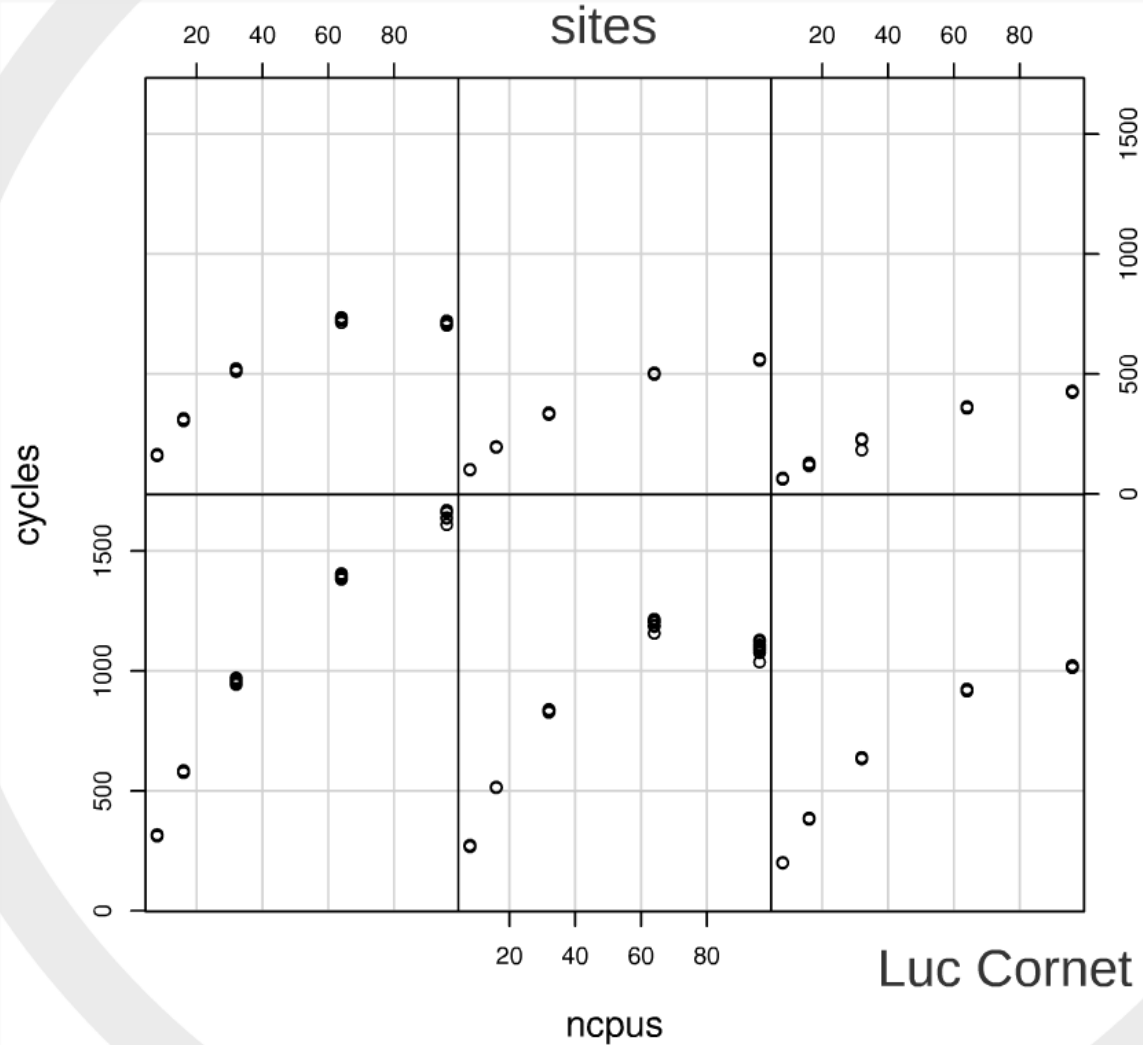
15,000

20,000

20 40 60 80

sites

20 40 60 80



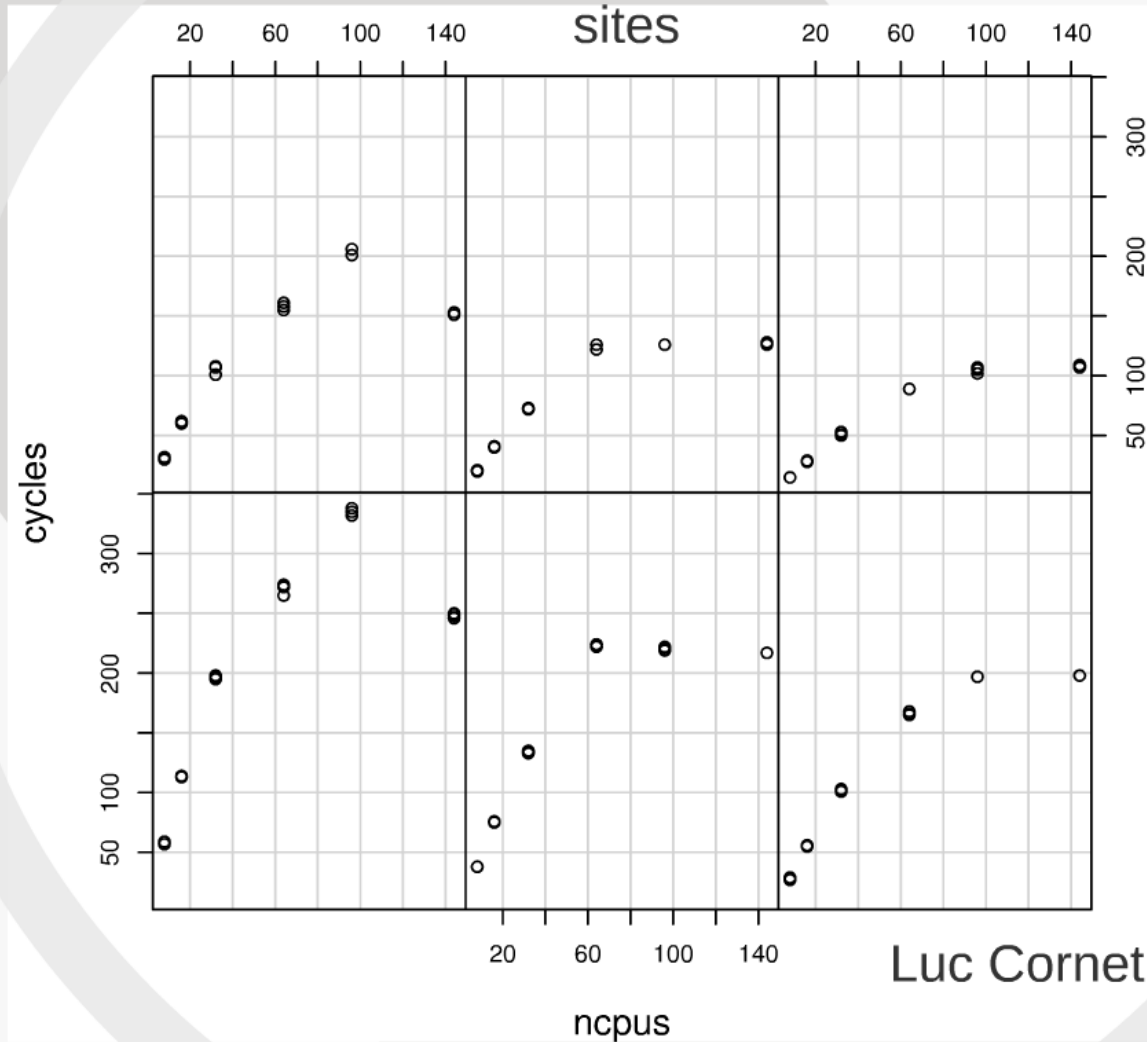
CAT

12 h

10,000

15,000

20,000



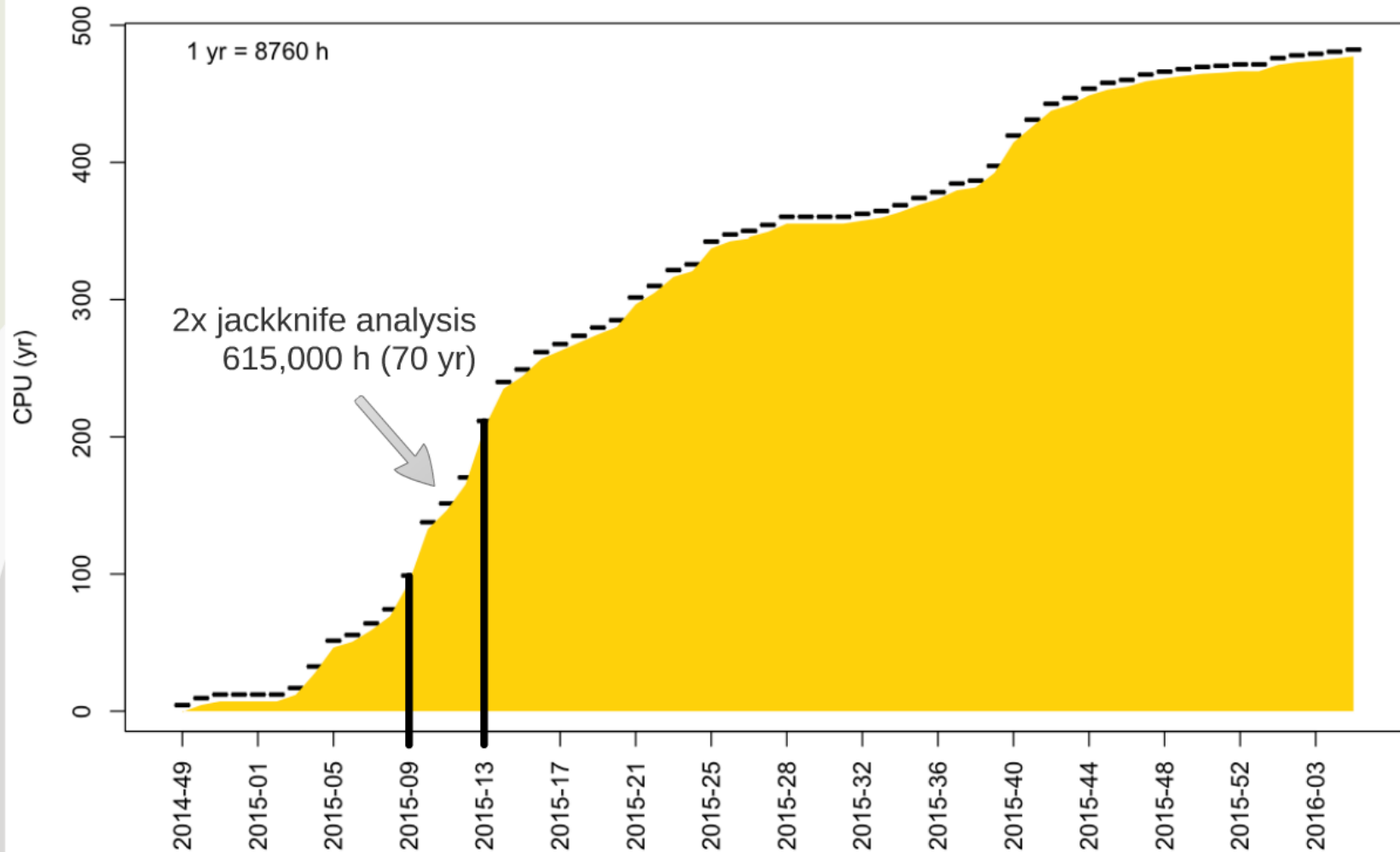
90 sp.

45 sp.

Luc Cornet

CAT-GTR

our CPU usage on the Zenobe Tier-1 supercomputer



Even if we do not use many cores per job (96), we still need quite a lot of computing power (hence HTC).

Acknowledgment

- **Liège**

- Luc Cornet
- Raphaël Léonard
- Damien Sirjacobs

- **CÉCI**

- David Colignon

- **Cenaero**

- Danielle Coulon

- **Région wallonne**

- Gr. agr. 1117545

- **Konstanz**

- Iker Irisarri

- **Braunschweig**

- Henner Brinkmann
- Jörn Petersen
- Miguel Vences

- **Montpellier**

- Frédéric Delsuc

- **Moulis / Montreal**

- Hervé Philippe



Let's have a drink now!