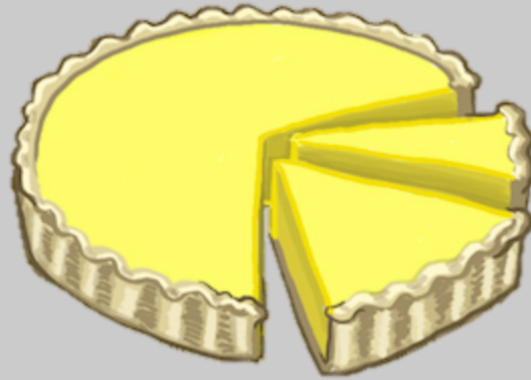


PartitionFinder



Paul B. Frandsen
Office of Research Information
Services



With Rob Lanfear and Brett Calcott

Google

SUMMER

OF

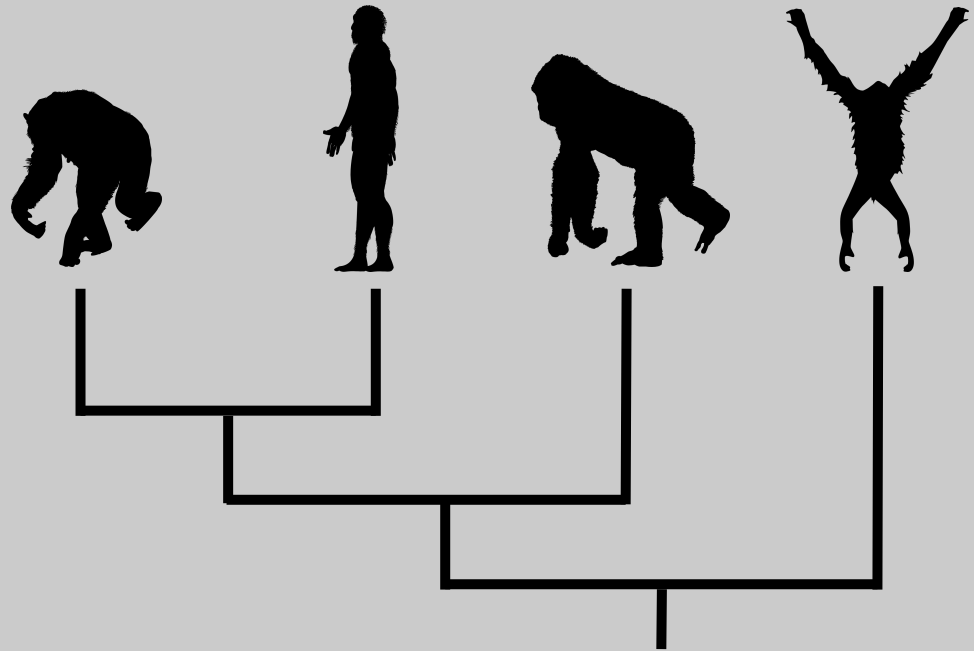
CODE

2013

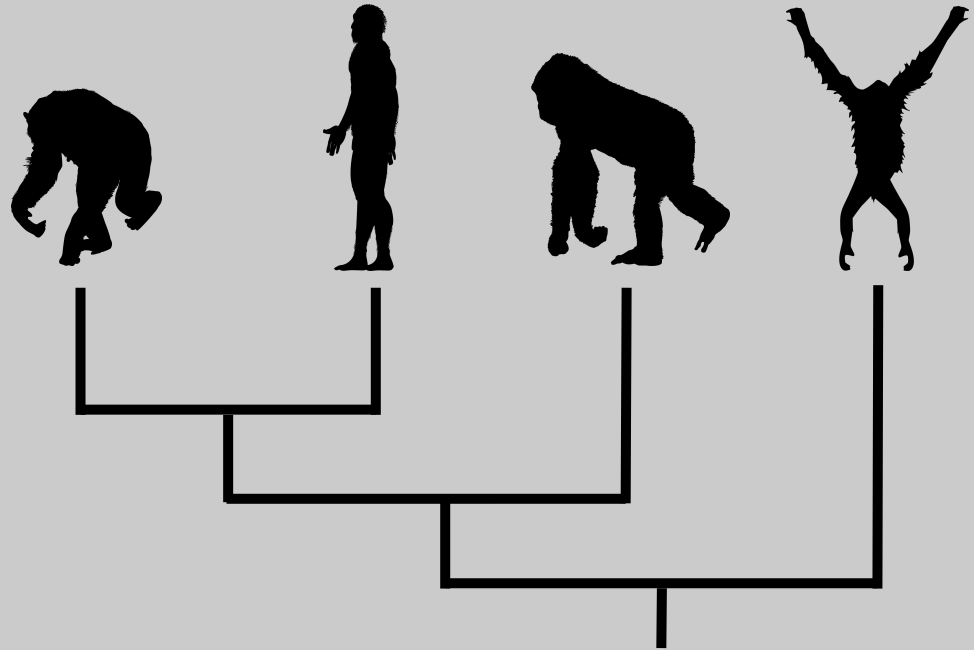
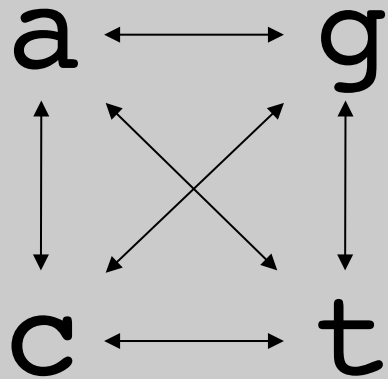


a c t g c c t g a c t g g c t g a c
g c t g a c g t c a t g g c t g a c
a g t g a c g t a g t g a g t g c g
a g c g c c t g c g c g a g c g c g

a c t g c c t g a c t g g c t g a c
g c t g a c g t c a t g g c t g a c
a g t g a c g t a g t g a g t g c g
a g c g c c t g c g c g a g c g c g

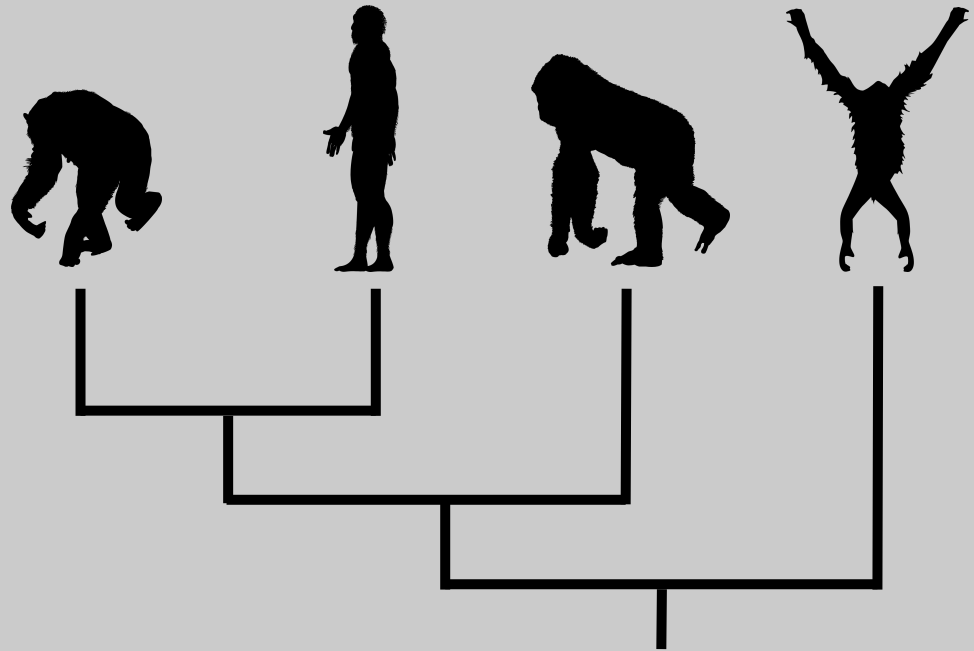
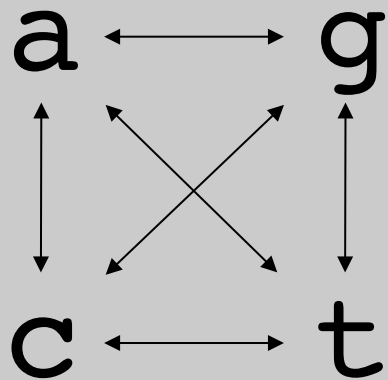


a c t g c c t g a c t g g c t g a c
g c t g a c g t c a t g g c t g a c
a g t g a c g t a g t g a g t g c g
a g c g c c t g c g c g a g c g c g

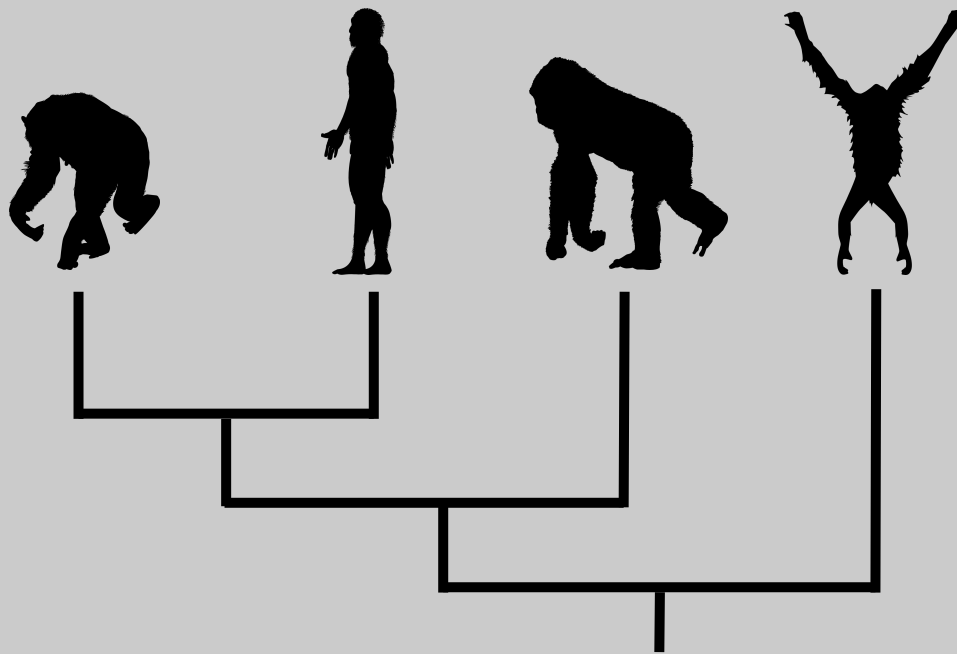
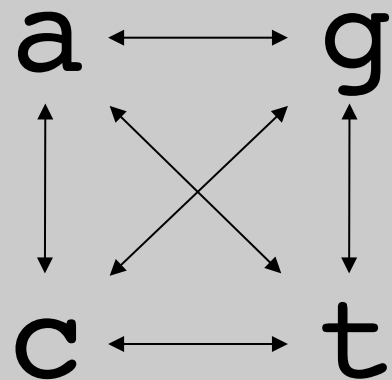


a c t g c c t g a c t g g c t g a c
g c t g a c g t c a t g g c t g a c
a g t g a c g t a g t g a g t g c g
a g c g c c t g c g c g a g c g c g

1 model

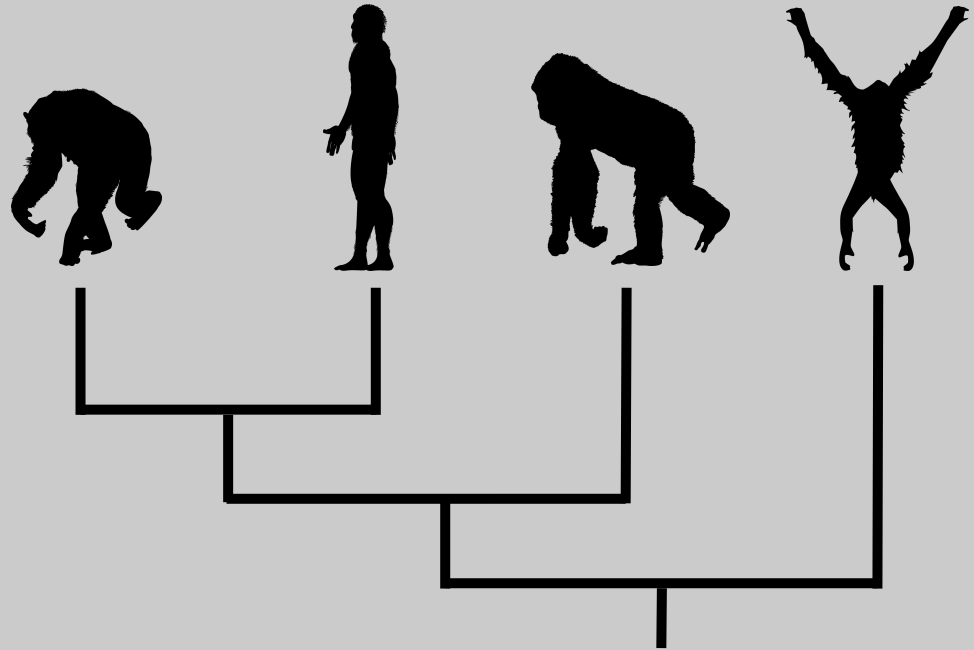
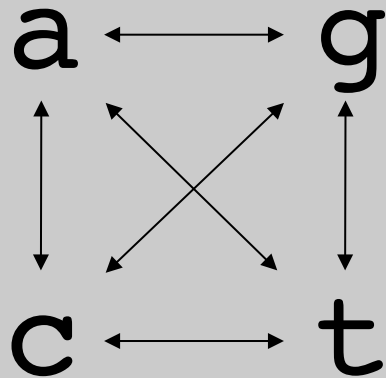


a c t g c c t g a c t g g c t g a c
g c t g a c g t c a t g g c t g a c
a g t g a c g t a g t g a g t g c g
a g c g c c t g c g c g a g c g c g



a	c	t	g	c	c	t	g	a	c	t	g	g	c	t	g	a	c
g	c	t	g	a	c	g	t	c	a	t	g	g	c	t	g	a	c
a	g	t	g	a	c	g	t	a	g	t	g	a	g	t	g	c	g
a	g	c	g	c	c	t	g	c	g	c	g	a	g	c	g	c	g

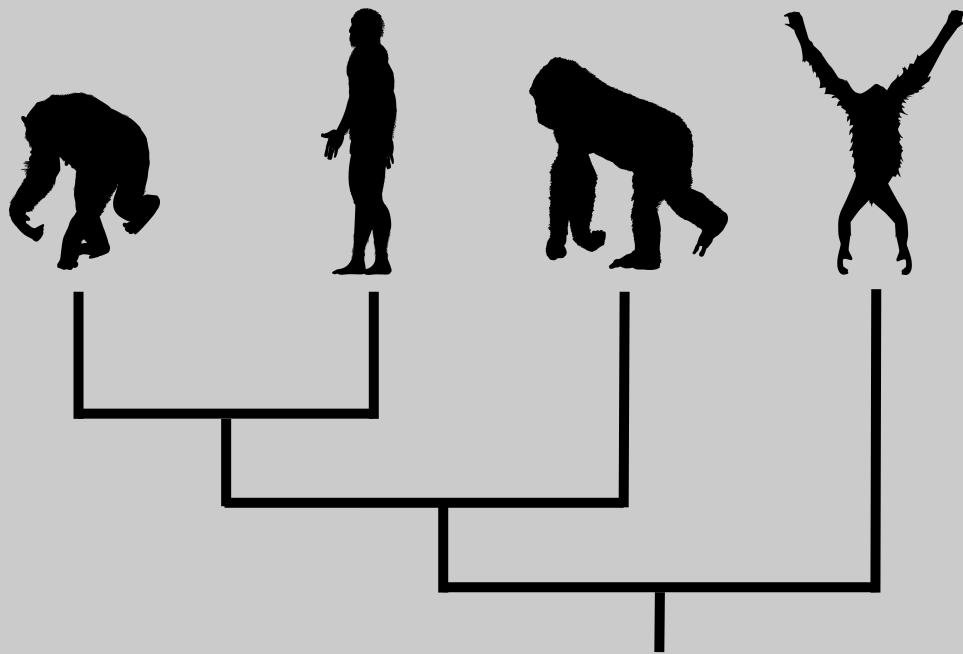
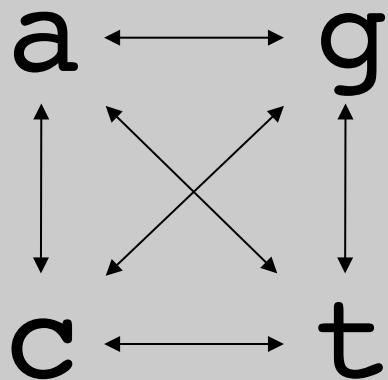
3 models

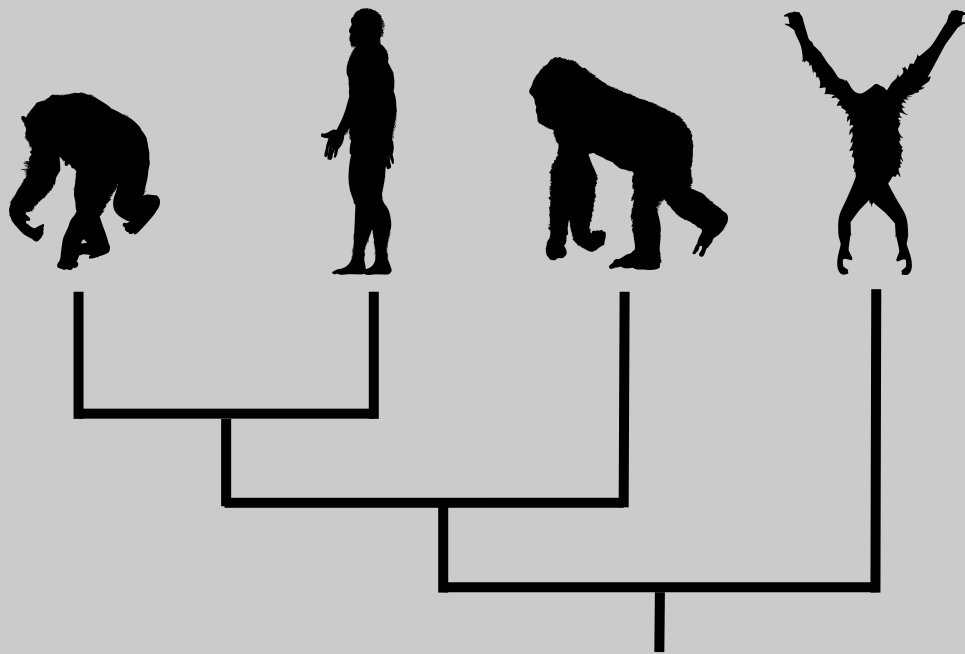
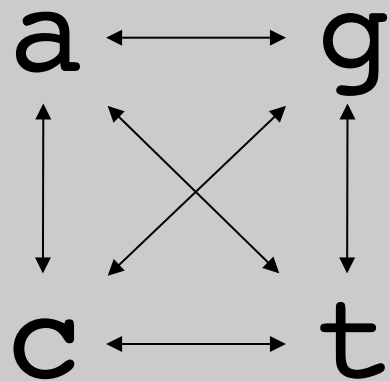
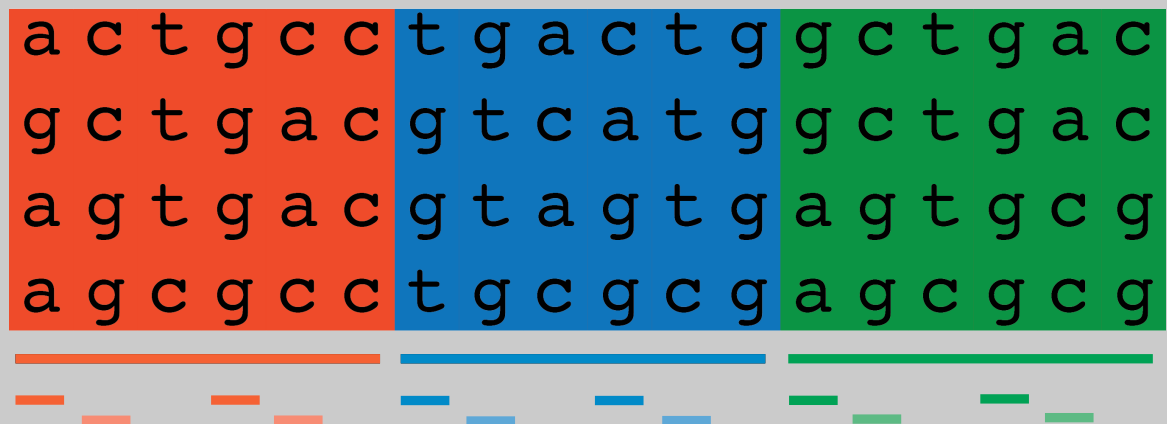


a	c	t	g	c	c	t	g	a	c	t	g	g	c	t	g	a	c
g	c	t	g	a	c	g	t	c	a	t	g	g	c	t	g	a	c
a	g	t	g	a	c	g	t	a	g	t	g	a	g	t	g	c	g
a	g	c	g	c	c	t	g	c	g	c	g	a	g	c	g	c	g

—
—
—

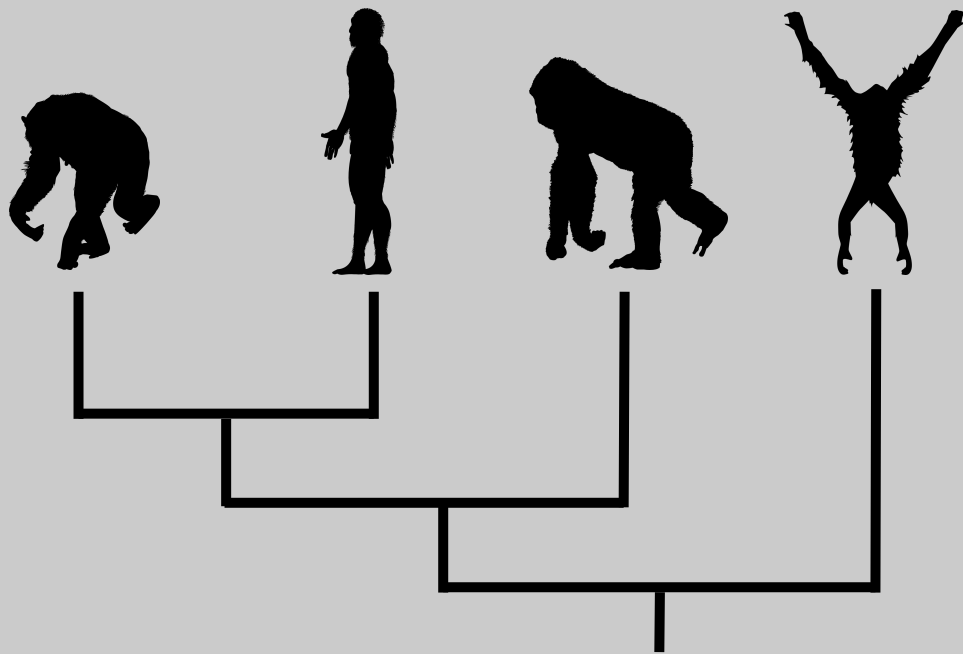
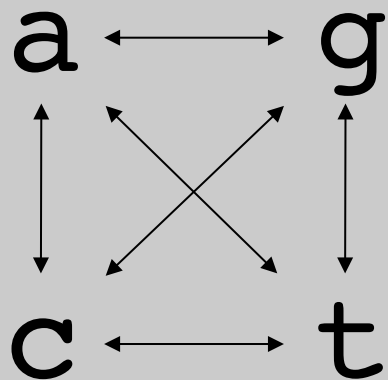
—
—
—

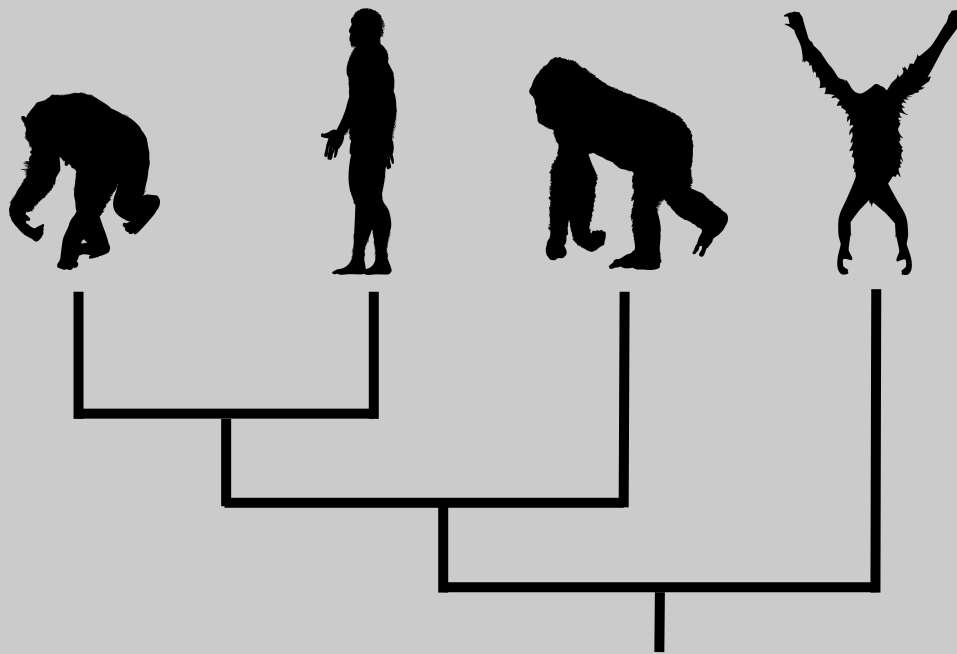
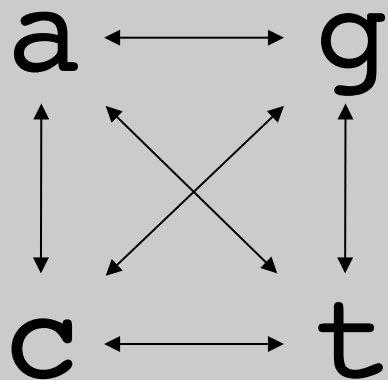
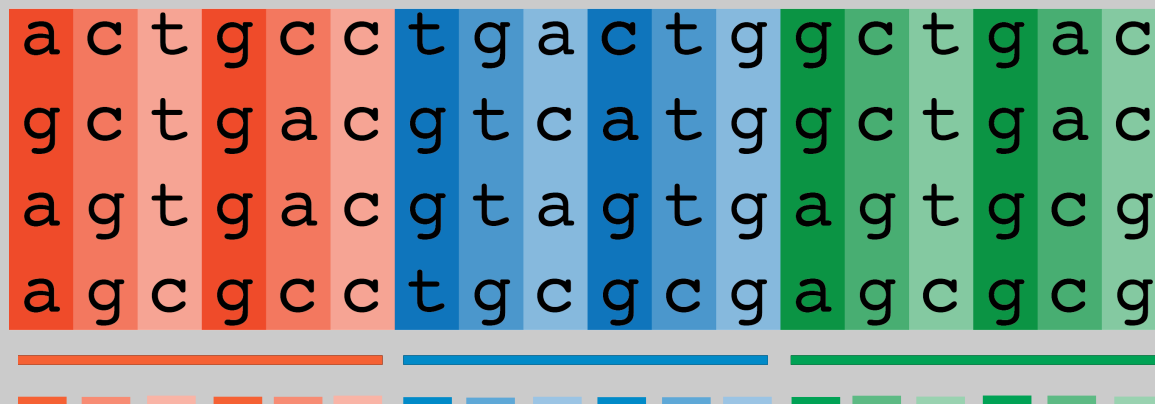




a	c	t	g	c	c	t	g	a	c	t	g	g	c	t	g	a	c
g	c	t	g	a	c	g	t	c	a	t	g	g	c	t	g	a	c
a	g	t	g	a	c	g	t	a	g	t	g	a	g	t	g	c	g
a	g	c	g	c	c	t	g	c	g	c	g	a	g	c	g	c	g

Below the sequence blocks are three horizontal bars: orange, blue, and green. Underneath these bars are several short, dashed horizontal lines in corresponding colors, representing a genomic map or track.





a c t g c c t g a c t g g c t g a c
g c t g a c g t c a t g g c t g a c
a g t g a c g t a g t g a g t g c g
a g c g c c t g c g c g a g c g c g

1

a c t g c c t g a c t g g c t g a c
g c t g a c g t c a t g g c t g a c
a g t g a c g t a g t g a g t g c g
a g c g c c t g c g c g a g c g c g

3

a c t g c c t g a c t g g c t g a c
g c t g a c g t c a t g g c t g a c
a g t g a c g t a g t g a g t g c g
a g c g c c t g c g c g a g c g c g

9

a	c	t	g	c	c	t	g	a	c	t	g	g	c	t	g	a	c
g	c	t	g	a	c	g	t	c	a	t	g	g	c	t	g	a	c
a	g	t	g	a	c	g	t	a	g	t	g	a	g	t	g	c	g
a	g	c	g	c	c	t	g	c	g	c	g	a	g	c	g	c	g

3

a	c	t	g	c	c	t	g	a	c	t	g	g	c	t	g	a	c
g	c	t	g	a	c	g	t	c	a	t	g	g	c	t	g	a	c
a	g	t	g	a	c	g	t	a	g	t	g	a	g	t	g	c	g
a	g	c	g	c	c	t	g	c	g	c	g	a	g	c	g	c	g

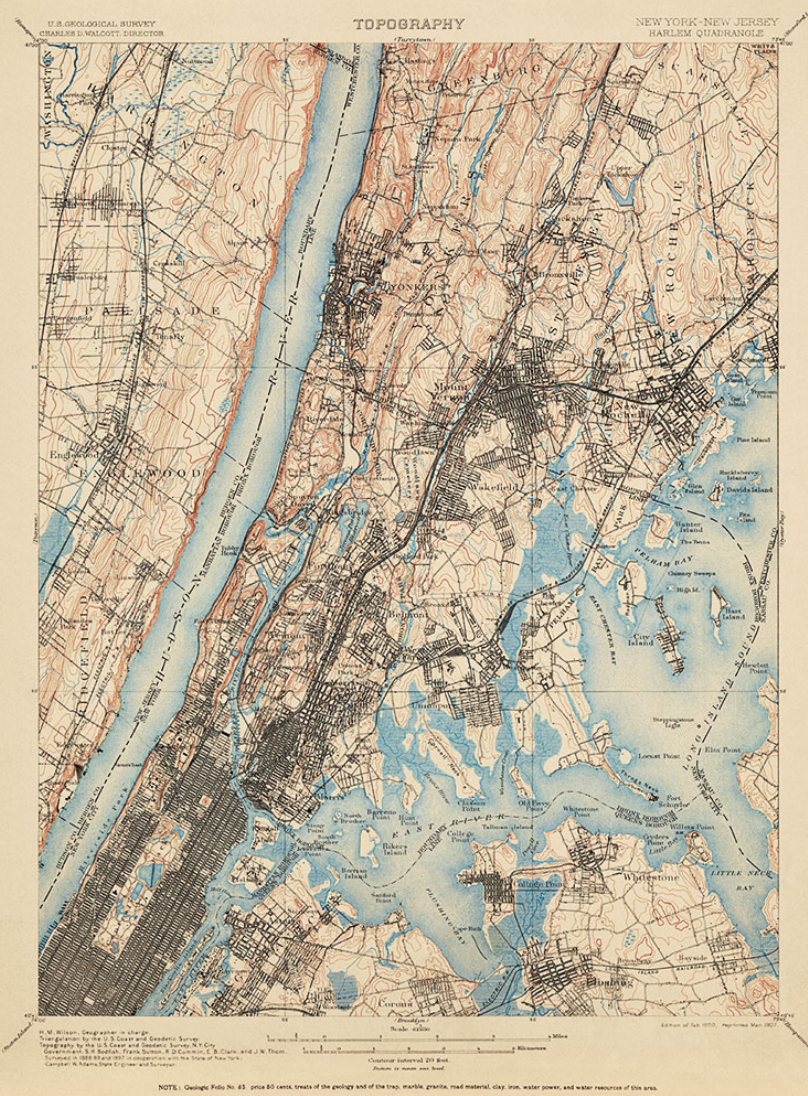
9

a	c	t	g	c	c	t	g	a	c	t	g	g	c	t	g	a	c
g	c	t	g	a	c	g	t	c	a	t	g	g	c	t	g	a	c
a	g	t	g	a	c	g	t	a	g	t	g	a	g	t	g	c	g
a	g	c	g	c	c	t	g	c	g	c	g	a	g	c	g	c	g

6

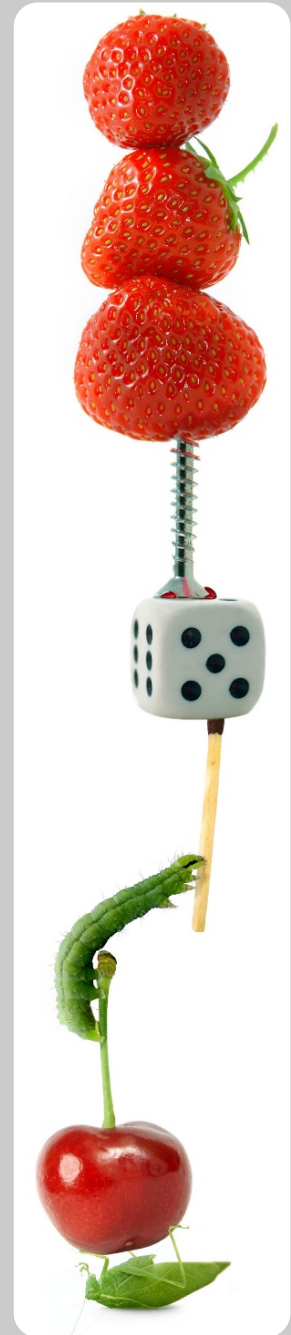
a	c	t	g	c	c	t	g	a	c	t	g	g	c	t	g	a	c
g	c	t	g	a	c	g	t	c	a	t	g	g	c	t	g	a	c
a	g	t	g	a	c	g	t	a	g	t	g	a	g	t	g	c	g
a	g	c	g	c	c	t	g	c	g	c	g	a	g	c	g	c	g

2

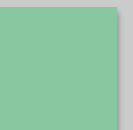
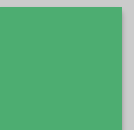
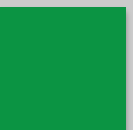
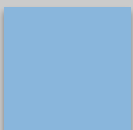
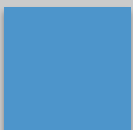
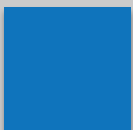
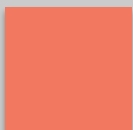
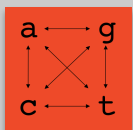


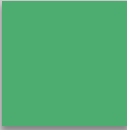
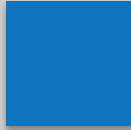
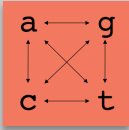
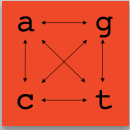
Information theoretic metrics

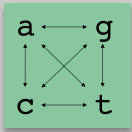
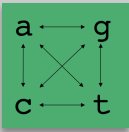
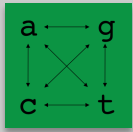
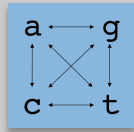
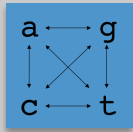
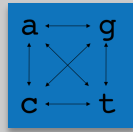
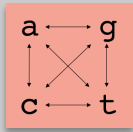
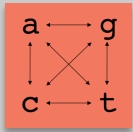
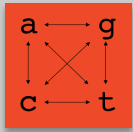
- AIC, AICc, BIC
- Balancing act between over-parameterization and under parameterization

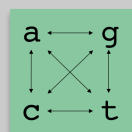
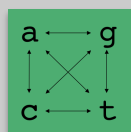
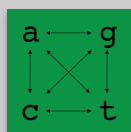
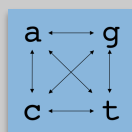
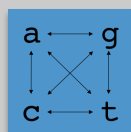
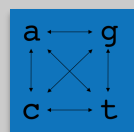
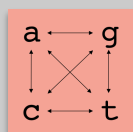
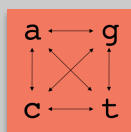
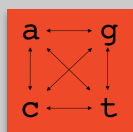




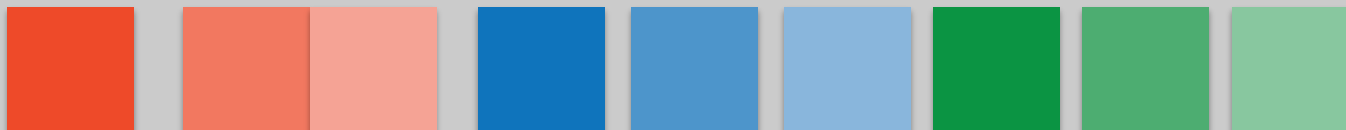
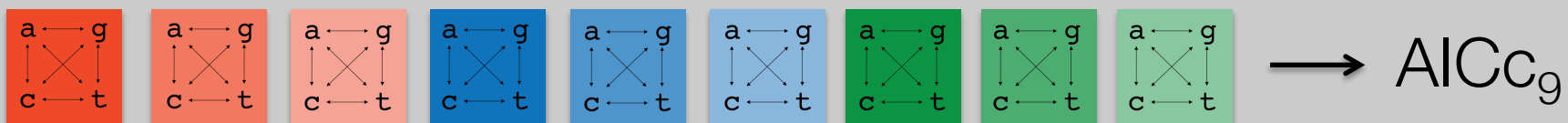


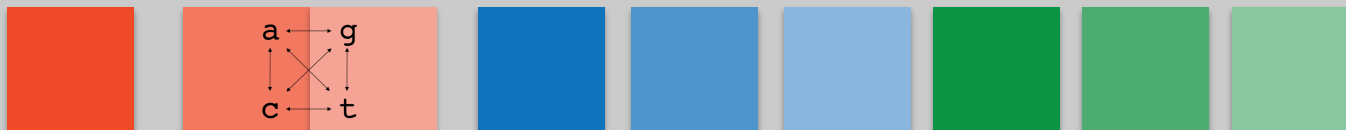
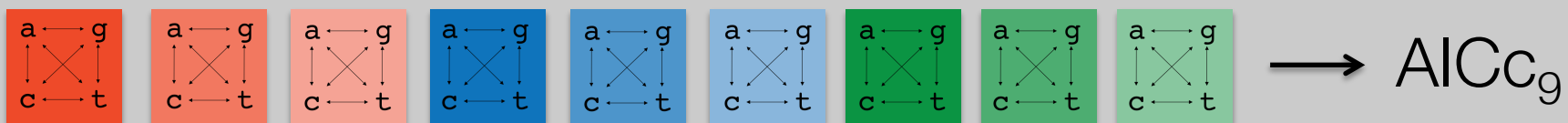


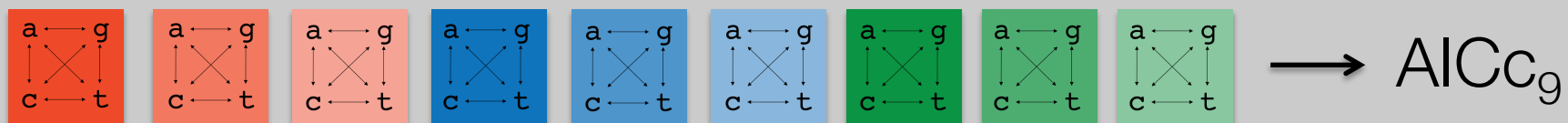


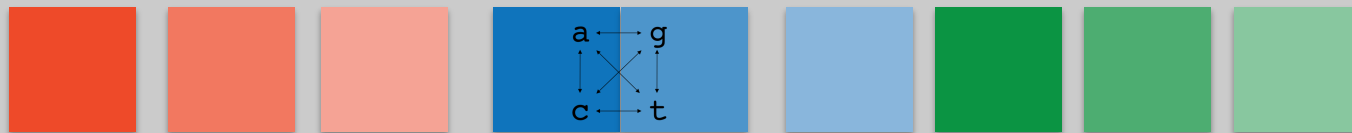
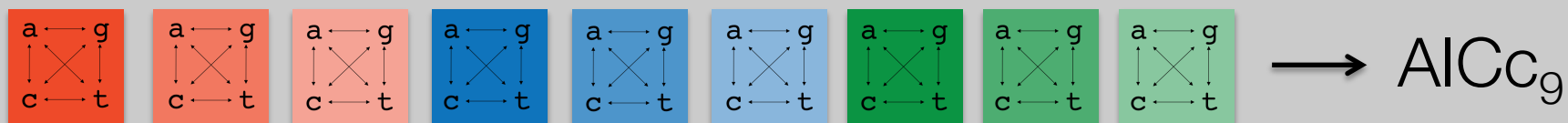


→ AlCc_9

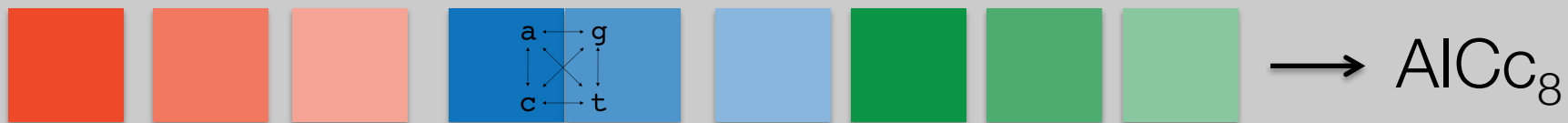
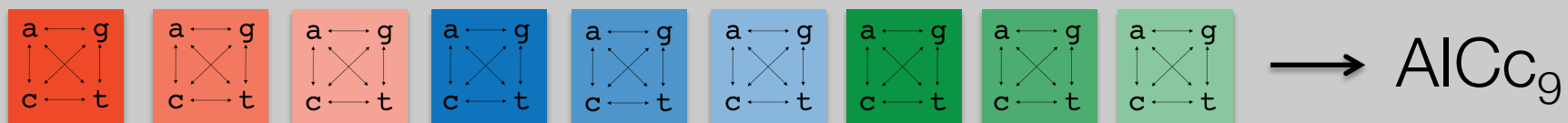


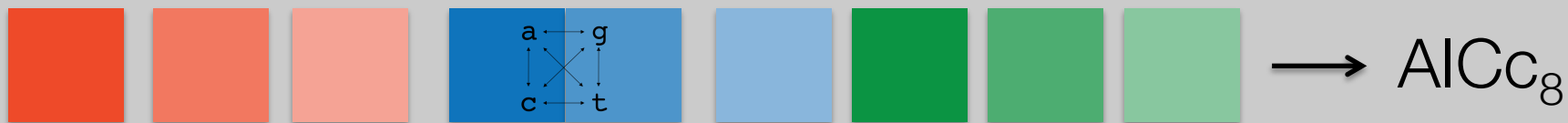
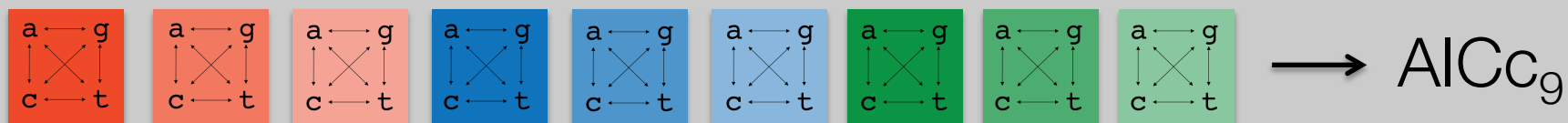




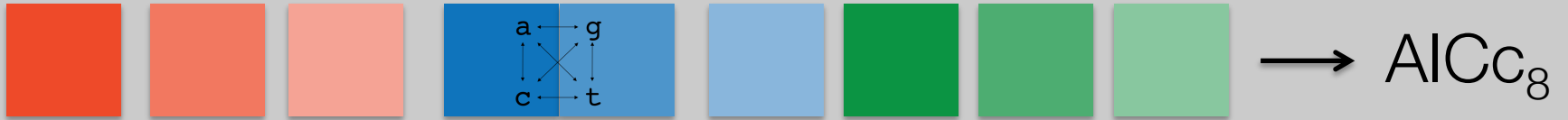
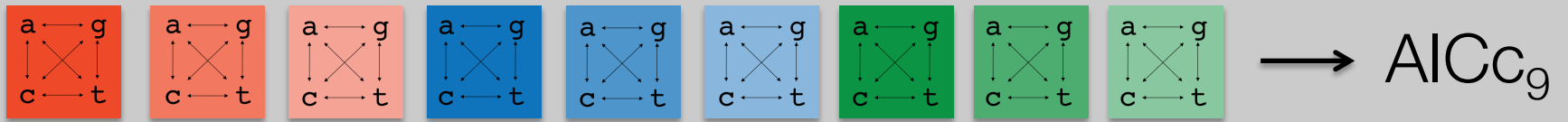


$$C_2^9 = 36$$



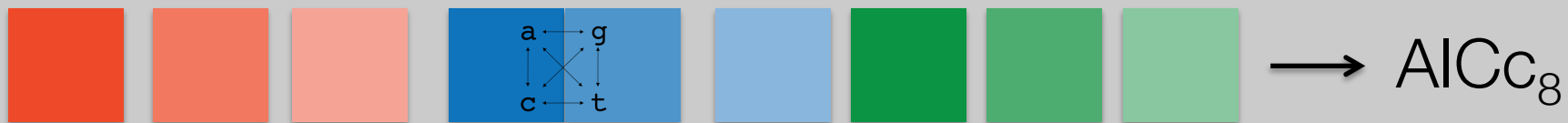
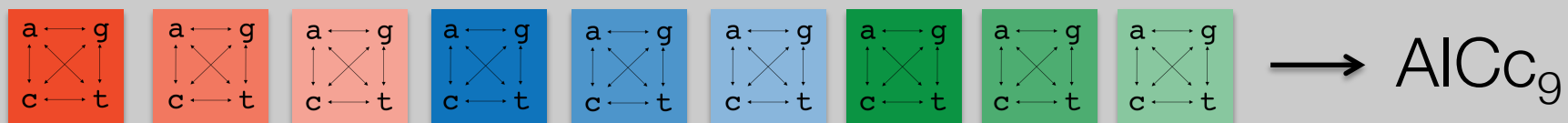


If: $AICc_8 < AICc_9$

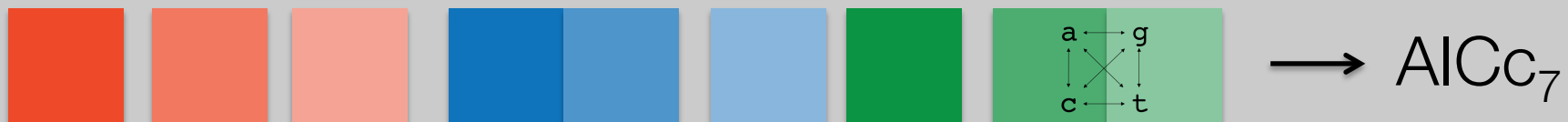


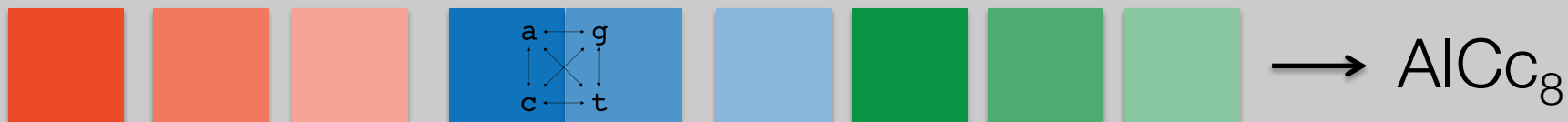
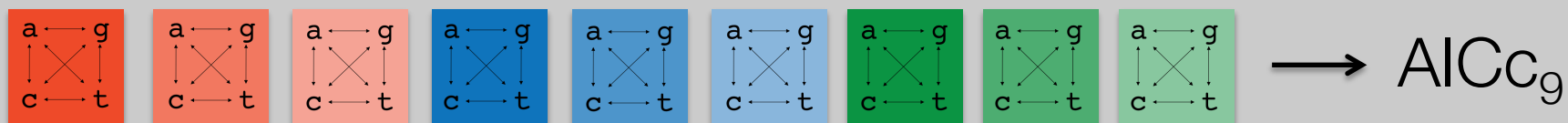
If: $AICc_8 < AICc_9$



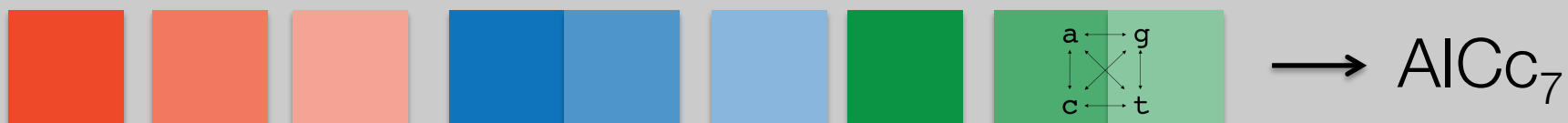


If: $AICc_8 < AICc_9$

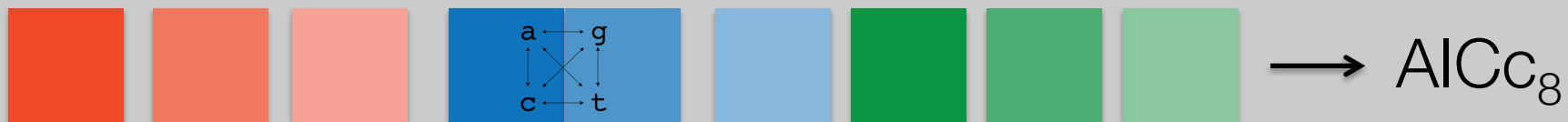
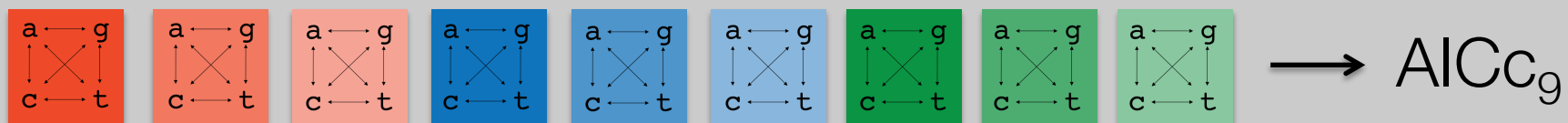




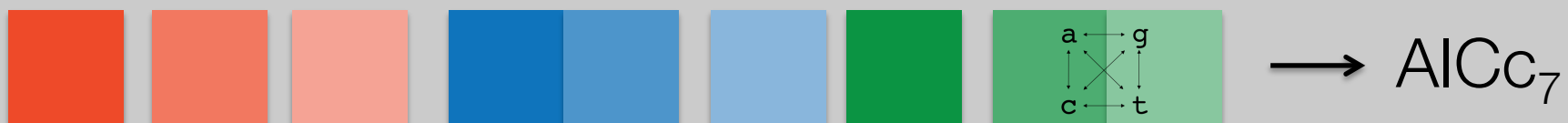
If: $AICc_8 < AICc_9$



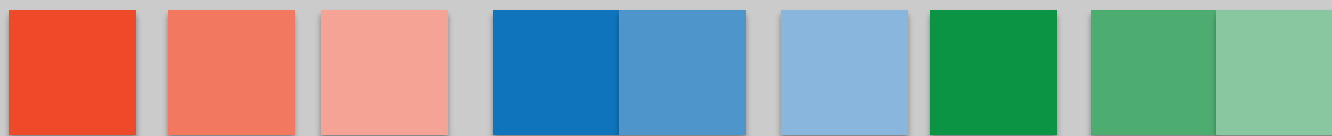
If: $AICc_7 < AICc_8$

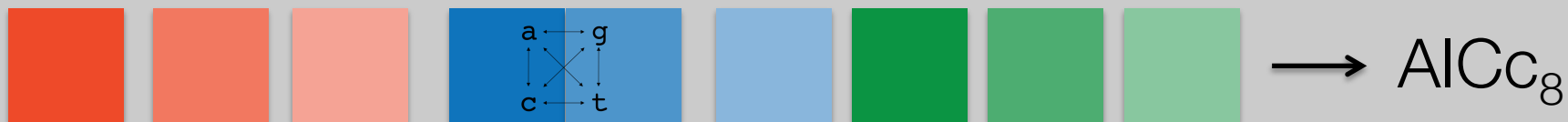
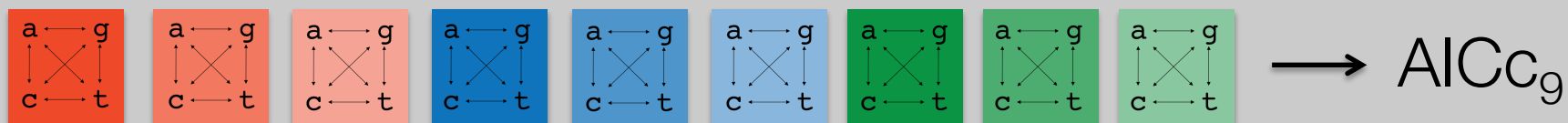


If: $AICc_8 < AICc_9$

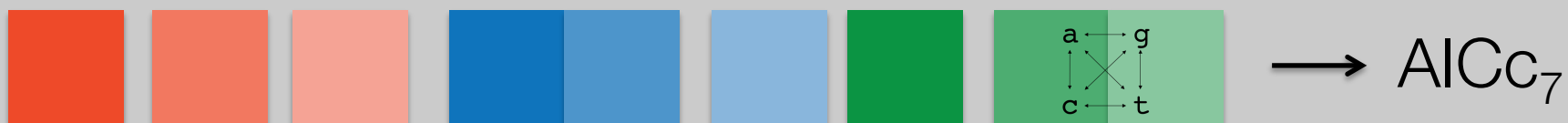


If: $AICc_7 < AICc_8$

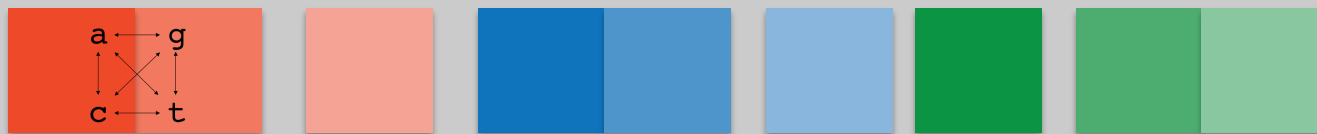


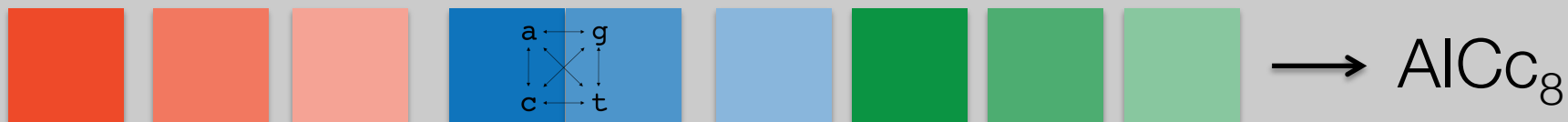
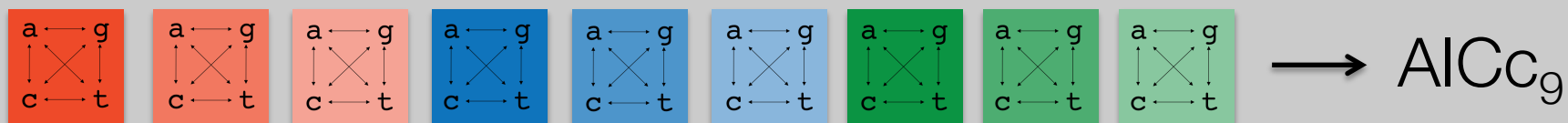


If: $AICc_8 < AICc_9$

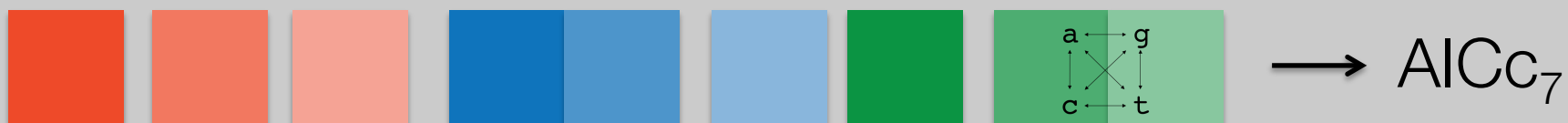


If: $AICc_7 < AICc_8$





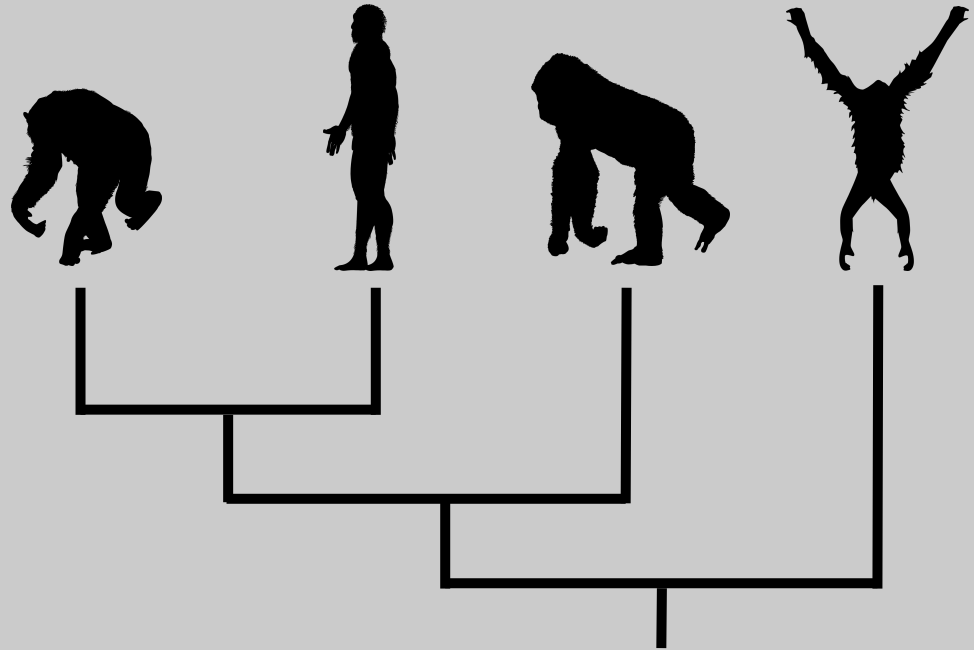
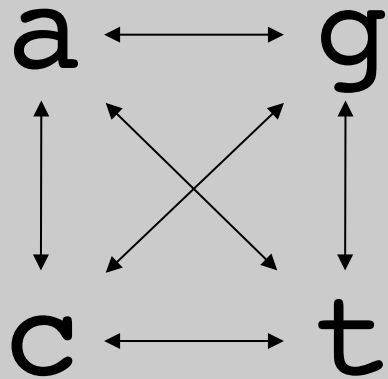
If: $AICc_8 < AICc_9$



If: $AICc_7 < AICc_8$

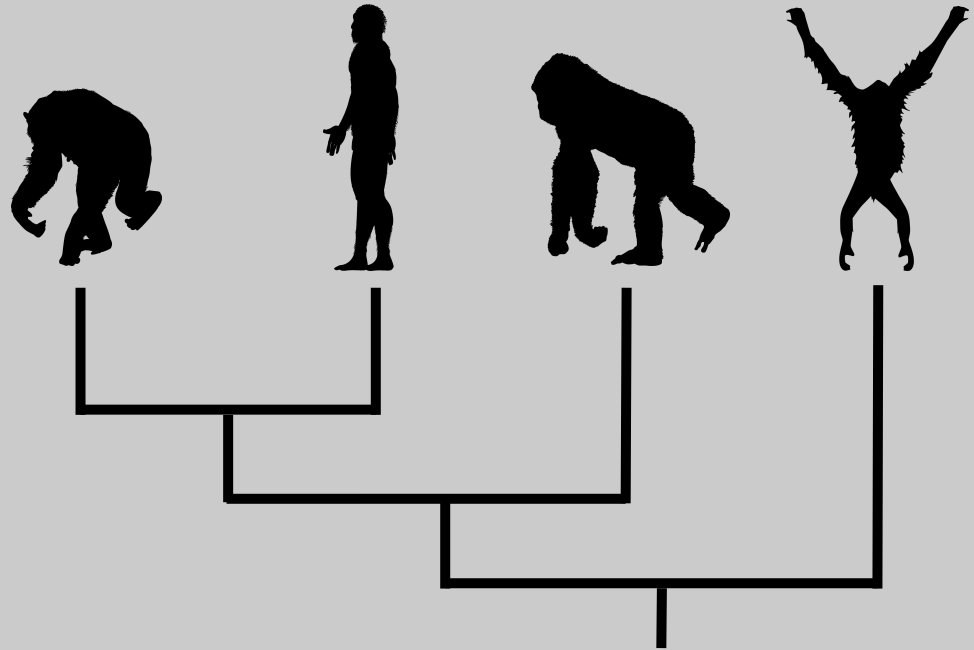
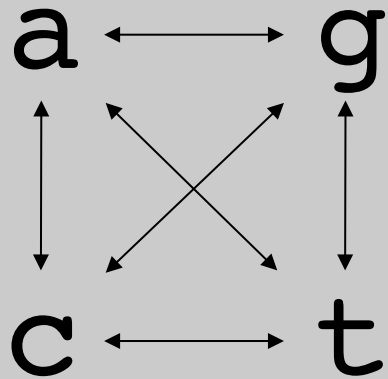


a c t g c c t g a c t g g c t g a c
g c t g a c g t c a t g g c t g a c
a g t g a c g t a g t g a g t g c g
a g c g c c t g c g c g a g c g c g

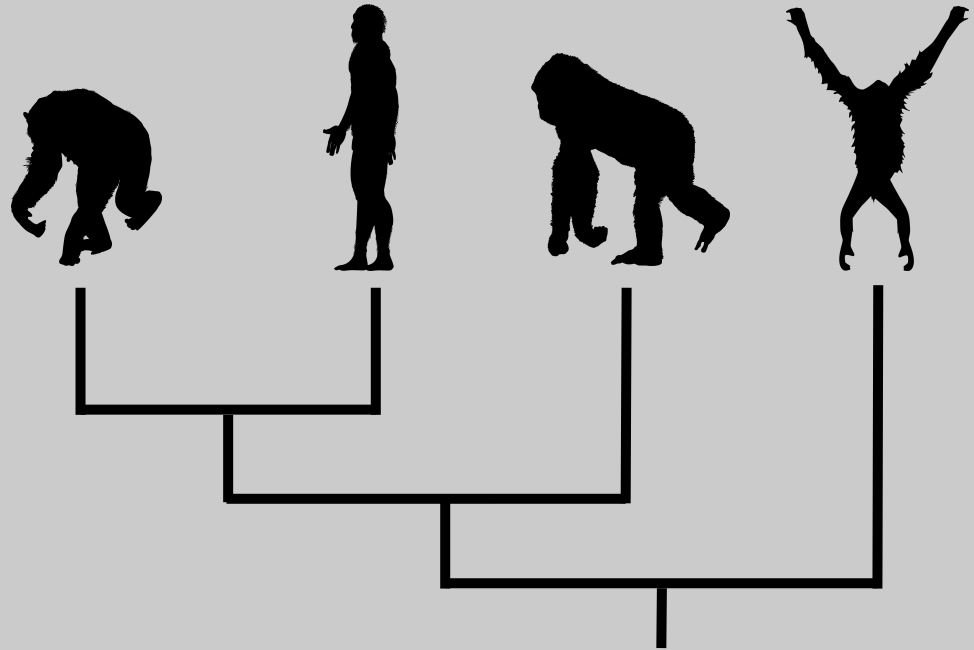
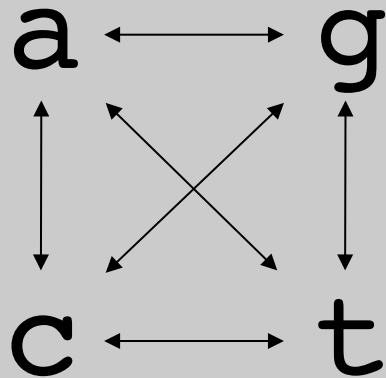


a c t g c c t g a c t g g c t g a c
 g c t g a c g t c a t g g c t g a c
 a g t g a c g t a g t g a g t g c g
 a g c g c c t g c g c g a g c g c g

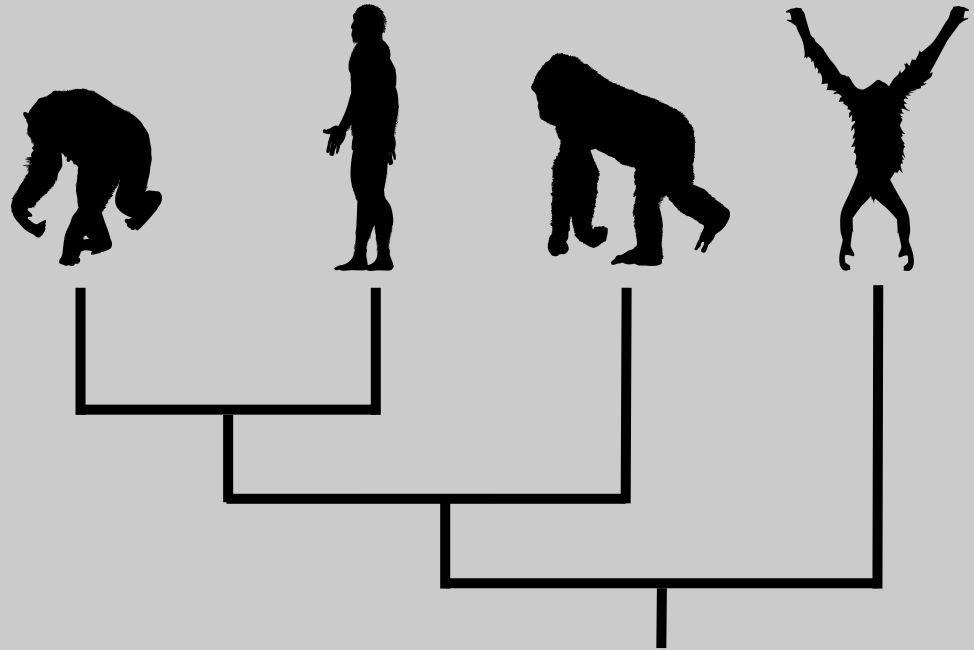
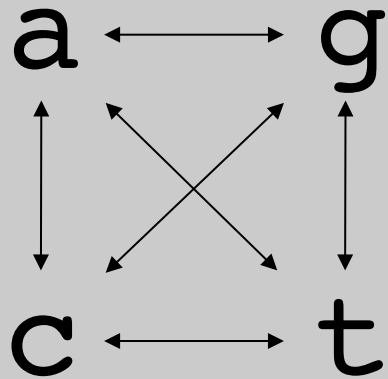
— — — — — — — —

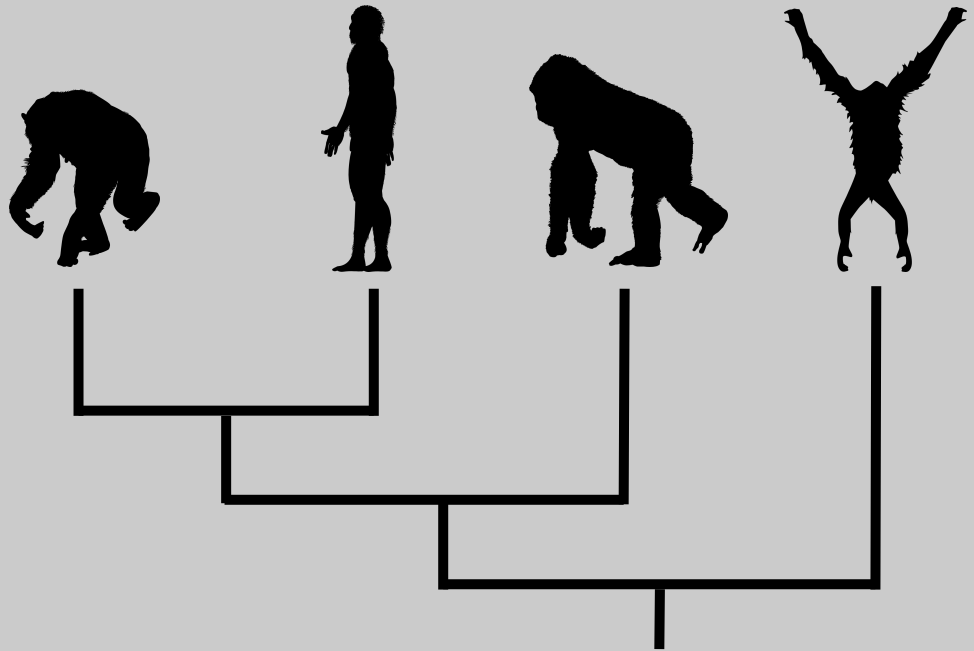


a c t g c c t g a c t g g c t g a c
 g c t g a c g t c a t g g c t g a c
 a g t g a c g t a g t g a g t g c g
 a g c g c c t g c g c g a g c g c g



a c t g c c t g a c t g g c t g a c
 g c t g a c g t c a t g g c t g a c
 a g t g a c g t a g t g a g t g c g
 a g c g c c t g c g c g a g c g c g





Algorithms available in PartitionFinder 2

Predefined data blocks:

user

greedy

hcluster

rcluster

rclusterf

Without predefined data blocks:

kmeans

PartitionFinder: Combined Selection of Partitioning Schemes and Substitution Models for Phylogenetic Analyses

Robert Lanfear,^{*}1 Brett Calcott,^{1,2} Simon Y. W. Ho,³ and Stephane Guindon⁴

greedy

METHODOLOGY ARTICLE

Open Access

Selecting optimal partitioning schemes for phylogenomic datasets

Robert Lanfear^{1,2*†}, Brett Calcott^{3†}, David Kainer¹, Christoph Mayer⁴ and Alexandros Stamatakis^{5,6}

rcluster, hcluster

METHODOLOGY ARTICLE

Open Access

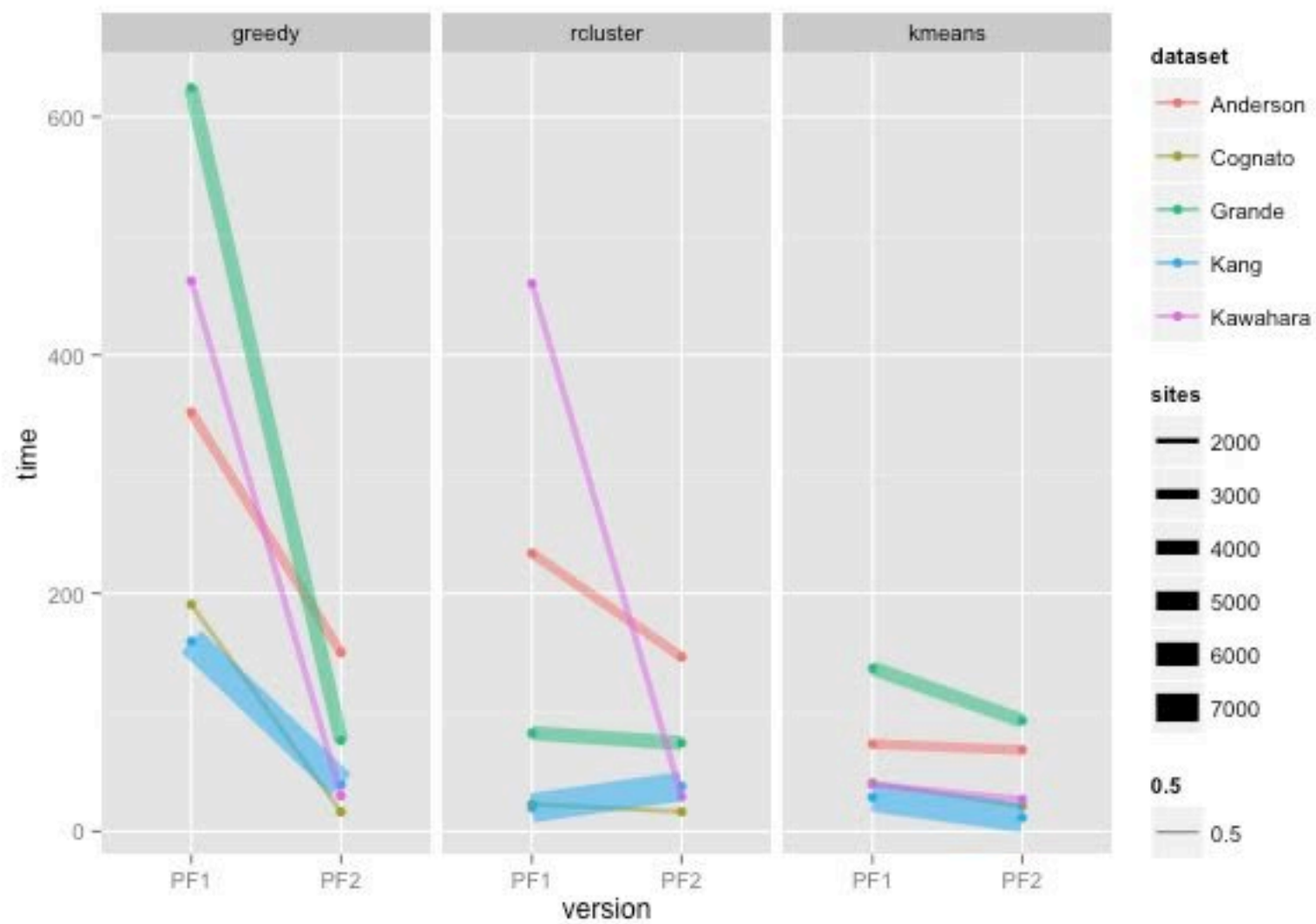
Automatic selection of partitioning schemes for phylogenetic analyses using iterative *k*-means clustering of site rates

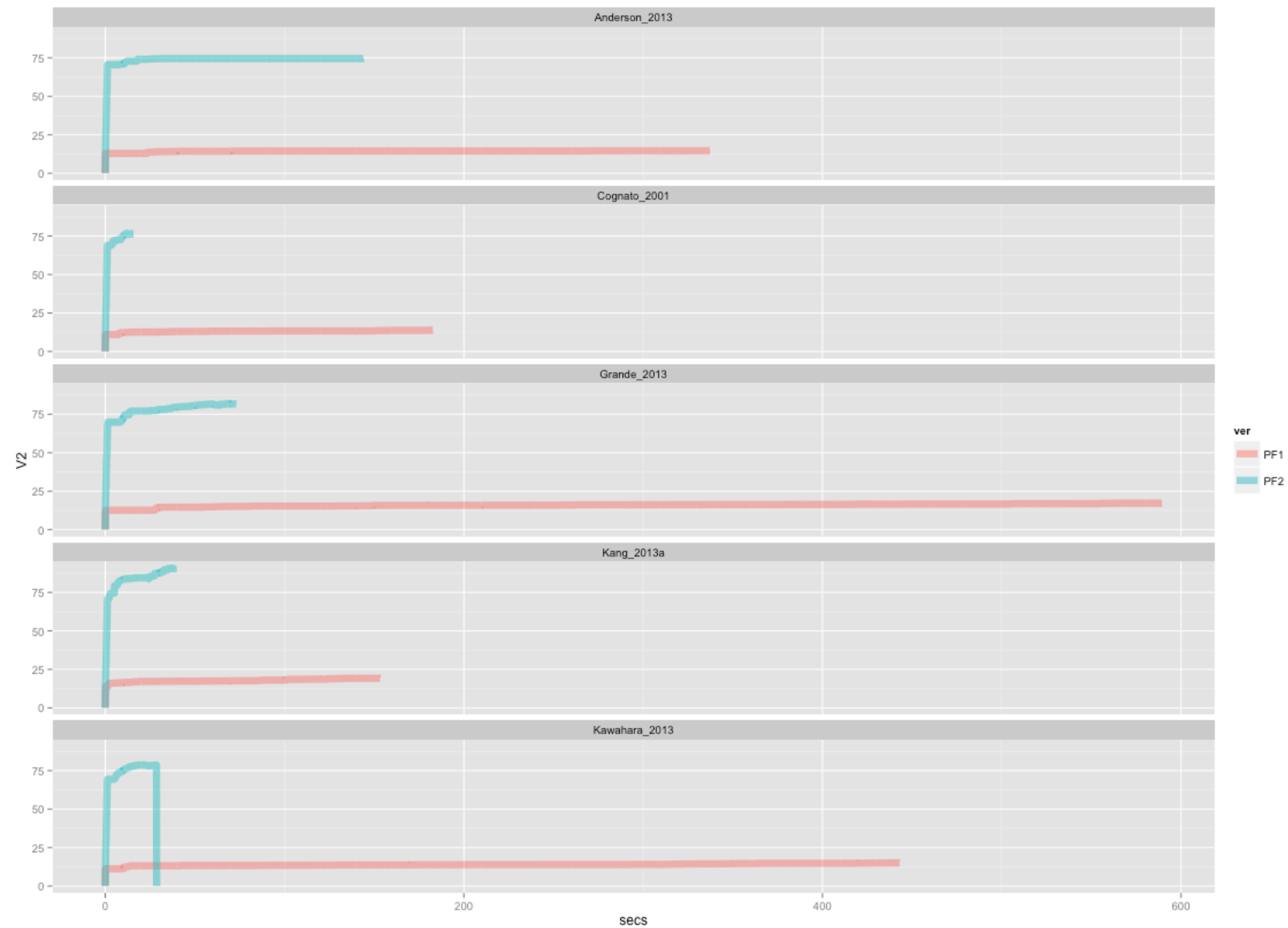
Paul B Frandsen^{1,2*}, Brett Calcott³, Christoph Mayer⁴ and Robert Lanfear^{5,6,7}

kmeans

PartitionFinder2

- Out now on GitHub!
- Installed on Hydra
- New features:
 - Uses NumPy for faster computation
 - kmeans and rclusterf algorithms are included
- And soon...
 - Morphology! (with help from April Wright of Iowa State University)





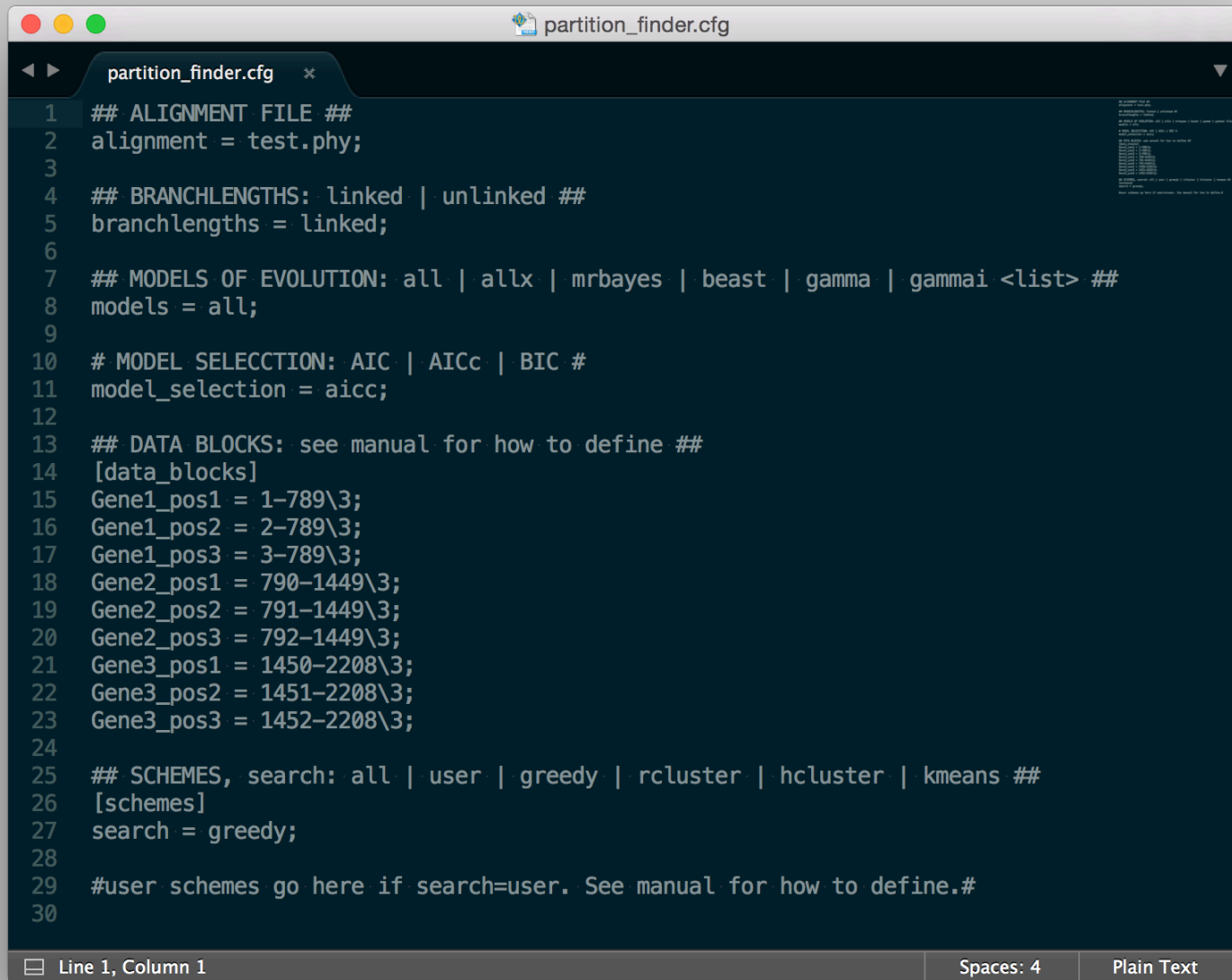
PartitionFinder support

- Excellent manual included with GitHub repo (<https://github.com/brettc/partitionfinder>)
- Google group: <https://groups.google.com/forum/#!forum/partitionfinder>
- Questions specific to Hydra installation: contact me

Tutorial

- ssh into Hydra
- Navigate to your /pool directory
- Copy /pool/genomics/frandsenp/examples/nucleotide to your /pool directory
- Change your directory into the nucleotide directory
- Inside should be a partition_finder.cfg file and an alignment called test.phy

The PartitionFinder config file



The image shows a screenshot of a text editor window titled "partition_finder.cfg". The editor has a dark blue background with light green text. The file content is as follows:

```
1  ## ALIGNMENT FILE ##
2  alignment = test.phy;
3
4  ## BRANCHLENGTHS: linked | unlinked ##
5  branchlengths = linked;
6
7  ## MODELS OF EVOLUTION: all | allx | mrbayes | beast | gamma | gammai <list> ##
8  models = all;
9
10 # MODEL SELECTION: AIC | AICc | BIC #
11 model_selection = aicc;
12
13 ## DATA BLOCKS: see manual for how to define ##
14 [data_blocks]
15 Gene1_pos1 = 1-789\3;
16 Gene1_pos2 = 2-789\3;
17 Gene1_pos3 = 3-789\3;
18 Gene2_pos1 = 790-1449\3;
19 Gene2_pos2 = 791-1449\3;
20 Gene2_pos3 = 792-1449\3;
21 Gene3_pos1 = 1450-2208\3;
22 Gene3_pos2 = 1451-2208\3;
23 Gene3_pos3 = 1452-2208\3;
24
25 ## SCHEMES, search: all | user | greedy | rcluster | hcluster | kmeans ##
26 [schemes]
27 search = greedy;
28
29 #user schemes go here if search=user. See manual for how to define.#
30
```

The status bar at the bottom of the editor shows "Line 1, Column 1", "Spaces: 4", and "Plain Text".

Tutorial

- Open the Hydra QSubGen web app (<https://hydra-3.si.edu/tools/QSubGen>)
- Select the short queue
- Select multi-thread and 4 CPUs
- Select “sh” for your shell
- Load the bioinformatics/partitionfinder/2.0pr13 module

Tutorial

- For PartitionFinder on Hydra you must call the PartitionFinder.py executable directly.
Note: this contradicts the instructions in the user manual
- Enter the command: PartitionFinder.py
nucleotide --raxml -p \$NSLOTS

Tutorial

- Enter the command: `PartitionFinder.py nucleotide --raxml -p $NSLOTS`
- `--raxml` specifies that you will use `raxml` for your likelihood calculations
- `-p` specifies the number of CPU threads to use. Note: this is very important to specify on Hydra because if you do not specify it, the default in `PartitionFinder` is to use all CPUs.

Tutorial

- Choose an informative job name
- Select -cwd
- Either download and upload your job file to Hydra or copy and paste it into your favorite text editor on Hydra. I save my job file as “pf_nuc.job”.