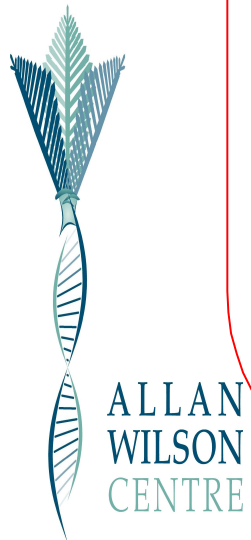# 'Realism' and 'Instrumentalism'

# in models of

# molecular evolution
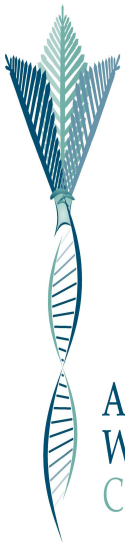
David Penny

Montpellier, June 08

Galileo

# Overview

sites free to vary

summing sources of error

'rates' of molecular evolution

estimates of time intervals

do we know anything? (flat priors)
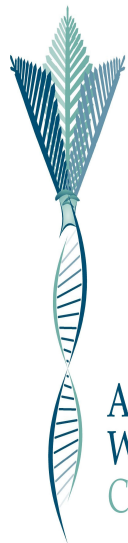
ALLAN
WILSON
CENTRE

# Human/chimp divergence

1) *Ramapithecus* = 12Ma → HC = 5±1Ma

But *Ramapithecus* in Asia, HCG in Africa.
Is 18-20Ma a better estimate for divergence?

2) *Ramapithecus* = 18Ma → HC = 7.5±1.5Ma

Or should we combine uncertainties?

*In this case, I would rather not – leave it as a conditional estimate – need both.*

# sites free to vary
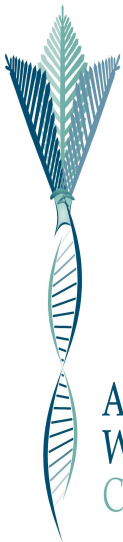
| | rate $k_{aa} \times 10^9$/yr |
|---|---|
| - fibrinopeptides | 8.3 |
| - lysozyme | 2.0 |
| - hemoglobin $\alpha$ | 1.2 |
| - cytochrome c | 0.3 |
| - histone H4 | 0.01 |

Dickerson, 1971

explained the differences by the proportion of sites 'free to vary'.

change of function should show a rate change

realism

# we use a tiny fraction
# of the information in the data

Alignment                          Reordered Alignment
original sequence order              shuffled/reordered

```
AIIFLNSALGPSPELFPIILATKVL    ASAGPSPPATPLLIIIILLFFNEKV
AIMFLNSALGPPTELFPVILATKVL    ASAGPPTPATPLLIMVILLFFNEKV
SIMFLNHTLNPTPELFPIILATETL    SHTNPTPPATPLLIMIILLFFNEET
TILFLNSSLGLQPEVTPTVLATKTL    TSSGLQPPATPLLILTVLVTFNEKT
TLLFLNSMLKPPSELFPIILATKTL    TSMKPPSPATPLLLLIILLFFNEKT
ALLFLNSTLNPPTELFPLILATKTL    ASTNPPTPATPLLLLLILLFFNEKT
AILFLNSFLNPPKEFFPIILATKIL    ASFNPPKPATPLLILIILFFFNEKI
```
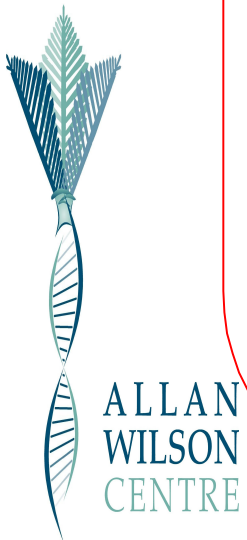
*c* columns

*c*! alignments

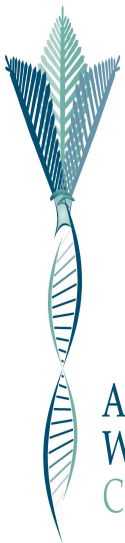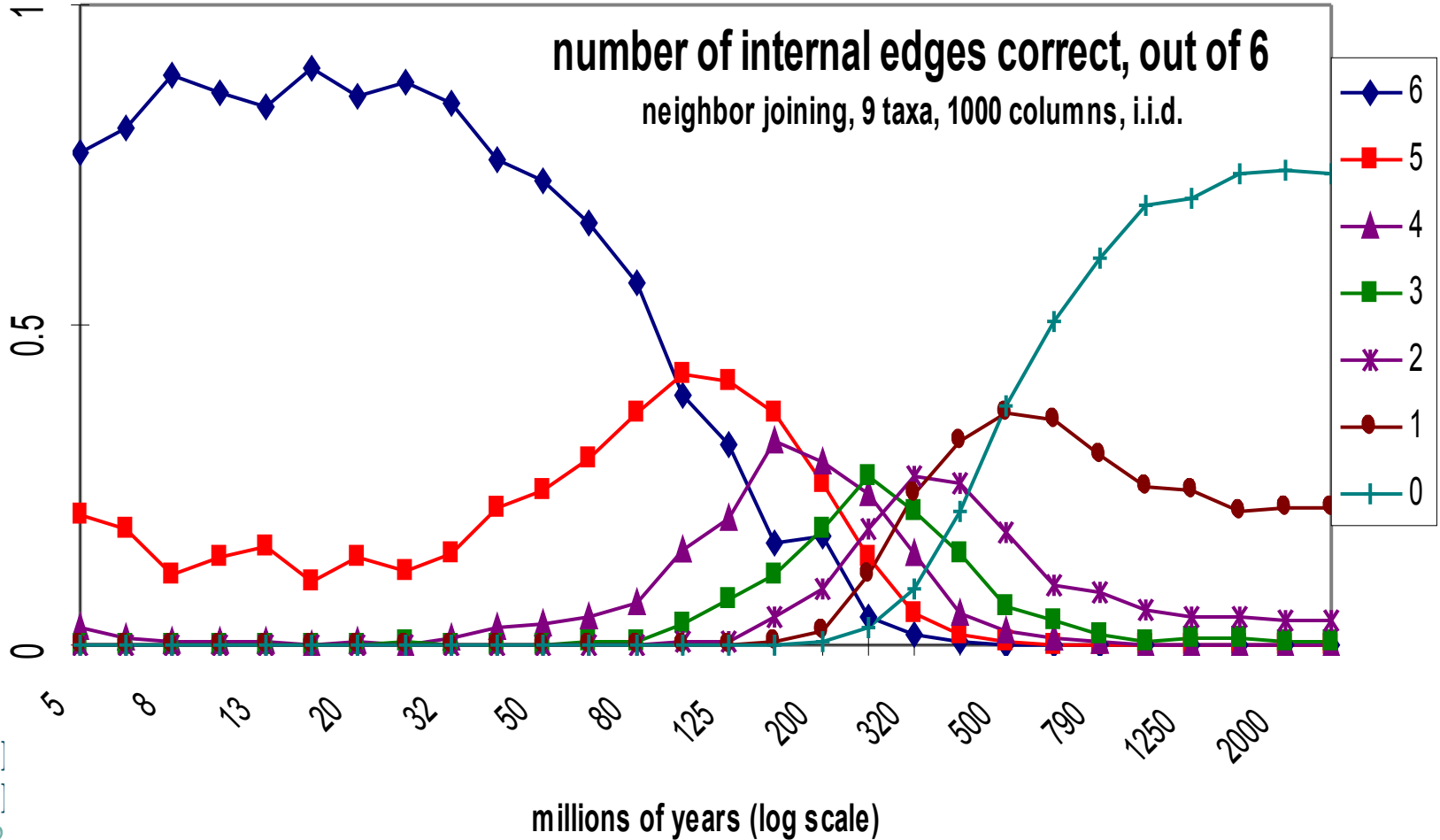If *c* = 1000, we use $\approx$ 1/ 1000! of the information

# sites change

X-ray crystallographers: the strongest conclusion we have is that the same sites in different species may be fixed, in others they are variable.
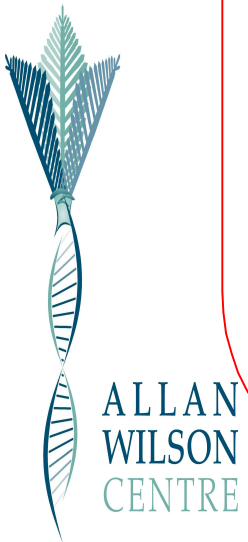
Molecular Phylogeneticists: Our methods (such as the Gamma distribution) assume sites are in the SAME rate class across the entire tree (AND, we only need one parameter- so there).

ALLAN
WILSON
CENTRE

# simulation results with standard model



**number of internal edges correct, out of 6**

**neighbor joining, 9 taxa, 1000 columns, i.i.d.**

Legend:
- 6
- 5
- 4
- 3
- 2
- 1
- 0

x-axis: **millions of years (log scale)** — 5, 8, 13, 20, 32, 50, 80, 125, 200, 320, 500, 790, 1250, 2000

y-axis: 0, 0.5, 1

ALLAN WILSON CENTR

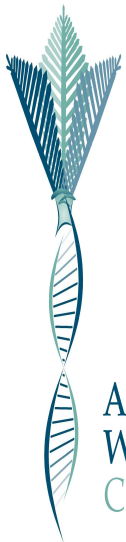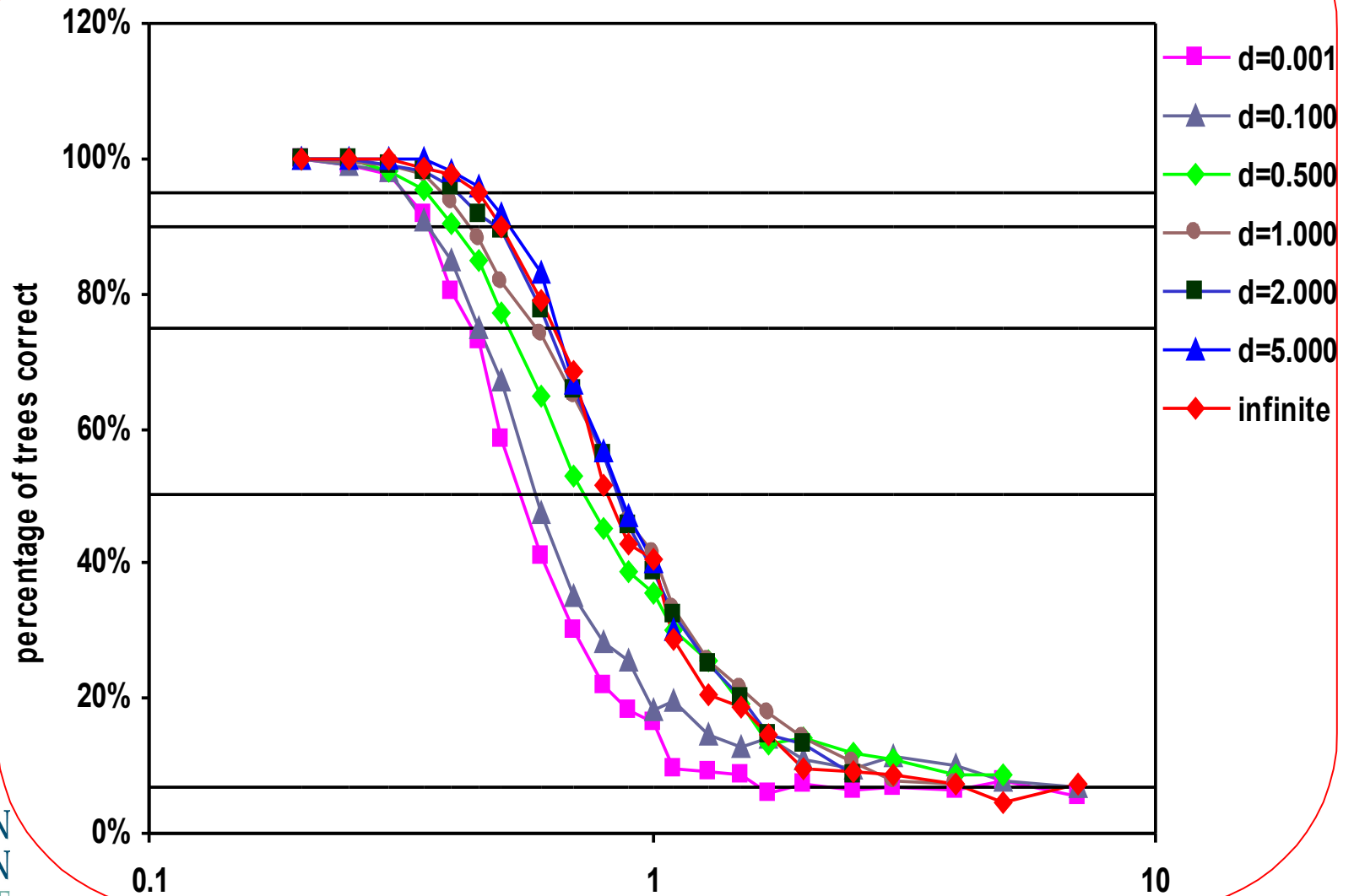# Calculated results, $\Delta \leq \frac{1}{4} + ne^{-qt}$

simulation results with covarion model

# do 'rates' exist !!!

We go ON

       and ON

             and ON

                  and ON

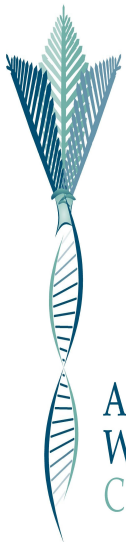About 'molecular clocks'.

Should we??

# not enough information to recover the full model

$$\begin{bmatrix} 1-\gamma & \gamma \\ \delta & 1-\delta \end{bmatrix}$$

$1\, (P_R,\, 1-P_R)$ *composition at root*

2

2

Seq 1

Seq 2

5 required,

3 available

# two taxa, two codes

Seq 1   -----------------------------------------------------------
Seq 2   -----------------------------------------------------------

$$
\text{Seq 1} \quad
\begin{array}{c}
R \\
\\
Y
\end{array}
\left[
\begin{array}{cc}
\alpha & \beta \\
\\
\gamma & *
\end{array}
\right]
\\
\qquad\quad R \qquad Y \\
\qquad\quad \text{Seq 2}
$$

| 1 | 2 | |
|---|---|---|
| R | R | $\alpha$ |
| R | Y | $\beta$ |
| Y | R | $\gamma$ |
| Y | Y | $*$ |

Divergence matrix, $F_{i,j}$

Three independent
parameters estimated

# three taxa

$$\begin{bmatrix} 1-\gamma & \gamma \\ \delta & 1-\delta \end{bmatrix}$$

1 $(P_R, 1-P_R)$

2    2    2

Seq 1    Seq 2    Seq 3

7 required

# four character states

$$\begin{bmatrix} * & \alpha & \beta & \gamma \\ \delta & * & \varepsilon & \phi \\ \eta & \iota & * & \varphi \\ \kappa & \lambda & \mu & * \end{bmatrix}$$

3    $(P_R, 1- P_R)$

12    12    12

Seq 1    Seq 2    Seq 3

39 required

ALLAN
WILSON
CENTRE

# tensor, 3D matrix

$$
\begin{bmatrix}
0.274950 & 0.007961 & 0.003838 & 0.000711 \\
0.009667 & 0.023742 & 0.002985 & 0.000426 \\
0.001848 & 0.001848 & 0.015496 & 0.000853 \\
0.000569 & 0.000142 & 0.001564 & 0.002132
\end{bmatrix}
$$

64 − 1 = 63 values, but a sparse matrix!

# primary diagonal

Gymnure,  Mole and Shrew

| | | T | C | A | G |
|---|---|---|---|---|---|
| **T** | **T** | **0.274950** | 0.007961 | 0.003838 | 0.000711 |
| T | C | 0.009667 | 0.023742 | 0.002985 | 0.000426 |
| T | A | 0.001848 | 0.001848 | 0.015496 | 0.000853 |
| T | G | 0.000569 | 0.000142 | 0.001564 | 0.002132 |
| | | | | | |
| C | T | 0.011231 | 0.006682 | 0.000995 | 0.000426 |
| **C** | **C** | 0.010520 | **0.188371** | 0.001564 | 0.000426 |
| C | A | 0.001137 | 0.002275 | 0.006682 | 0.000426 |
| C | G | 0.000284 | 0.000569 | 0.000853 | 0.000995 |
| | | | | | |
| A | T | 0.007819 | 0.002701 | 0.004265 | 0.000284 |
| A | C | 0.002985 | 0.009383 | 0.004407 | 0.000426 |
| **A** | **A** | 0.003838 | 0.004834 | **0.201166** | 0.003554 |
| A | G | 0.000426 | 0.000853 | 0.005118 | 0.007819 |
| | | | | | |
| G | T | 0.001279 | 0.000071 | 0.000071 | 0.000853 |
| G | C | 0.000142 | 0.001990 | 0.000284 | 0.000284 |
| G | A | 0.000284 | 0.000284 | 0.004691 | 0.001137 |
| **G** | **G** | 0.000995 | 0.000711 | 0.001279 | **0.143588** |
| | | **T** | **C** | **A** | **G** |

ALLAN
WILSON
CENTRE

# secondary diagonals

|   |   | Gymnure(moon rat) | Mole, | Shrew |   |
|---|---|---|---|---|---|
| **T** | T | **0.274950** | 0.007961 | 0.003838 | 0.000711 |
| **T** | C | 0.009667 | **0.023742** | 0.002985 | 0.000426 |
| **T** | A | 0.001848 | 0.001848 | **0.015496** | 0.000853 |
| **T** | G | 0.000569 | 0.000142 | 0.001564 | **0.002132** |
| **C** | T | **0.011231** | 0.006682 | 0.000995 | 0.000426 |
| **C** | C | 0.010520 | **0.188371** | 0.001564 | 0.000426 |
| **C** | A | 0.001137 | 0.002275 | **0.006682** | 0.000426 |
| **C** | G | 0.000284 | 0.000569 | 0.000853 | **0.000995** |
| **A** | T | **0.007819** | 0.002701 | 0.004265 | 0.000284 |
| **A** | C | 0.002985 | **0.009383** | 0.004407 | 0.000426 |
| **A** | A | 0.003838 | 0.004834 | **0.201166** | 0.003554 |
| **A** | G | 0.000426 | 0.000853 | 0.005118 | **0.007819** |
| **G** | T | **0.001279** | 0.000071 | 0.000071 | 0.000853 |
| **G** | C | 0.000142 | **0.001990** | 0.000284 | 0.000284 |
| **G** | A | 0.000284 | 0.000284 | **0.004691** | 0.001137 |
| **G** | G | 0.000995 | 0.000711 | 0.001279 | **0.143588** |
|   |   | **T** | **C** | **A** | **G** |

# moon rat, 1+2

|   | T | C | A | G |
|---|---|---|---|---|
| **T** | 0.955 | 0.148 | 0.087 | 0.028 |
| **C** | 0.025 | 0.803 | 0.025 | 0.009 |
| **A** | 0.018 | 0.043 | 0.876 | 0.076 |
| **G** | 0.002 | 0.006 | 0.012 | 0.887 |

|   | T | C | A | G |
|---|---|---|---|---|
| **T** | .955 ±.004 | .150 ±.013 | .087 ±.009 | .029 ±.008 |
| **C** | .025 ±.003 | .800 ±.014 | .025 ±.005 | .009 ±.003 |
| **A** | .018 ±.003 | .044 ±.006 | .877 ±.011 | .077 ±.011 |
| **G** | .002 ±.001 | .006 ±.002 | .012 ±.002 | .886 ±.015 |

therefore we believe in symmetric models
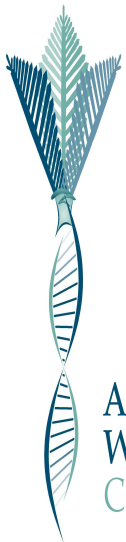
# mole, shrew and moon rat

| mole | | T | C | A | G |
|------|---|-------|-------|-------|-------|
| | T | 0.976 | **0.062** | 0.021 | 0.013 |
| | C | 0.017 | 0.931 | 0.020 | 0.007 |
| | A | **0.006** | 0.006 | 0.948 | **0.012** |
| | G | 0.001 | 0.001 | 0.010 | 0.968 |
| | | T | C | A | G |

| shrew | | T | C | A | G |
|-------|---|-------|-------|-------|-------|
| | T | 0.977 | **0.038** | 0.024 | 0.011 |
| | C | 0.020 | 0.951 | 0.020 | 0.003 |
| | A | **0.002** | 0.009 | 0.942 | **0.011** |
| | G | 0.001 | 0.001 | 0.015 | 0.976 |

| moon rat | | T | C | A | G |
|----------|---|-------|-------|-------|-------|
| | T | 0.955 | **0.148** | 0.087 | 0.028 |
| | C | 0.025 | 0.803 | 0.025 | 0.009 |
| | A | **0.018** | 0.043 | 0.876 | **0.076** |
| | G | 0.002 | 0.006 | 0.012 | 0.887 |
| | | T | C | A | G |

ALLAN
WILSON
CENTRE

# change in rate

$$\begin{bmatrix} * & \alpha & \beta & \gamma \\ \delta & * & \varepsilon & \phi \\ \eta & \iota & * & \varphi \\ \kappa & \lambda & \mu & * \end{bmatrix}$$

$$\begin{bmatrix} * & \alpha & \beta & \gamma \\ \delta & * & \varepsilon & \phi \\ \eta & \iota & * & \varphi \\ \kappa & \lambda & \mu & * \end{bmatrix}$$

$$\begin{bmatrix} * & \alpha & \beta & \gamma \\ \delta & * & \varepsilon & \phi \\ \eta & \iota & * & \varphi \\ \kappa & \lambda & \mu & * \end{bmatrix}$$

$$\begin{bmatrix} * & \alpha & \beta \\ \gamma & & \\ \delta & * & \varepsilon & \phi \\ \eta & \iota & * & \varphi \\ \kappa & \lambda & \mu & * \end{bmatrix}$$

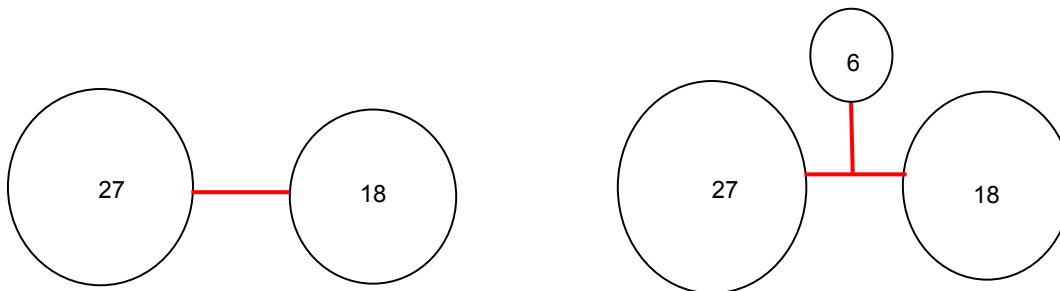change in process

# do we know anything?
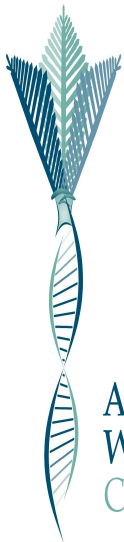
the curse of 'flat priors'

the 'we know nothing syndrome'

ALLAN
WILSON
CENTRE
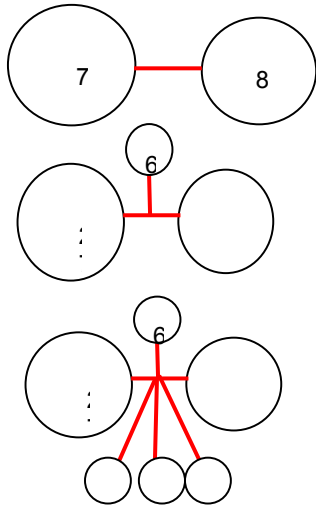
# binary trees, b(n) = (2n-5)!!
= 1 x 3 x 5 x 7 … 2n-5.

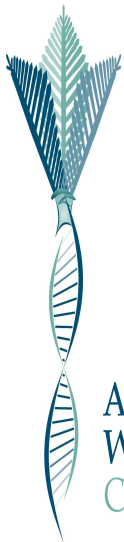$5.68 \times 10^{-18}$

# binary trees, $b(n) = (2n-5)!!$

$$= 1 \times 3 \times 5 \times 7 \ldots 2n-5.$$
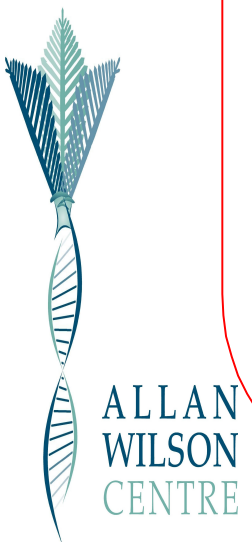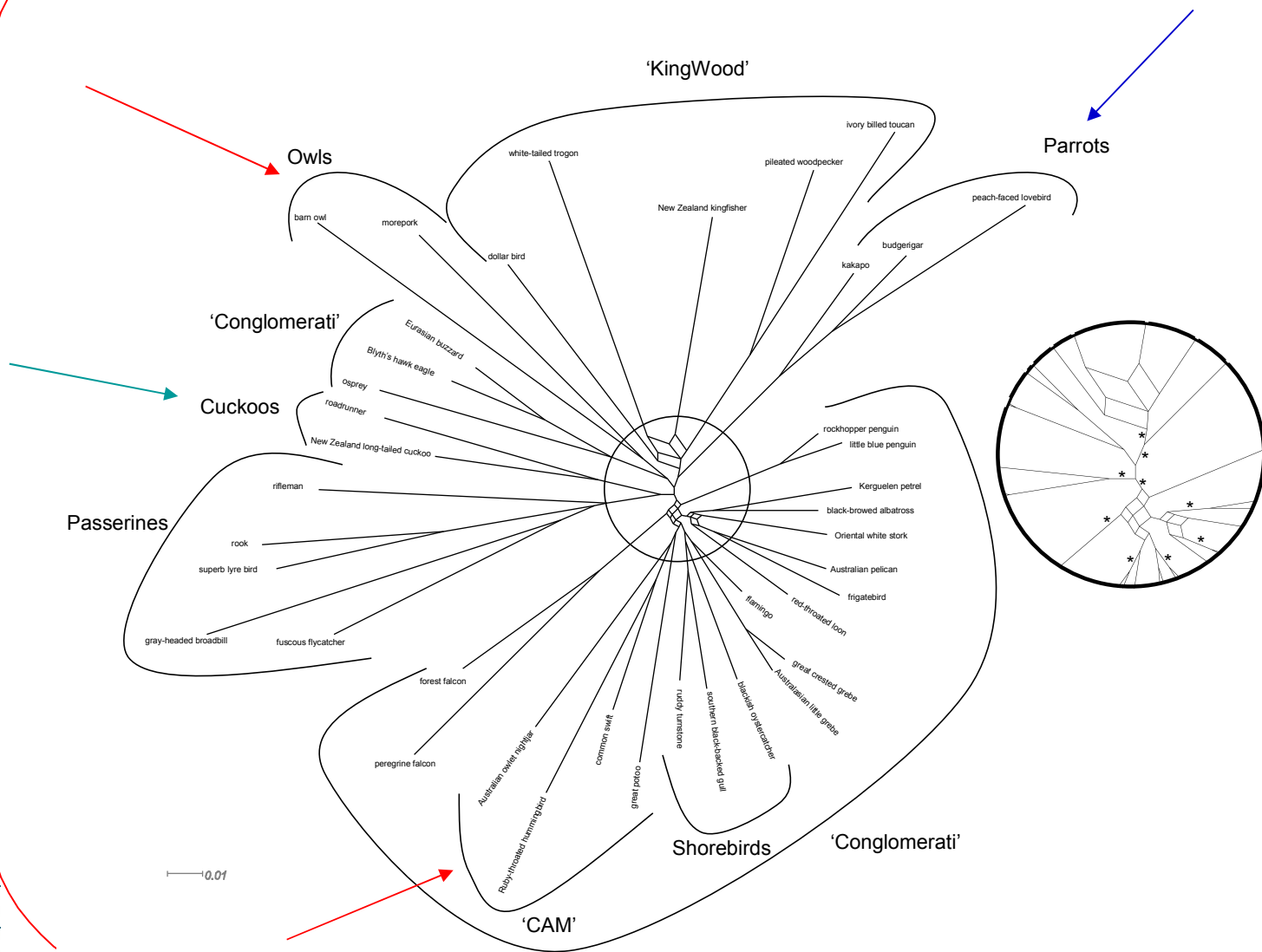


$b(n_1+1).b(n_2+1) / b(n_t)$

$b(n_1+1).b(n_2+1).b(n_3+1) / b(n_t)$

$b(n_1+1).b(n_2+1) \ldots b(n_i+1) / b(n_t)$

# 40 birds

'KingWood'

Owls

Parrots

ivory billed toucan

white-tailed trogon

pileated woodpecker

New Zealand kingfisher

peach-faced lovebird

barn owl          morepork

budgerigar

dollar bird

kakapo

'Conglomerati'

Eurasian buzzard

Blyth's hawk eagle

Cuckoos

osprey

roadrunner

New Zealand long-tailed cuckoo

rockhopper penguin

little blue penguin

Passerines

rifleman

Kerguelen petrel

black-browed albatross

rook

Oriental white stork

superb lyre bird

Australian pelican

frigatebird

gray-headed broadbill          fuscous flycatcher

flamingo          red-throated loon

great crested grebe

forest falcon

Australasian little grebe

peregrine falcon

Australian owlet-nightjar

common swift

ruddy turnstone

southern black-backed gull

blackish oystercatcher

Ruby-throated hummingbird

great potoo

0.01

Shorebirds

'Conglomerati'

'CAM'

ALLAN
WILSON
CENTRE

$$P(n,k) = \frac{R(k) \times B(n-k+1)}{B(n)}$$

probability with $n$ taxa of observing a prespecified clade of size $k$.
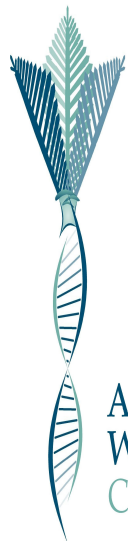
with $n$ = 40 and

$k$ = 2, P ≈ 0.013  cuckoo,roadrunner

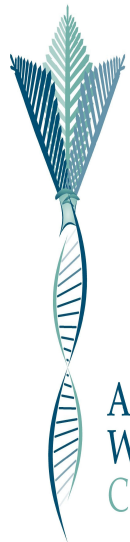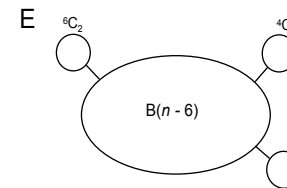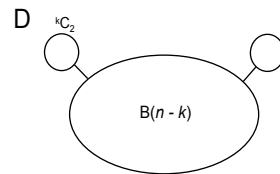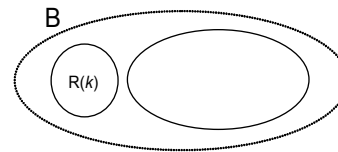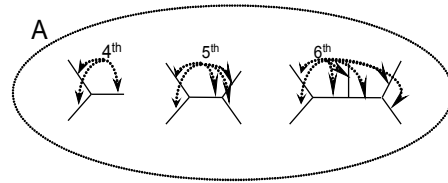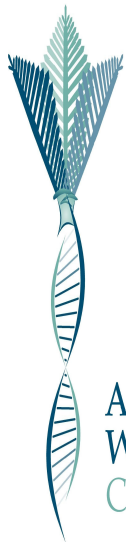$k$ = 3, P ≈ 0.0026  parrots

$k$ = 4, P ≈ $7.12 \times 10^{-6}$,

$k$ = 5, P ≈ $5.84 \times 10^{-8}$.

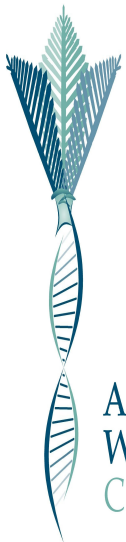potoo, owlet-nightjar, owl, barn owl, swift, hummingbird (6)

# Where next in Phylogeny?

allow realism in phylogeny

set the biological question
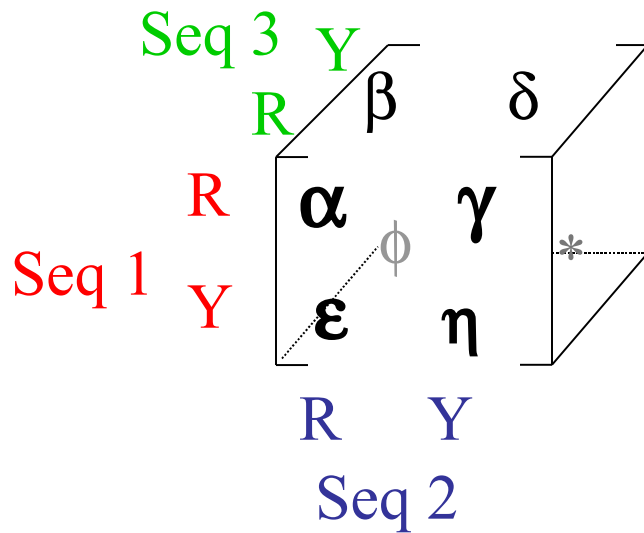
we have some bad failures

we need a range of alternatives

ALLAN
WILSON
CENTRE

Belief is the curse of the thinking class

# tensor, 2-states

Seq 1 --------R------------------------

Seq 2 --------R------------------------

Seq 3 --------R------------------------

Seq 3  Y
       R    β        δ

R          **α**      **γ**
                φ              *
Seq 1
Y          **ε**      **η**

       R        Y

       Seq 2

R R R α
R R Y β
R Y R γ
R Y Y δ
Y R R ε
Y R Y φ
Y Y R η
Y Y Y *
1  2  3

**7 available !**