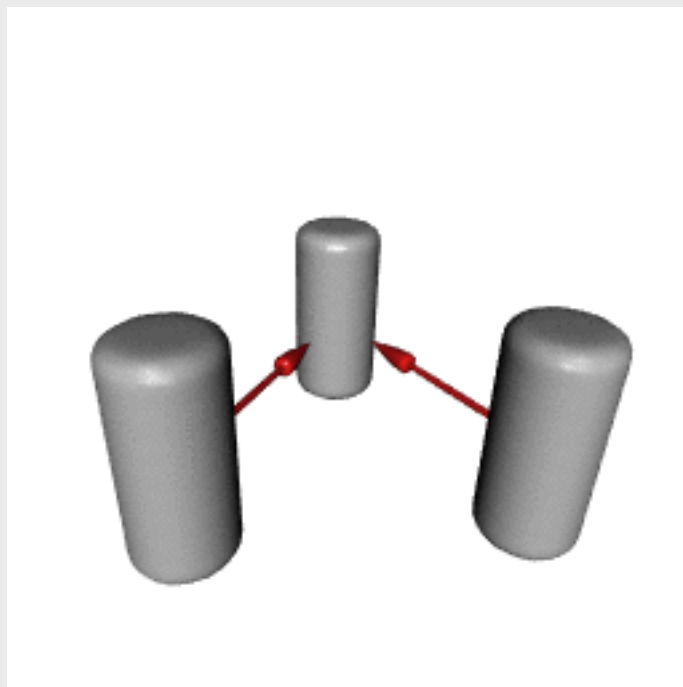


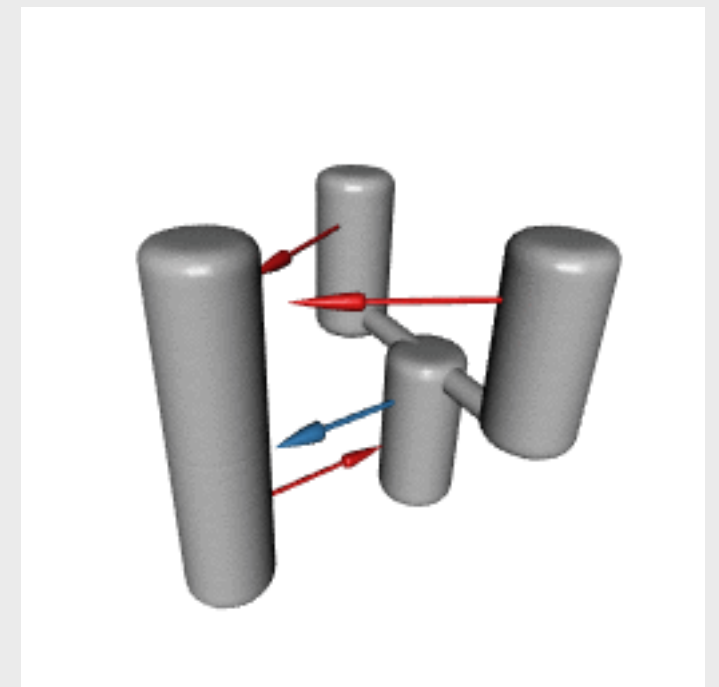
Simulated likelihood for species delimitation and phylogeography

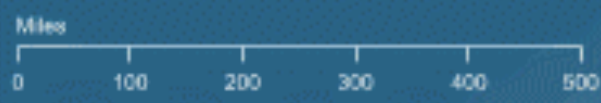


Brian O'Meara
UT Knoxville



Bryan Carstens
Ohio State U

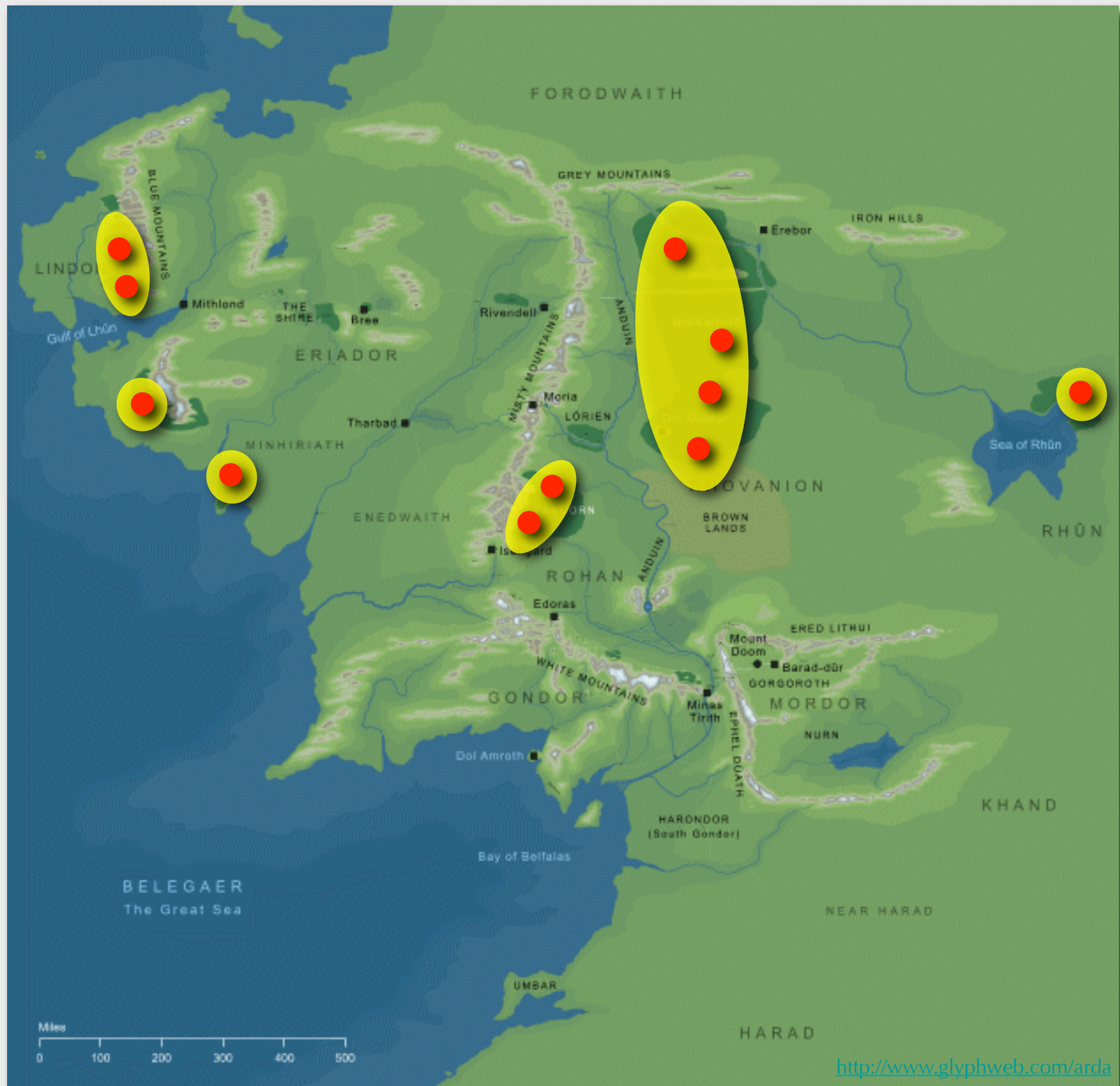


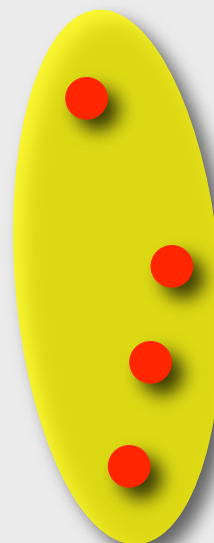
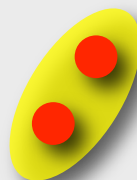
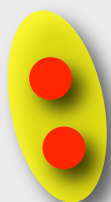


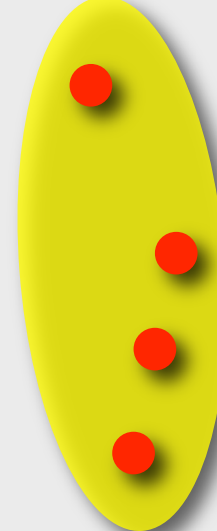
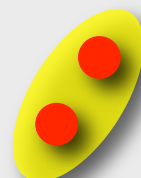
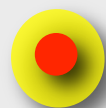
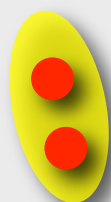


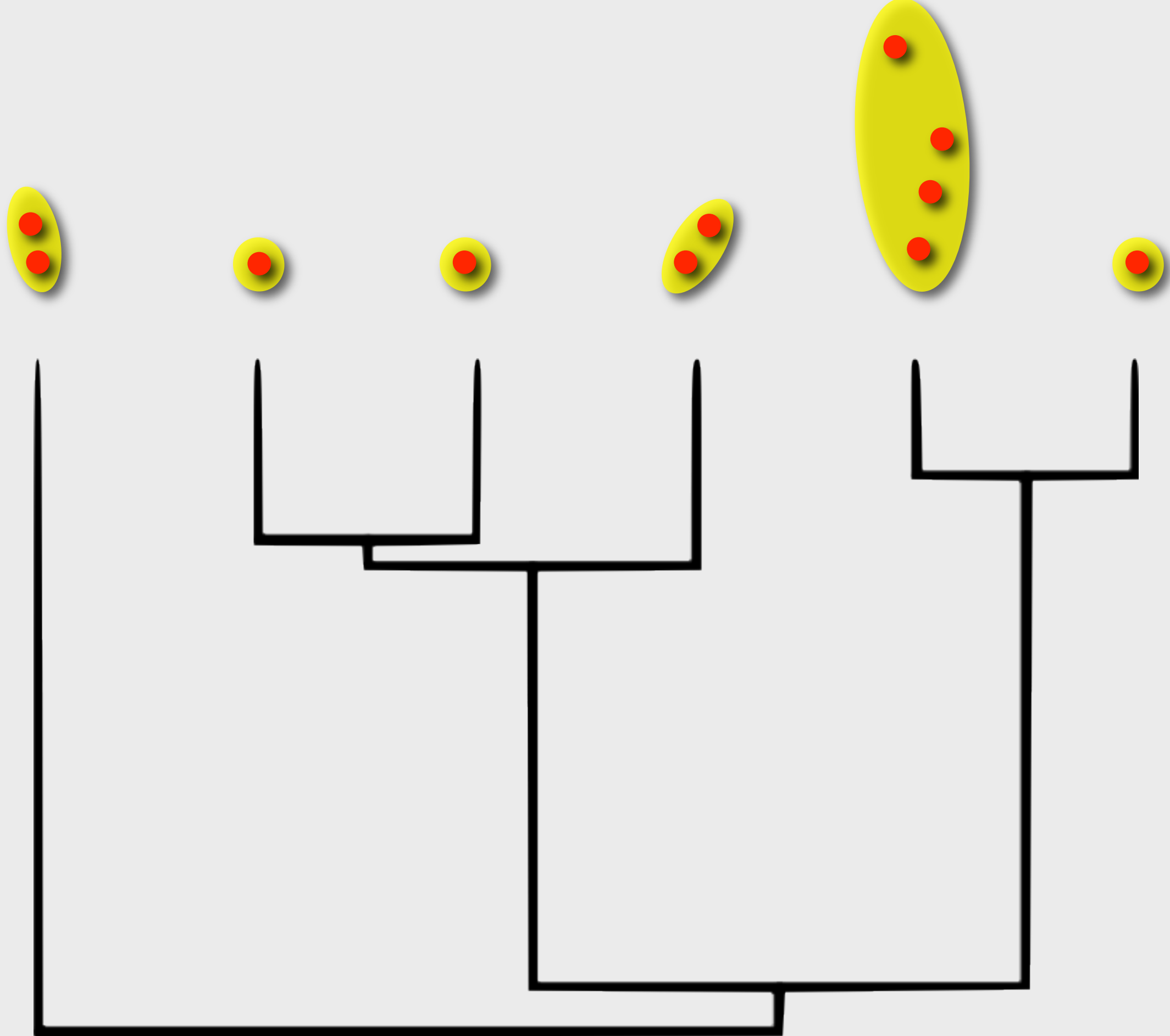


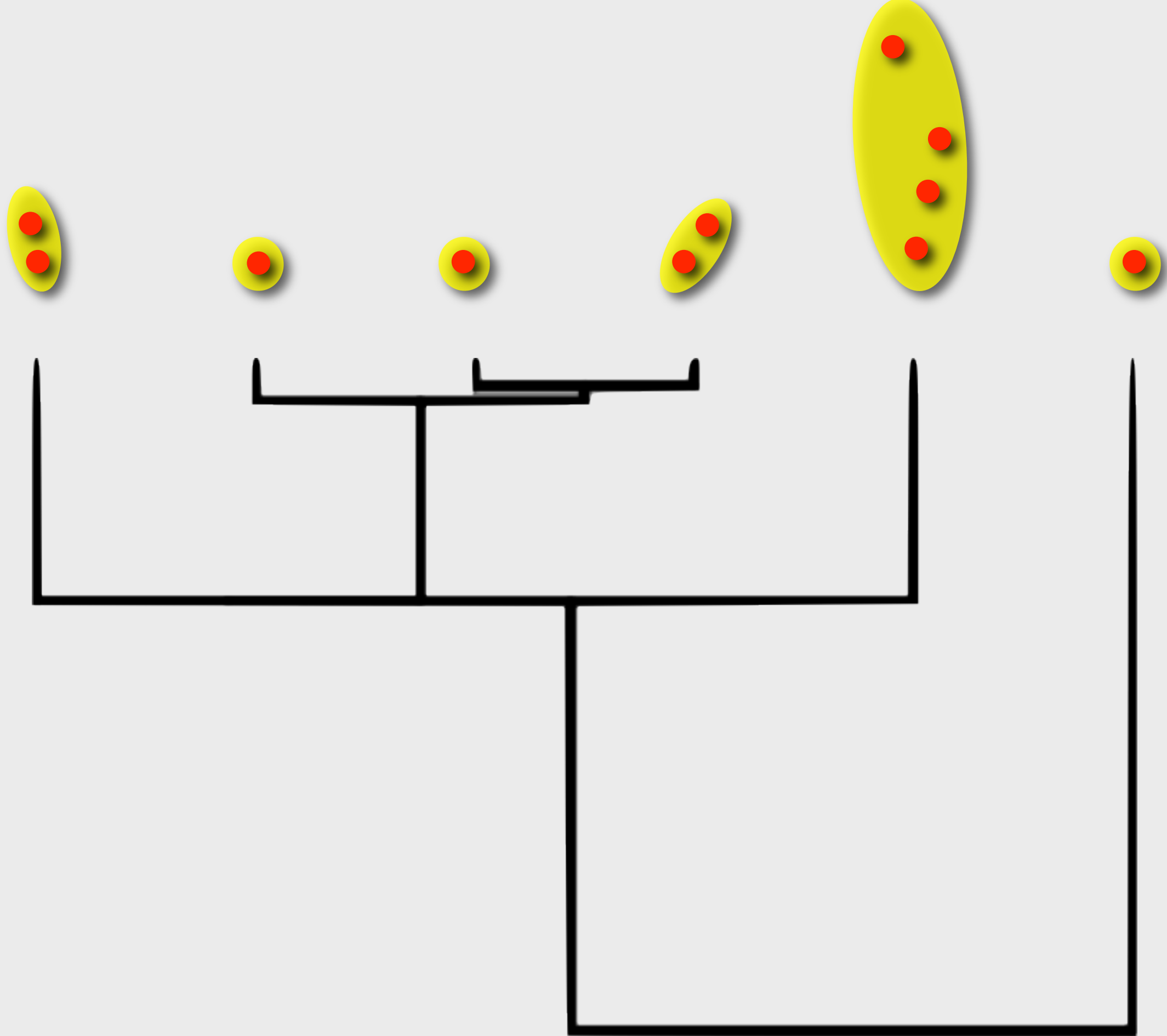


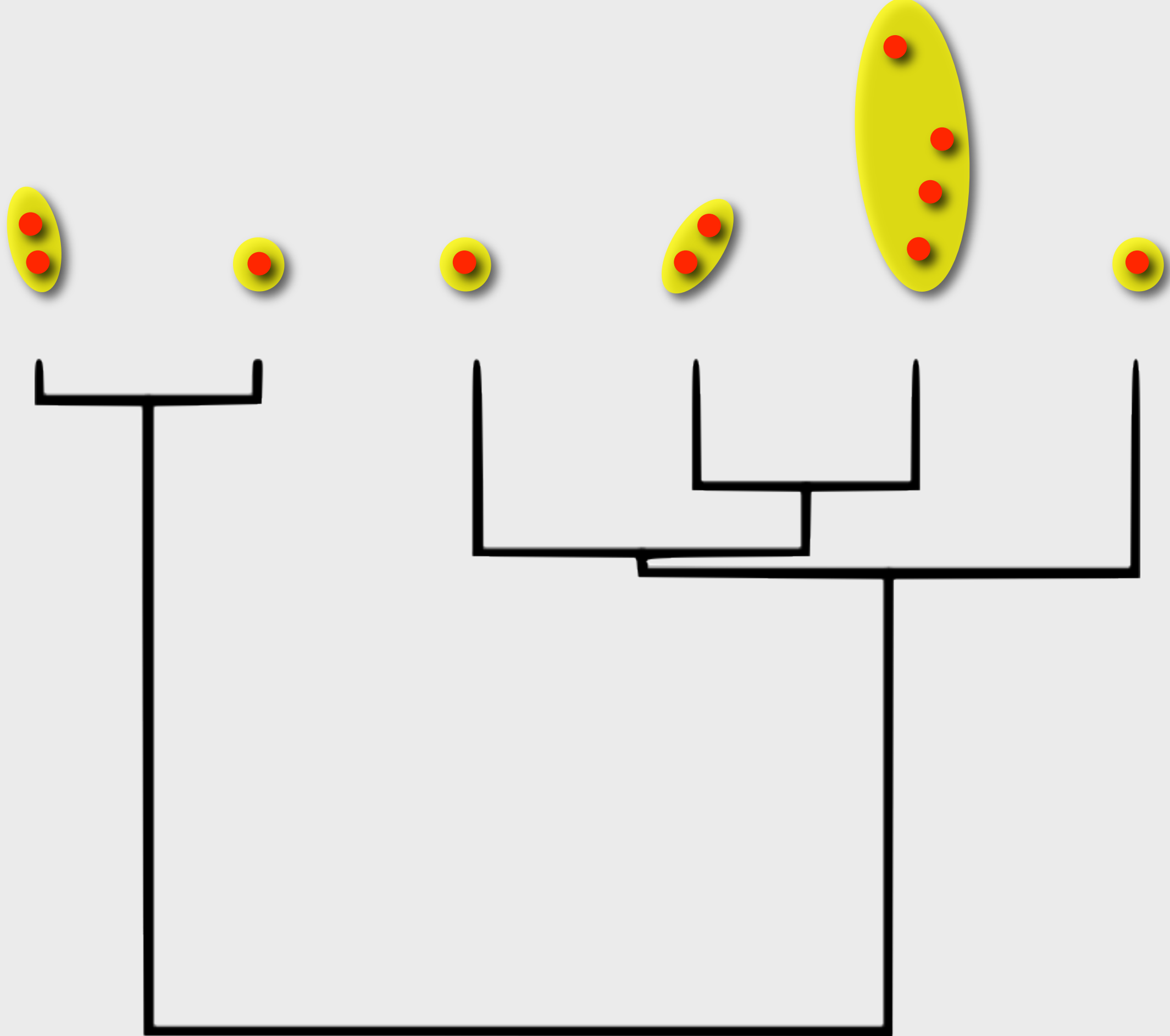


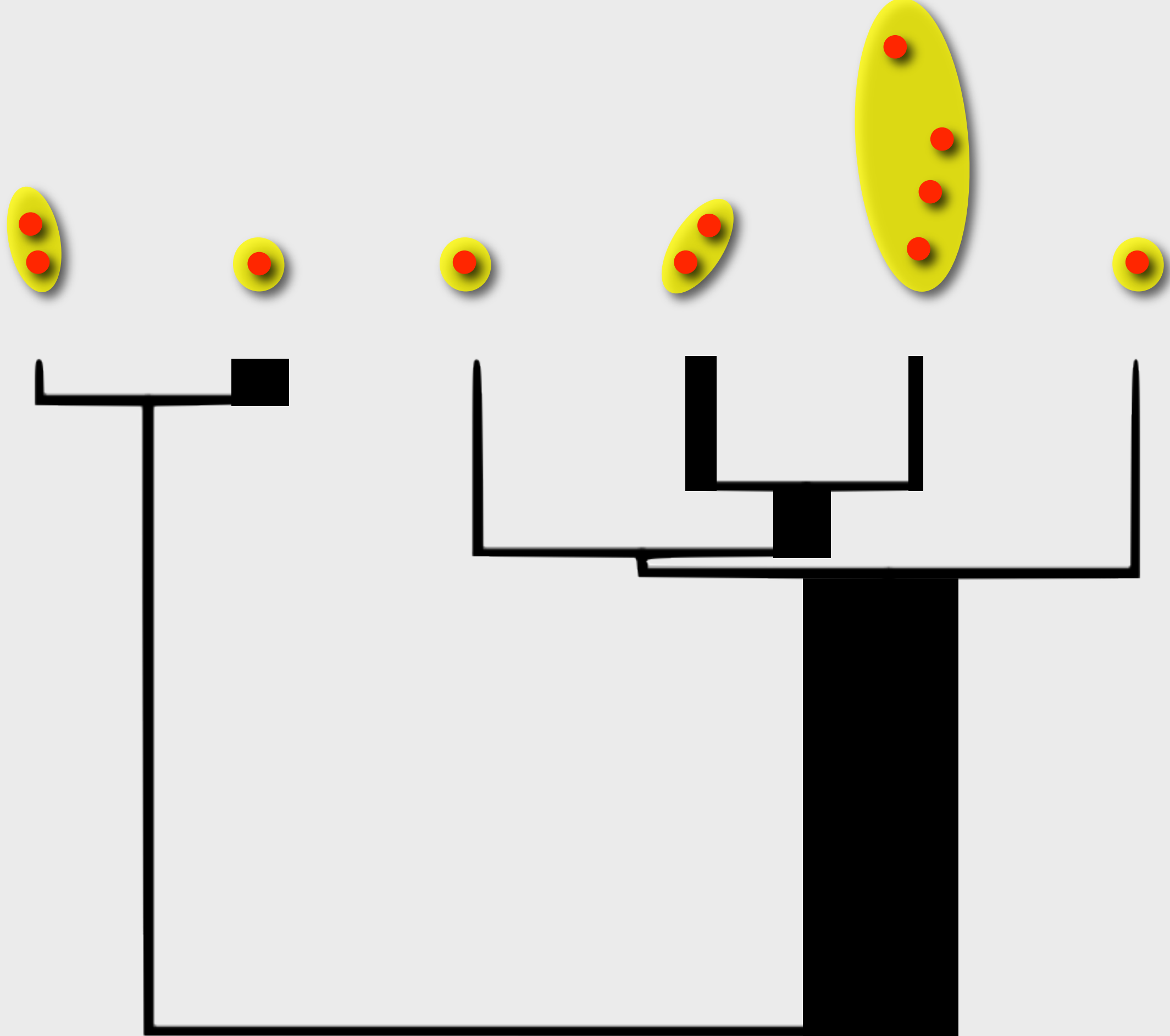


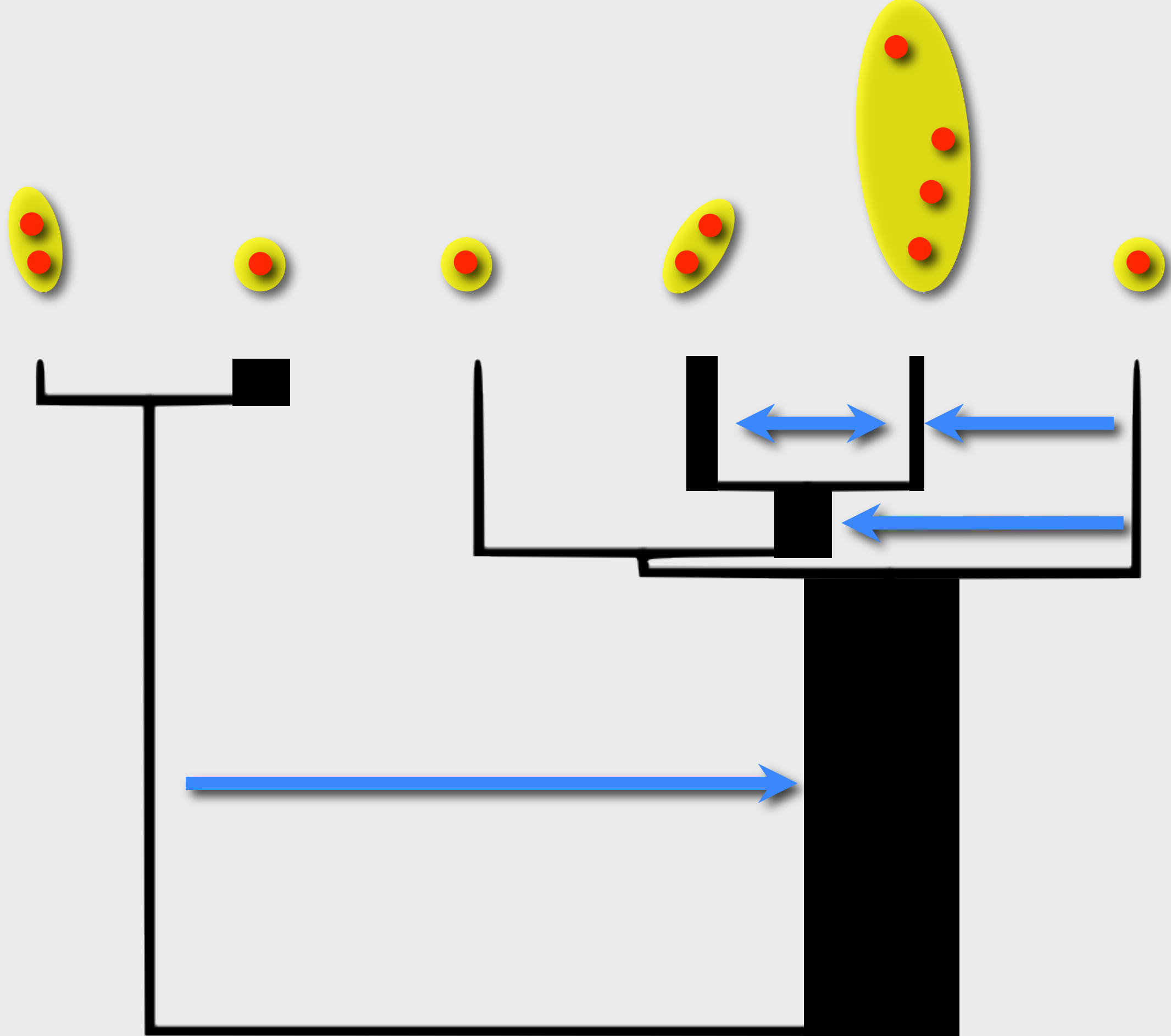


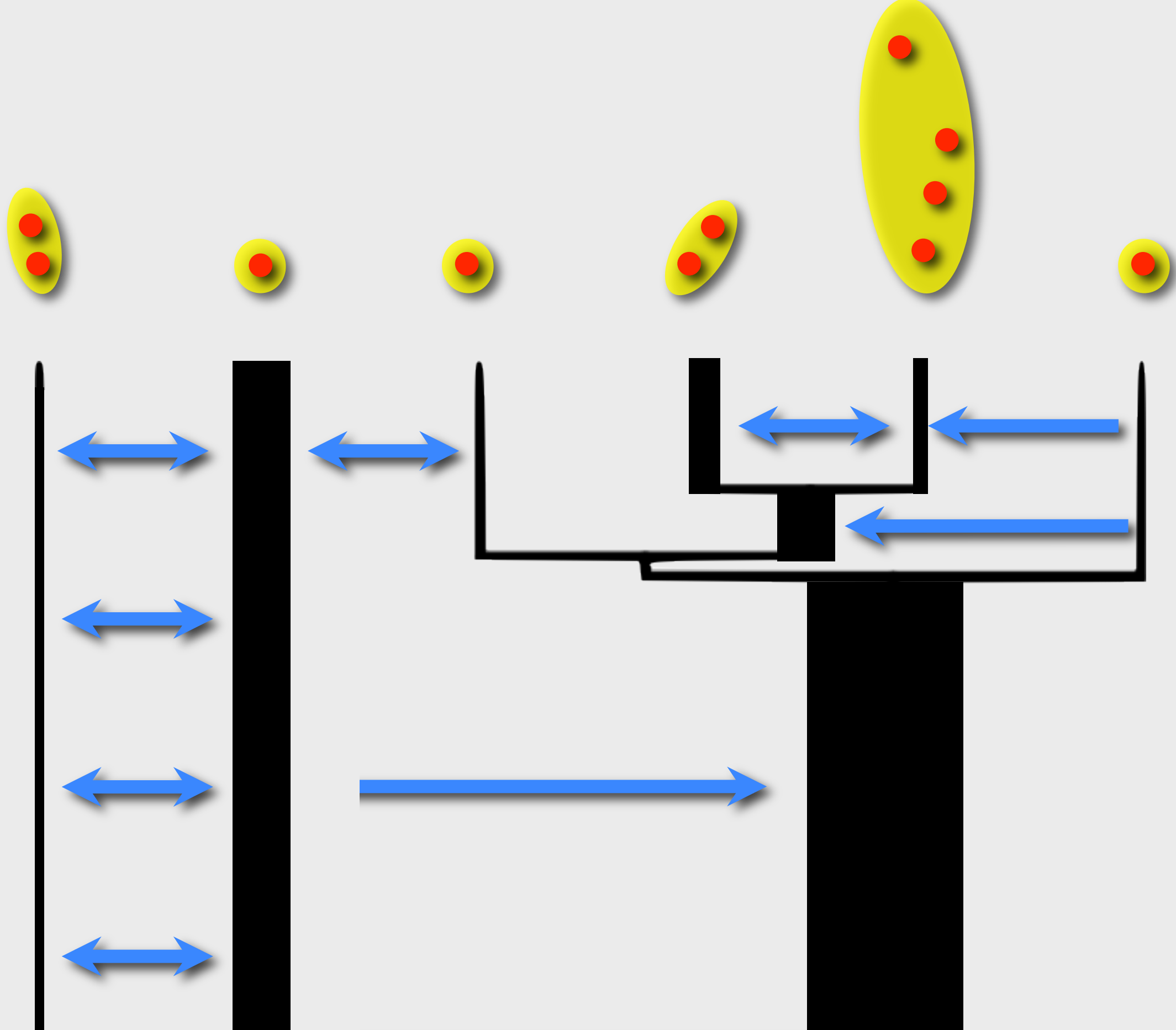


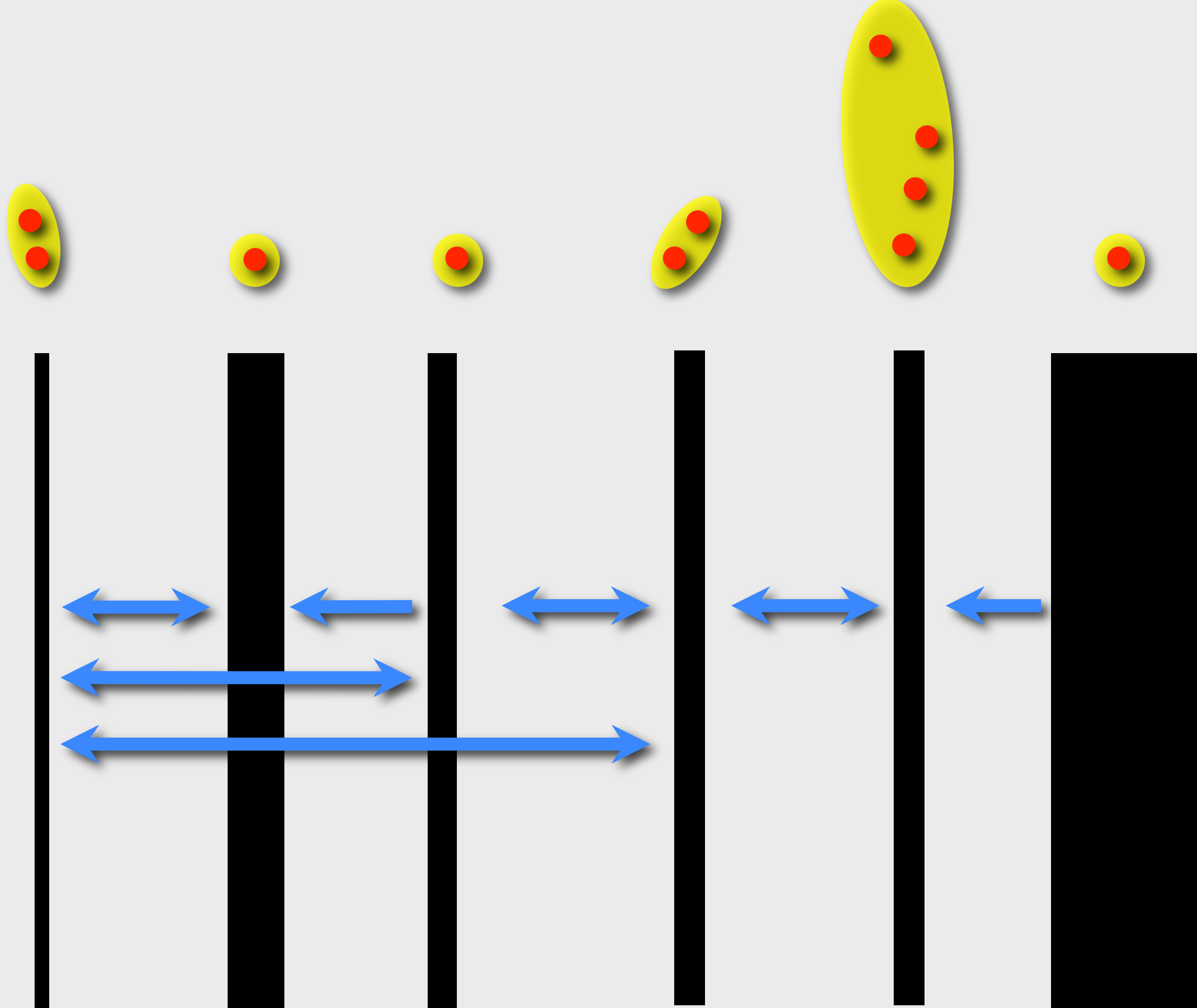












Brownie for species delim

Now starting the search proper.

A ">" before a score indicates that calculation of that score was aborted once the score for that move exceeded the best local score

Rep	Moves	#Spp	Type	Qual	CombScore	GTP	Struct	Local	Global	NTrees	Remaining
1	0	2		=G	770.620	52.000	1489.241	770.620	770.620	1	_aid_
1	1	2->3	i	*G	56.716	113.000	0.431	56.716	56.716	1	_aid_
1	2	3->2	d		> 56.716	0.000	> 113.432	56.716	56.716	1	_ai__
1	3	3->3	a		> 56.831	> 113.432	> 0.230	56.716	56.716	1	_ai__
1	4	3->3	a		> 56.831	> 113.432	> 0.230	56.716	56.716	1	_ai__
1	5	3->3	a		> 56.918	> 113.432	> 0.403	56.716	56.716	1	_ai__
1	6	3->4	i		> 56.831	> 113.432	> 0.230	56.716	56.716	1	_ai__
2	169	3->3	a		> 33.547	> 57.714	> 9.379	28.857	28.857	1	_a___
2	170	3->3	a		> 34.075	> 57.714	> 10.436	28.857	28.857	1	_a___
2	171	3->3	a		> 34.075	> 57.714	> 10.436	28.857	28.857	1	_a___
2	172	3->3	a		35.954	54.000	17.907	28.857	28.857	1	_a___
2	173	3->3	a		34.240	49.000	19.481	28.857	28.857	1	_____

Best trees overall

```

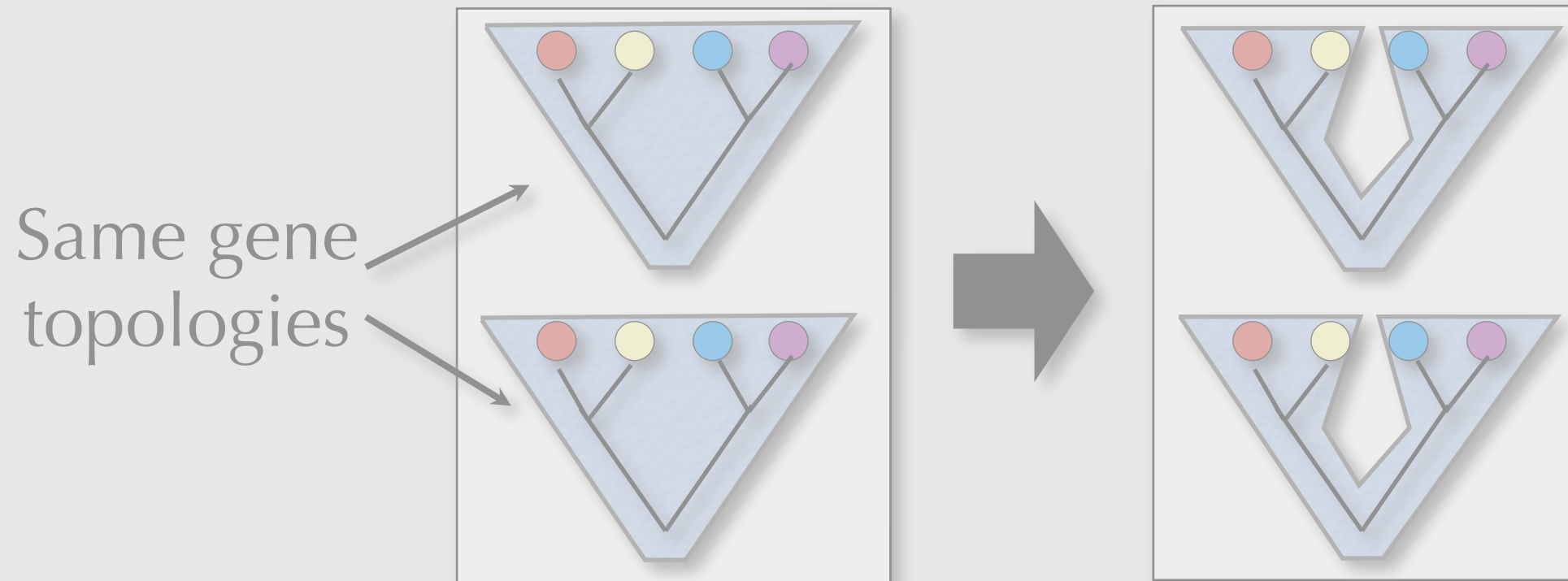
+-- (DpersimilisMather27,DpersimilisMather37,DpersimilisMather41,DpersimilisMatherB,DpersimilisMSH1,DpersimilisMSH3,DpersimilisMSH42,DpersimilisMSH7,DpersimilisSalem)
|
|+-- (DbogotanaPotosi2,DbogotanaSusa3,DbogotanaSutatausa1,DbogotanaSutatausa3,DbogotanaToro4,DbogotanaToro6,DbogotanaToro7)
+|
+-- (DpseudoobscuraAF2,DpseudoobscuraAFC12,DpseudoobscuraAFC3,DpseudoobscuraAFC7,DpseudoobscuraMather48,DpseudoobscuraMSH32)

```

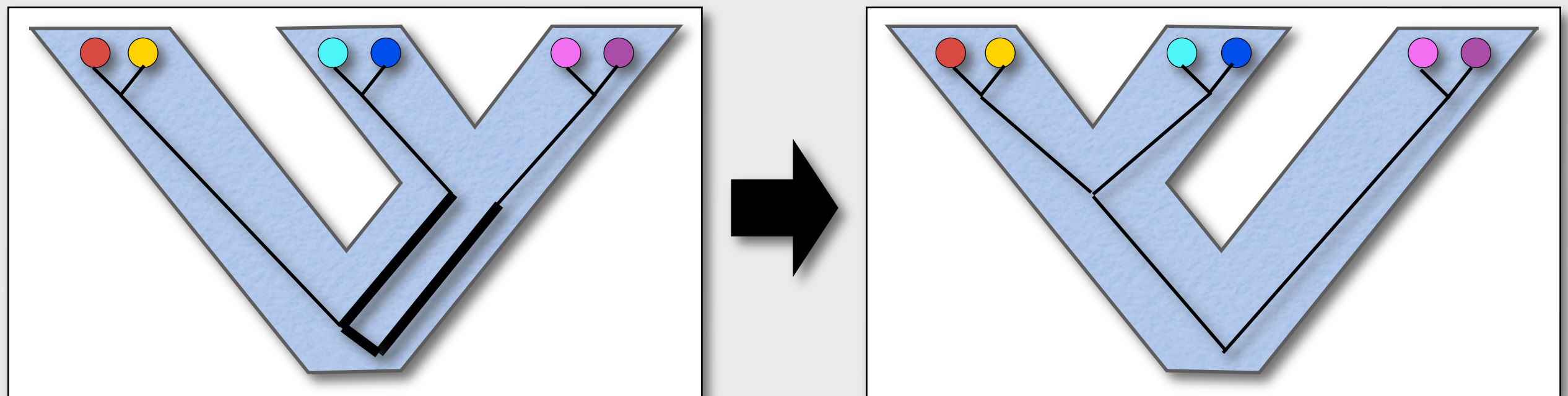
O'Meara 2010

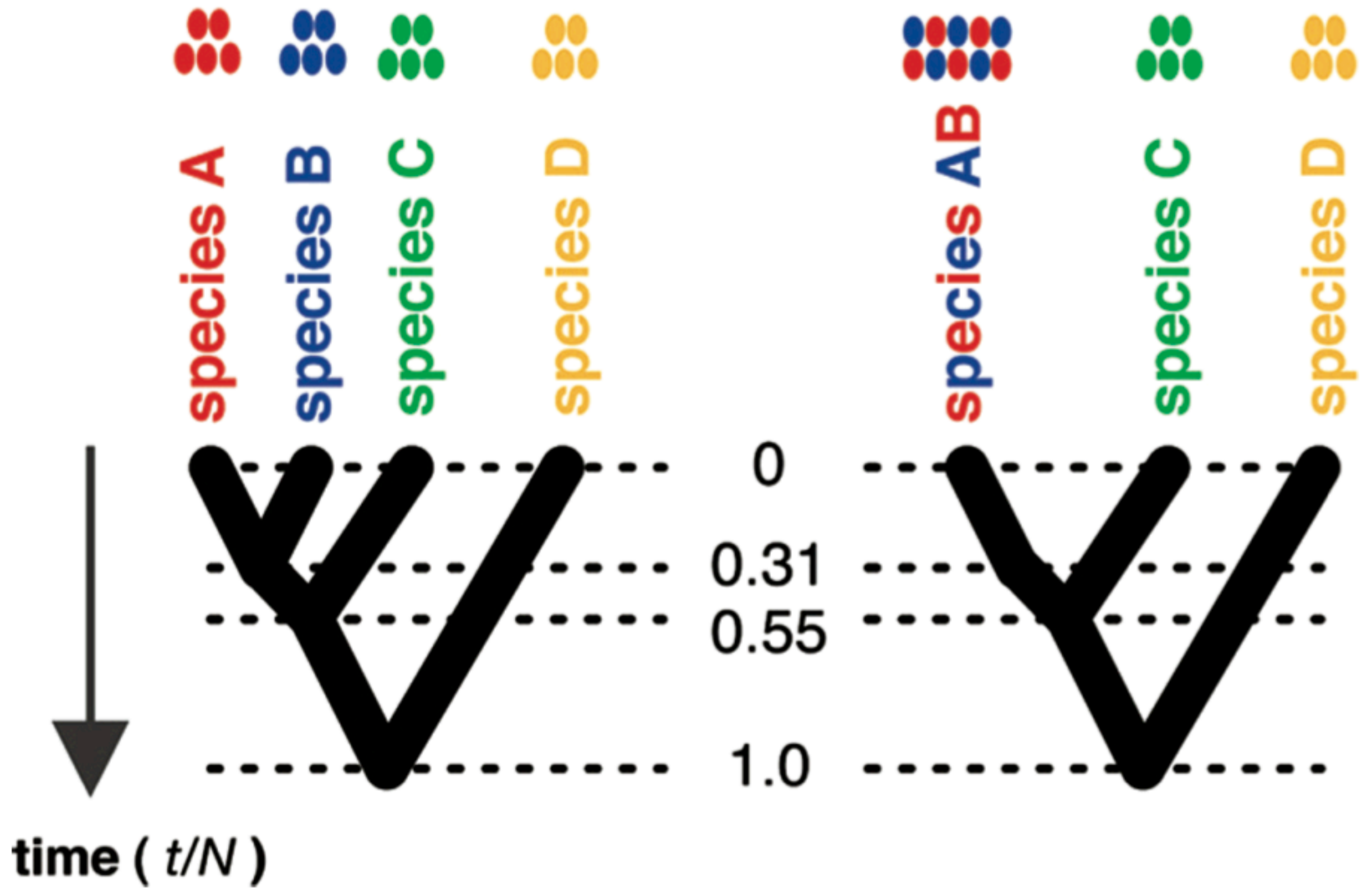
Infers species assignments and species tree simultaneously, using two approaches:

Minimize excess intraspecific structure



Minimize gene tree-species tree conflict





Brownie for species delim

Now starting the search proper.

A ">" before a score indicates that calculation of that score was aborted once the score for that move exceeded the best local score

Rep	Moves	#Spp	Type	Qual	CombScore	GTP	Struct	Local	Global	NTrees	Remaining
1	0	2		=G	770.620	52.000	1489.241	770.620	770.620	1	_aid_
1	1	2->3	i	*G	56.716	113.000	0.431	56.716	56.716	1	_aid_
1	2	3->2	d		> 56.716	0.000	> 113.432	56.716	56.716	1	_ai__
1	3	3->3	a		> 56.831	> 113.432	> 0.230	56.716	56.716	1	_ai__
1	4	3->3	a		> 56.831	> 113.432	> 0.230	56.716	56.716	1	_ai__
1	5	3->3	a		> 56.918	> 113.432	> 0.403	56.716	56.716	1	_ai__
1	6	3->4	i		> 56.831	> 113.432	> 0.230	56.716	56.716	1	_ai__
2	169	3->3	a		> 33.547	> 57.714	> 9.379	28.857	28.857	1	_a___
2	170	3->3	a		> 34.075	> 57.714	> 10.436	28.857	28.857	1	_a___
2	171	3->3	a		> 34.075	> 57.714	> 10.436	28.857	28.857	1	_a___
2	172	3->3	a		35.954	54.000	17.907	28.857	28.857	1	_a___
2	173	3->3	a		34.240	49.000	19.481	28.857	28.857	1	-----

Best trees overall

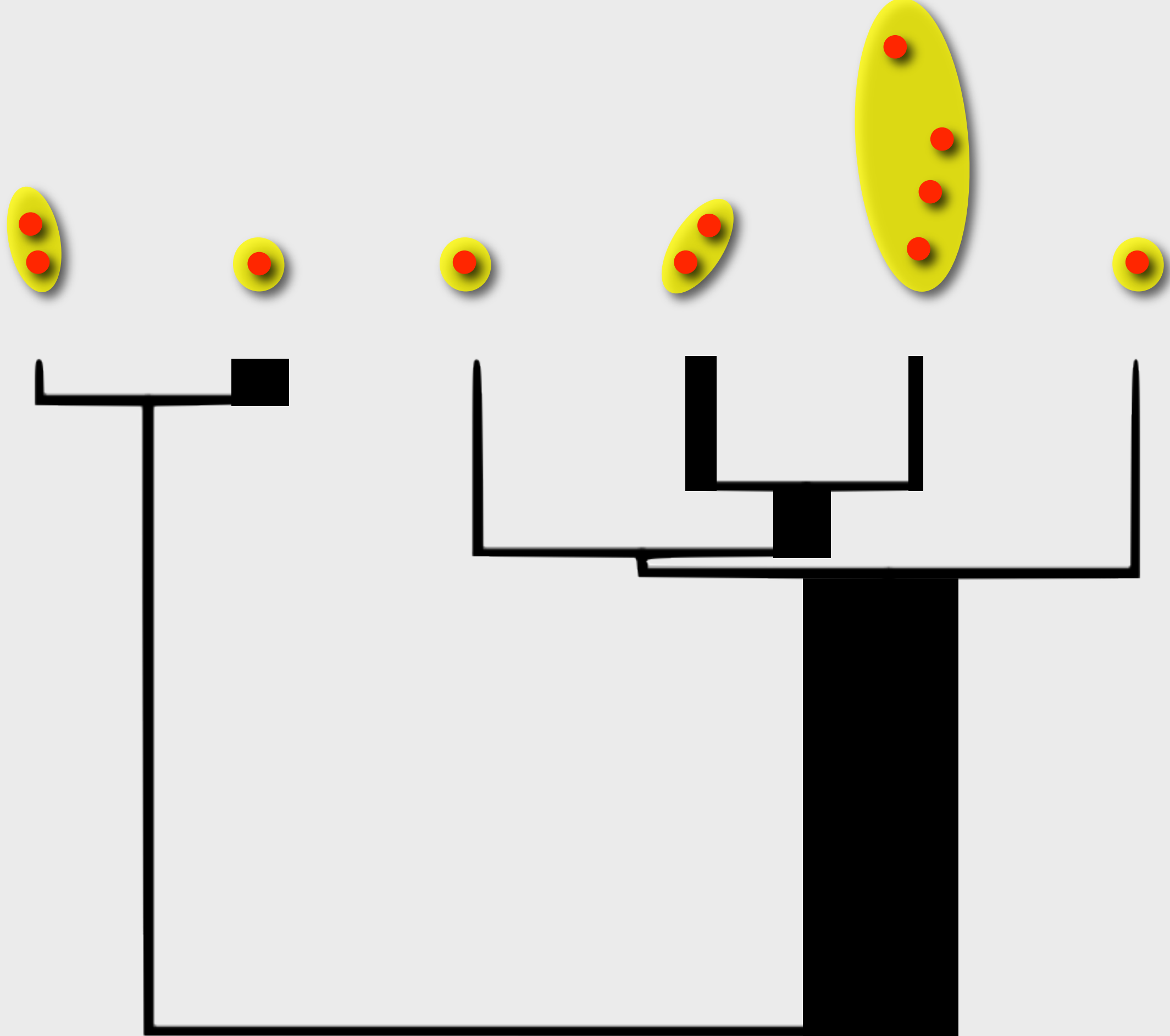
```

+-- (DpersimilisMather27,DpersimilisMather37,DpersimilisMather41,DpersimilisMatherB,DpersimilisMSH1,DpersimilisMSH3,DpersimilisMSH42,DpersimilisMSH7,DpersimilisSalem)
|
|+-- (DbogotanaPotosi2,DbogotanaSusa3,DbogotanaSutatausa1,DbogotanaSutatausa3,DbogotanaToro4,DbogotanaToro6,DbogotanaToro7)
+|
+-- (DpseudoobscuraAF2,DpseudoobscuraAFC12,DpseudoobscuraAFC3,DpseudoobscuraAFC7,DpseudoobscuraMather48,DpseudoobscuraMSH32)

```

O'Meara 2010

Infers species assignments and species tree simultaneously. Returns a set of equally good trees.



Brownie

Single point estimate, no confidence

Tries a wide range of models

No migration

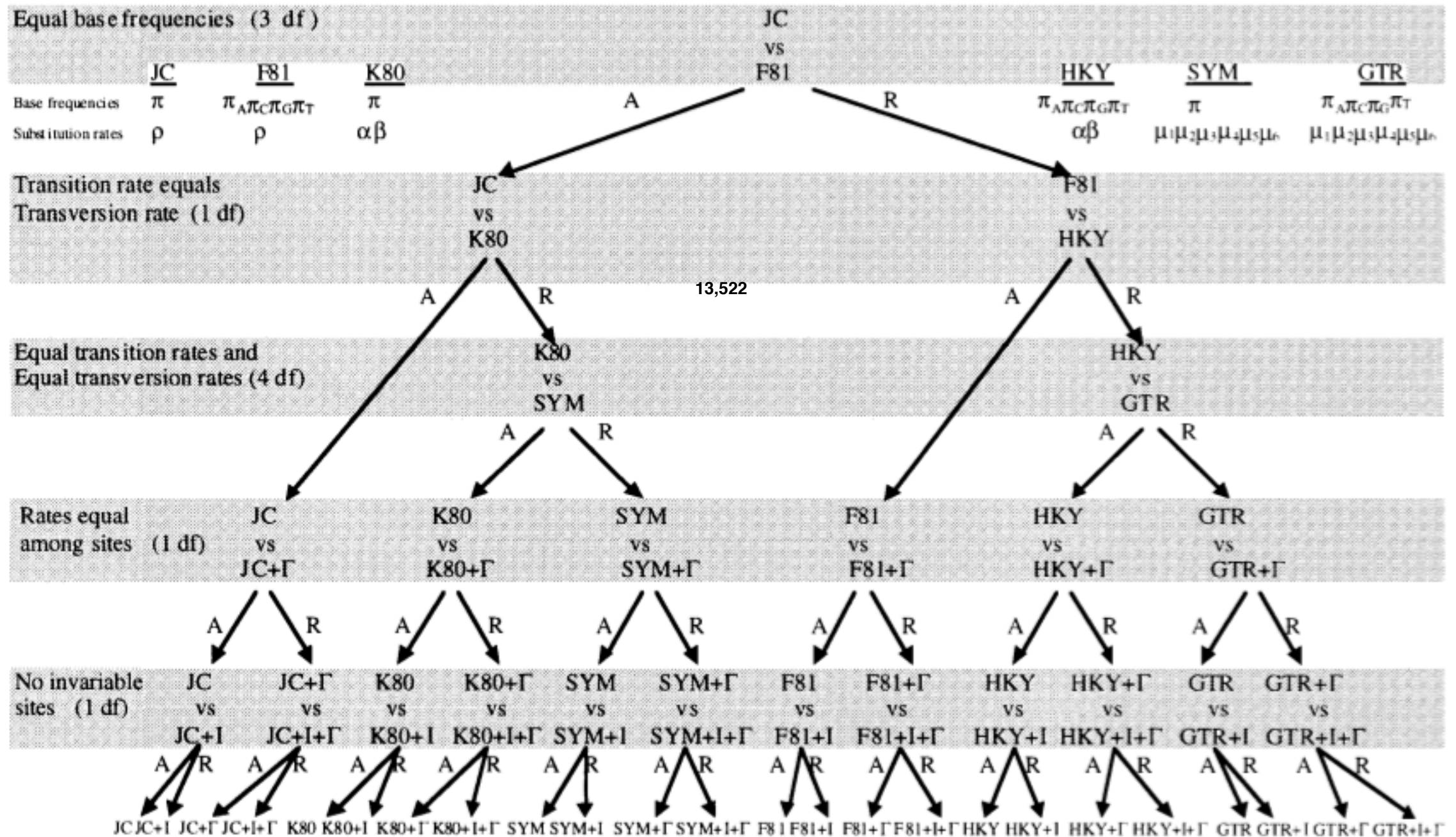
Many other approaches

Migration only (e.g. Migrate) or population
coalescence only (e.g. BPP, gmyc)

Often specify two or a few models (backbone tree,
specific ABC model, etc.)



13,522 citations of first ModelTest paper (Posada and Crandall 1998) alone



Question	Specified hypothesis method	Try all method
DNA model	PAML	ModelTest
AA model	PAML	ProtTest
Continuous trait model	OUCH, Brownie (traits), OUWie	SURFACE, auteur
Diversification model	laser	Medusa
Phylogeographic	various ABC approaches	phrapl

phrapl **phylogeographic inference using approximated likelihoods**

Input:

gene topologies

assignments of samples to populations

Jargon

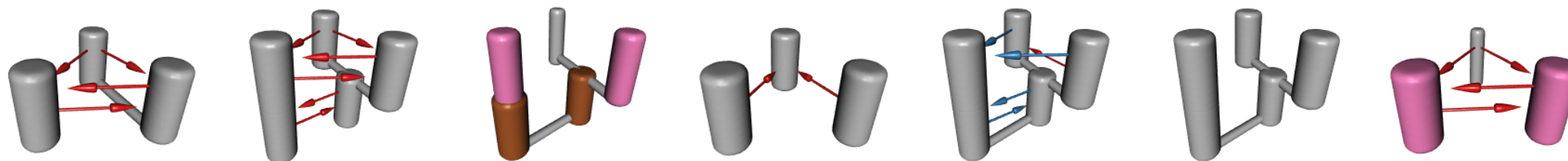
N = number of populations

K = number of free parameters (as in AIC)

1

Generate all possible models

Given N populations, $\leq K$ free parameters



Filter

(only tree models, no more than two migration rates, etc.)

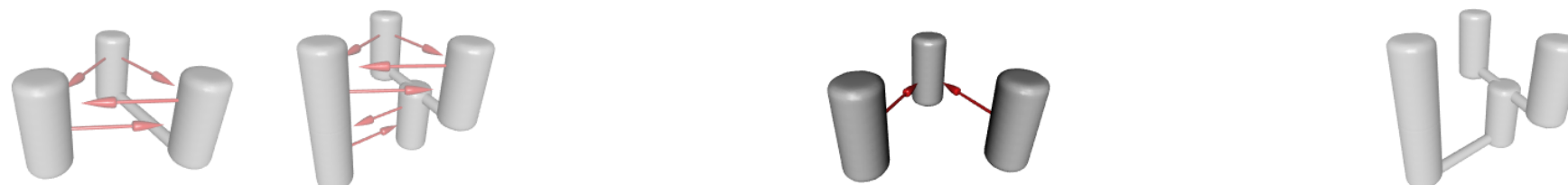
2



Analyze

(find best, find AIC for all)

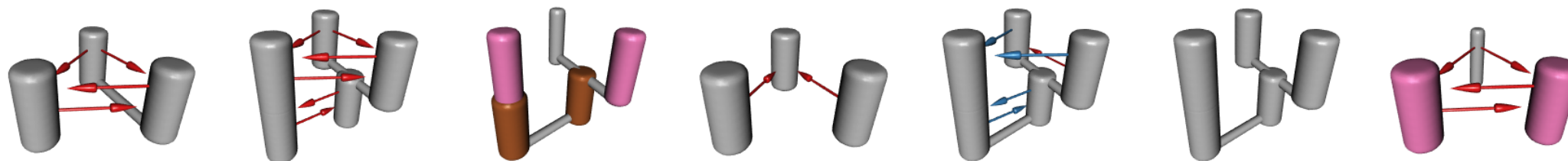
3



1

Generate all possible models

Given N populations, $\leq K$ free parameters



Filter

(only tree models, no more than two migration rates, etc.)

2



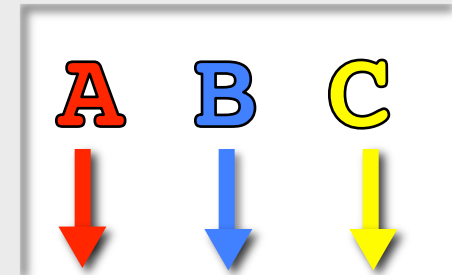
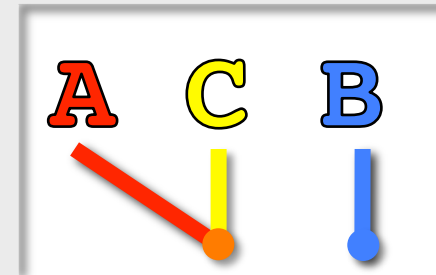
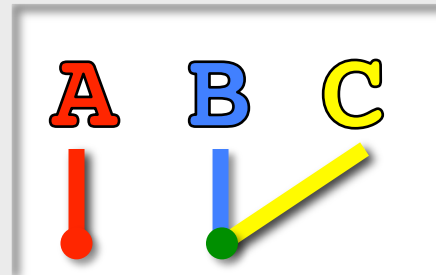
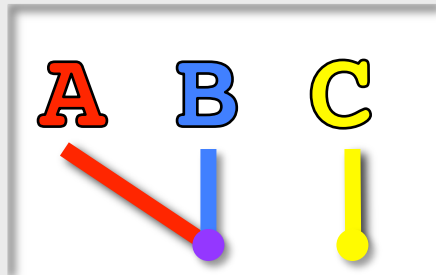
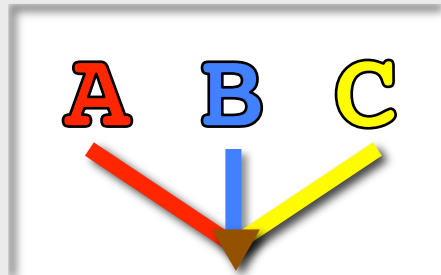
Analyze

(find best, find AIC for all)

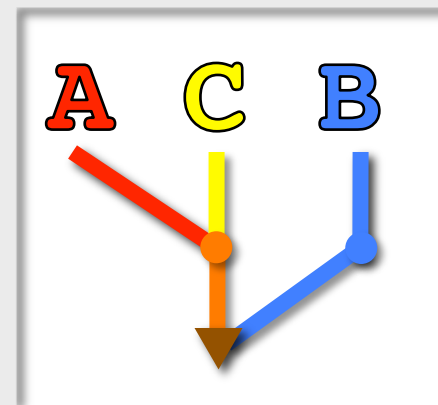
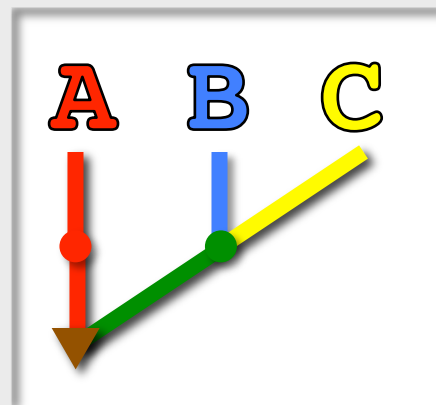
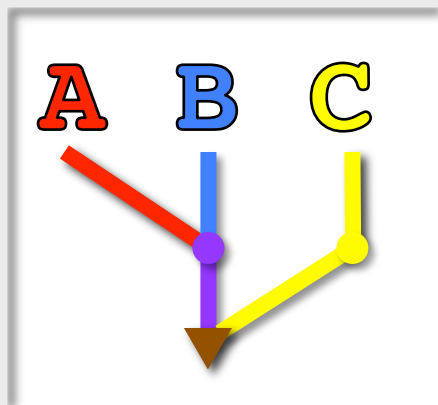
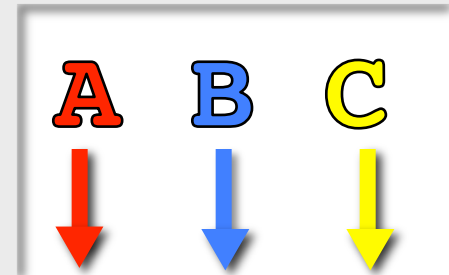
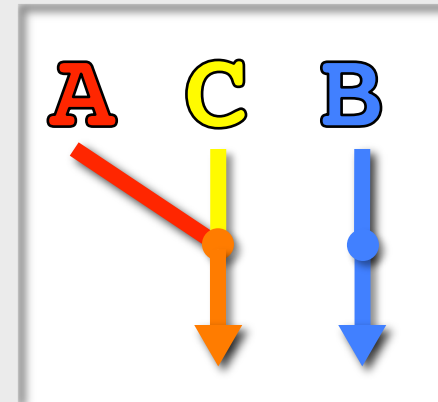
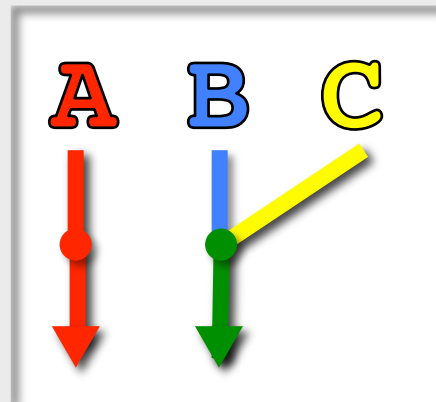
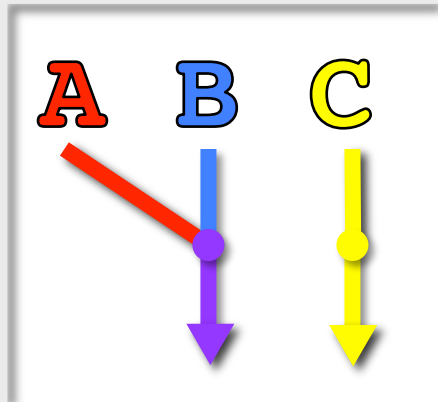
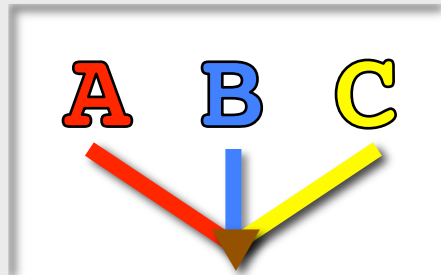
3



Coalescence of populations

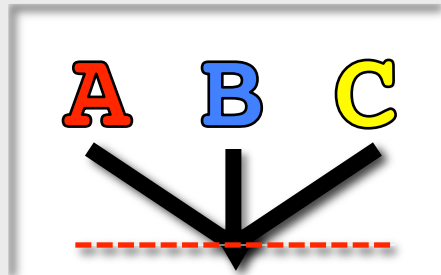


Coalescence of populations

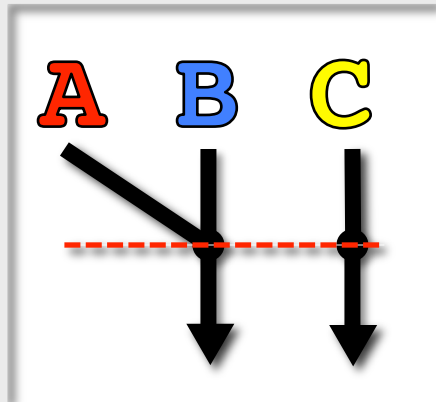


Coalescence of populations

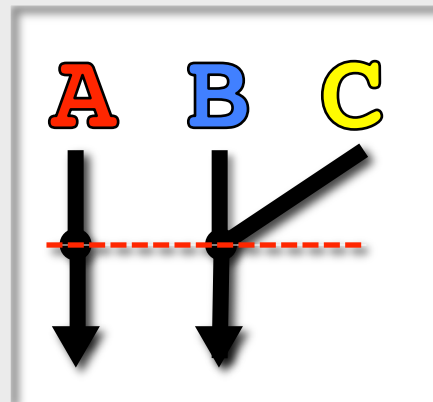
K=1



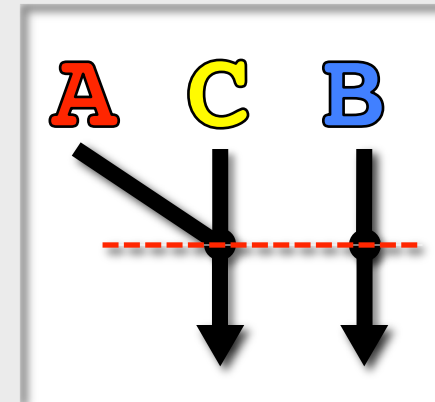
K=1



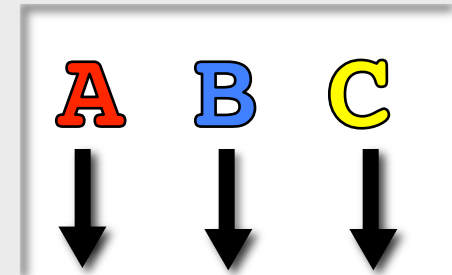
K=1



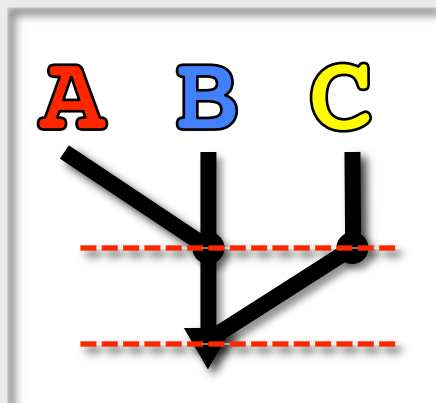
K=1



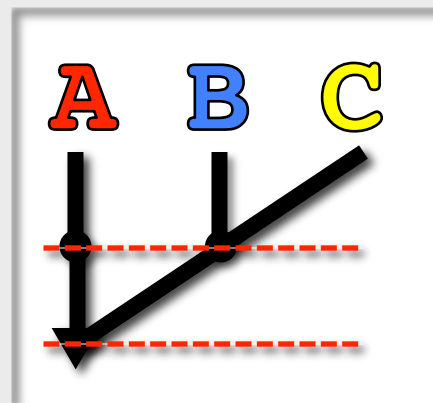
K=0



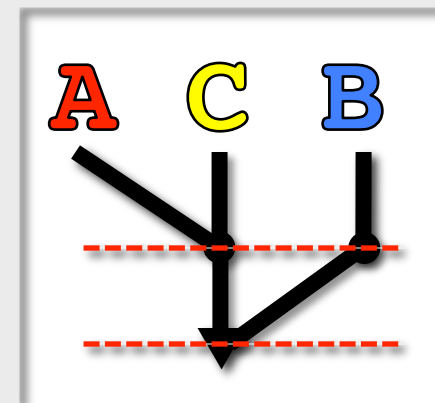
K=2



K=2



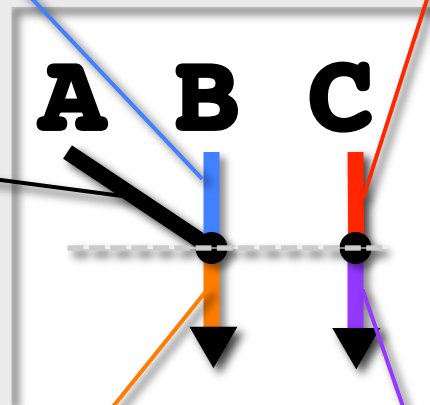
K=2



B: Same or
different as A

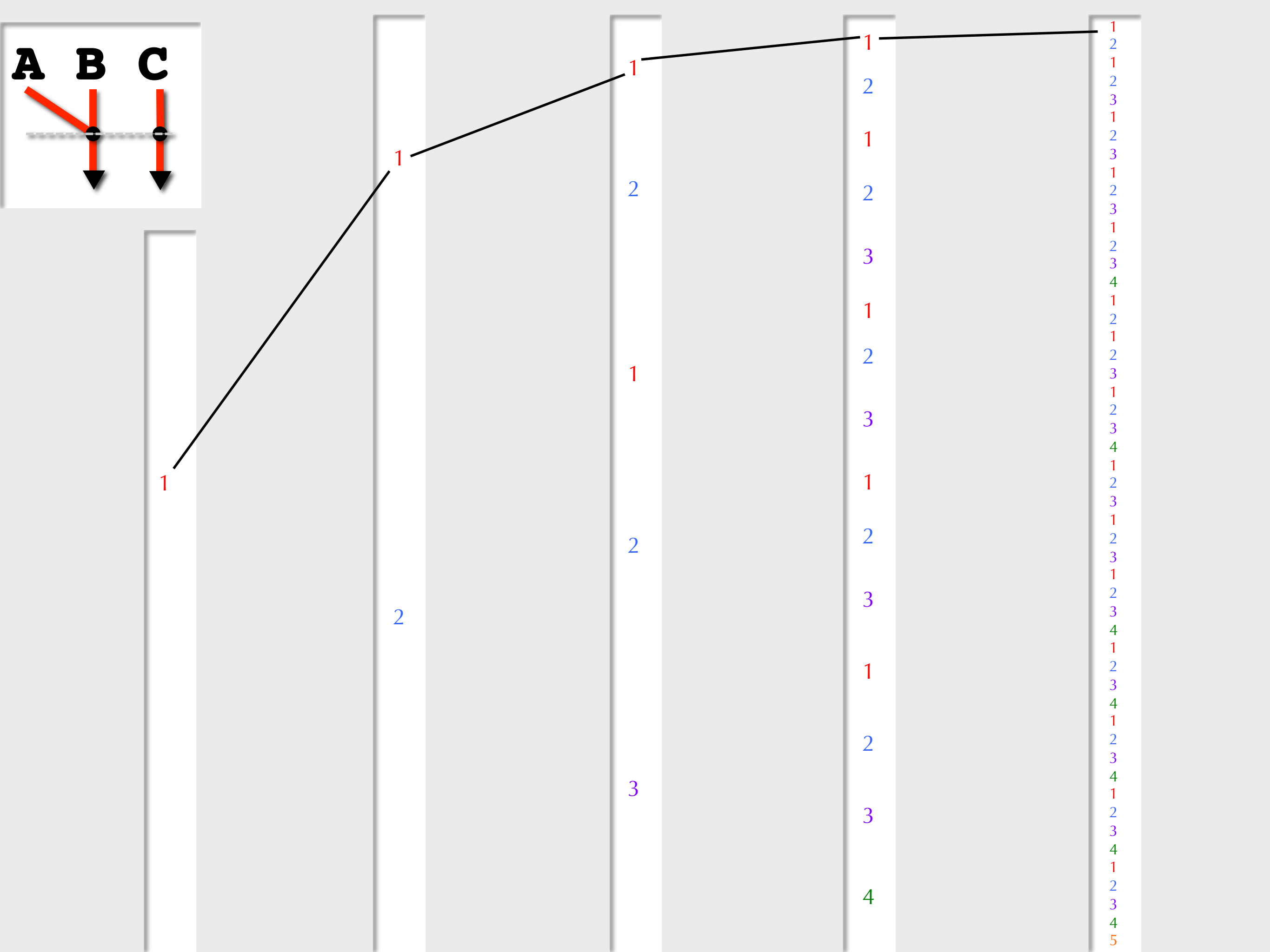
C1: Same or
different as A
and/or B and/or AB

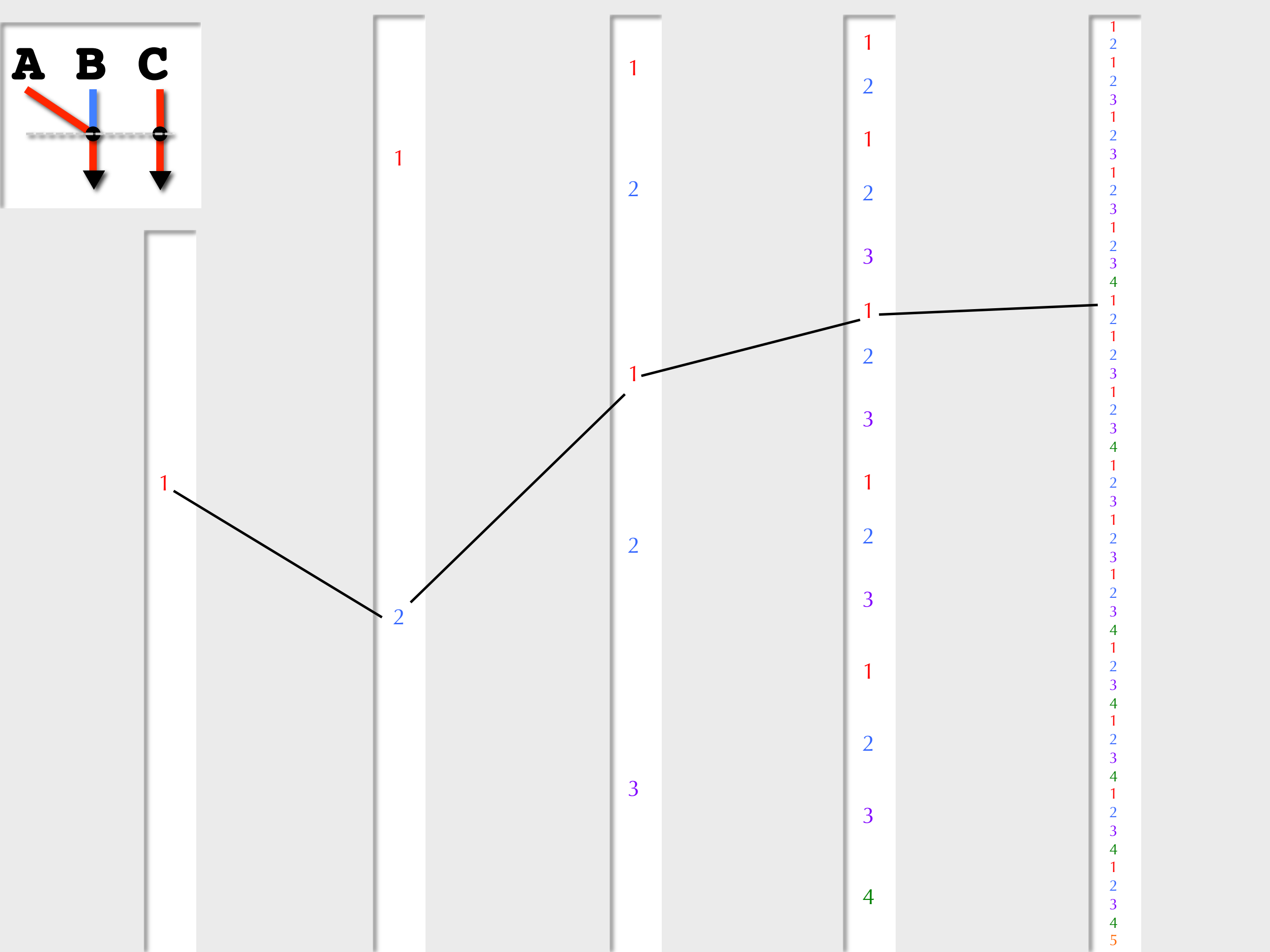
A: N_1

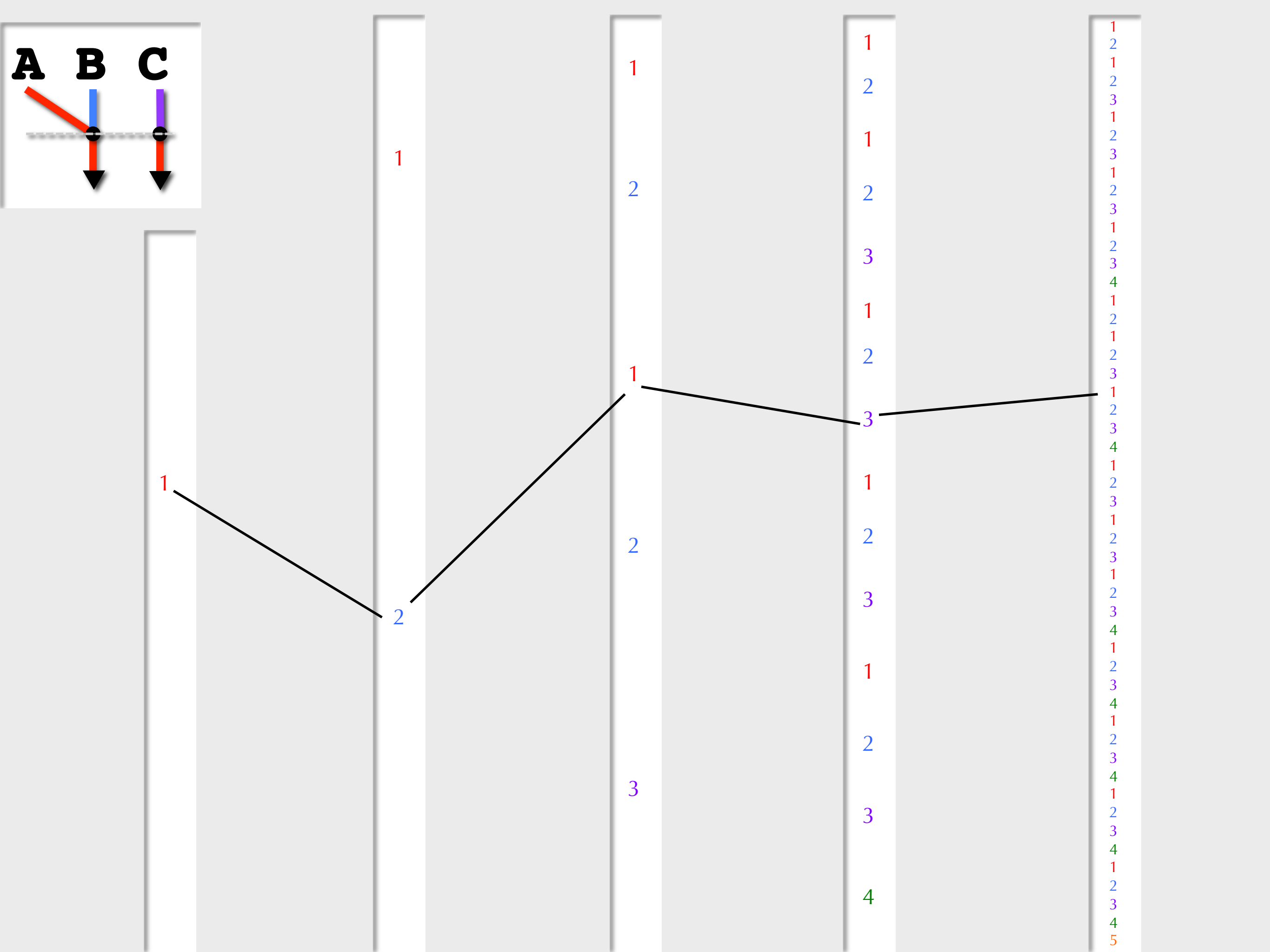


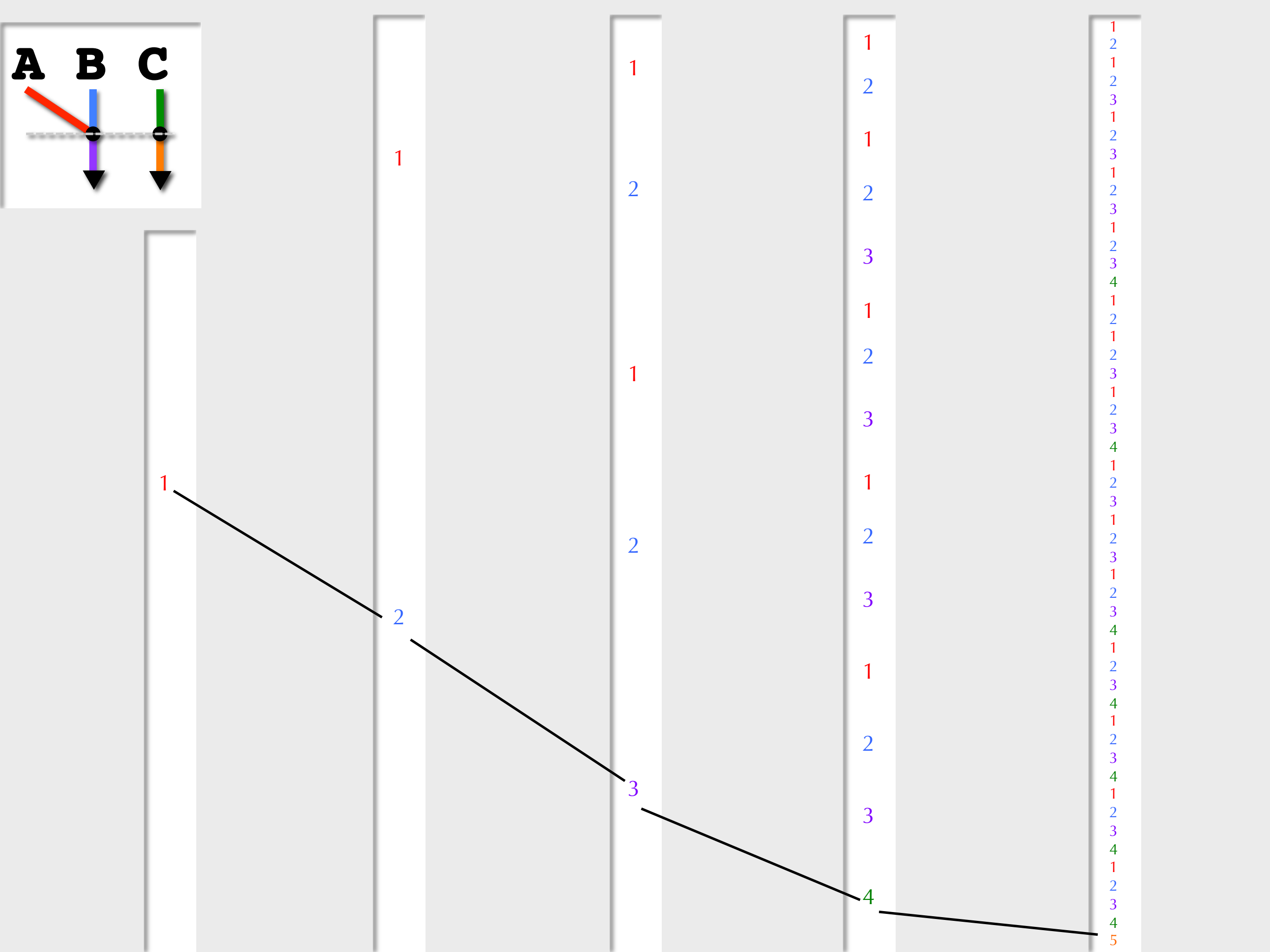
AB: Same or
different as A
and/or B

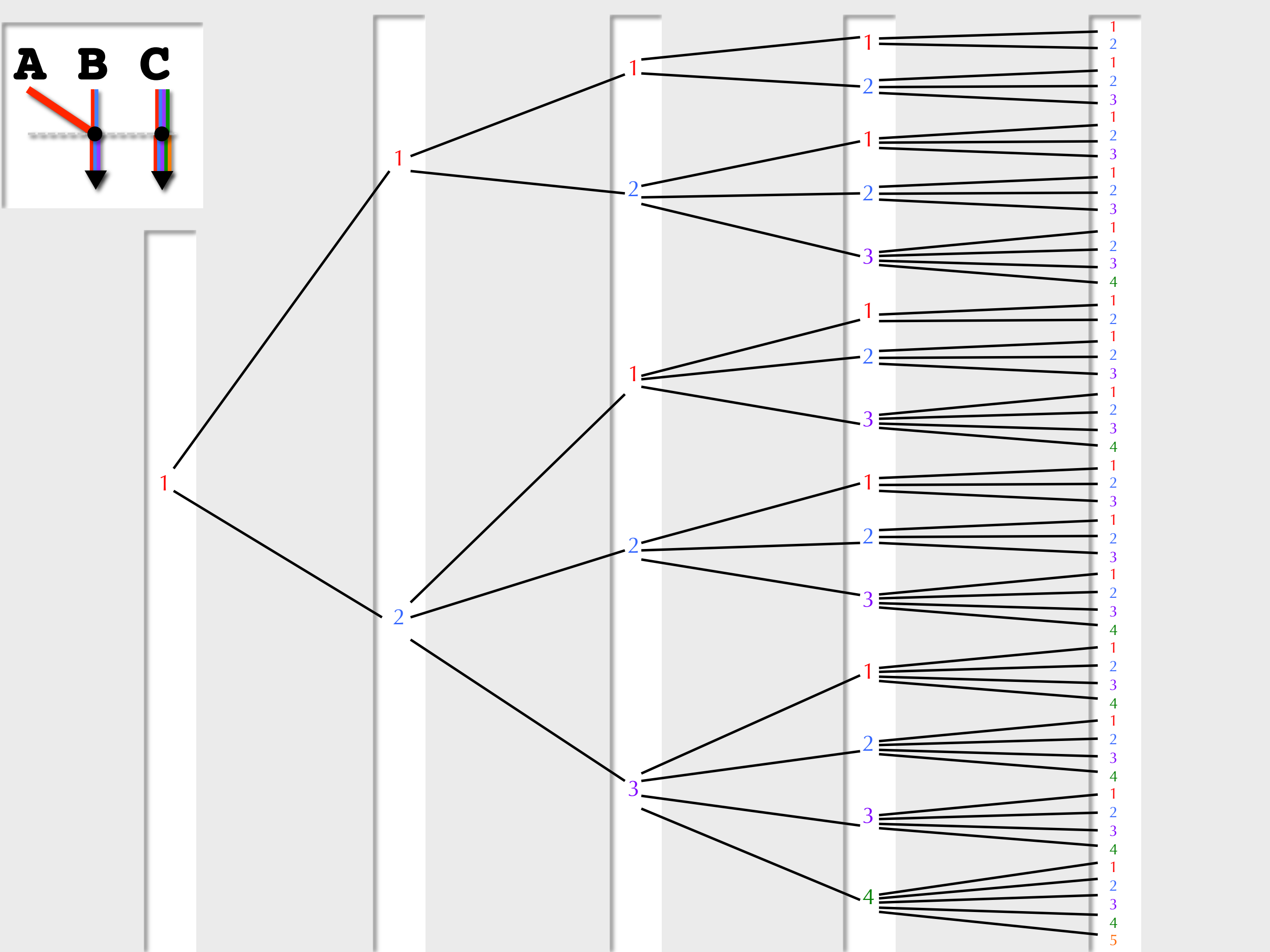
C2: Same or
different as A and/or
B and/or AB and/or
C1

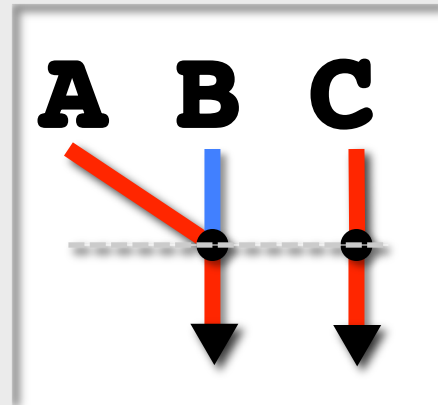




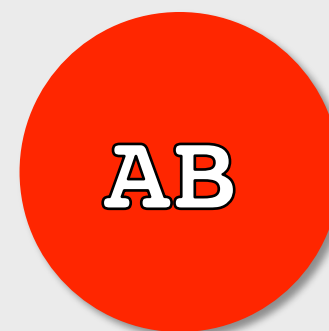
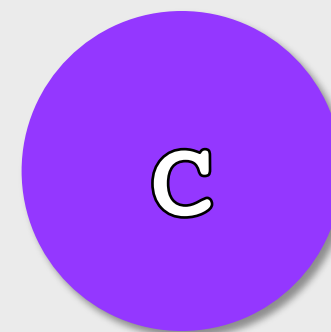
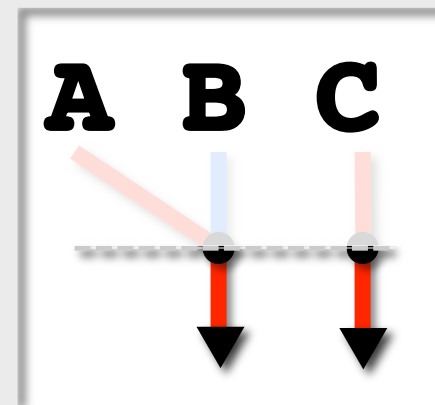
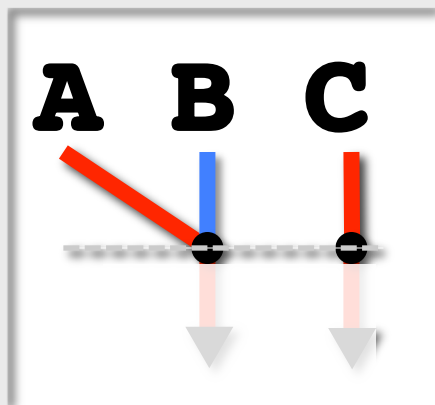


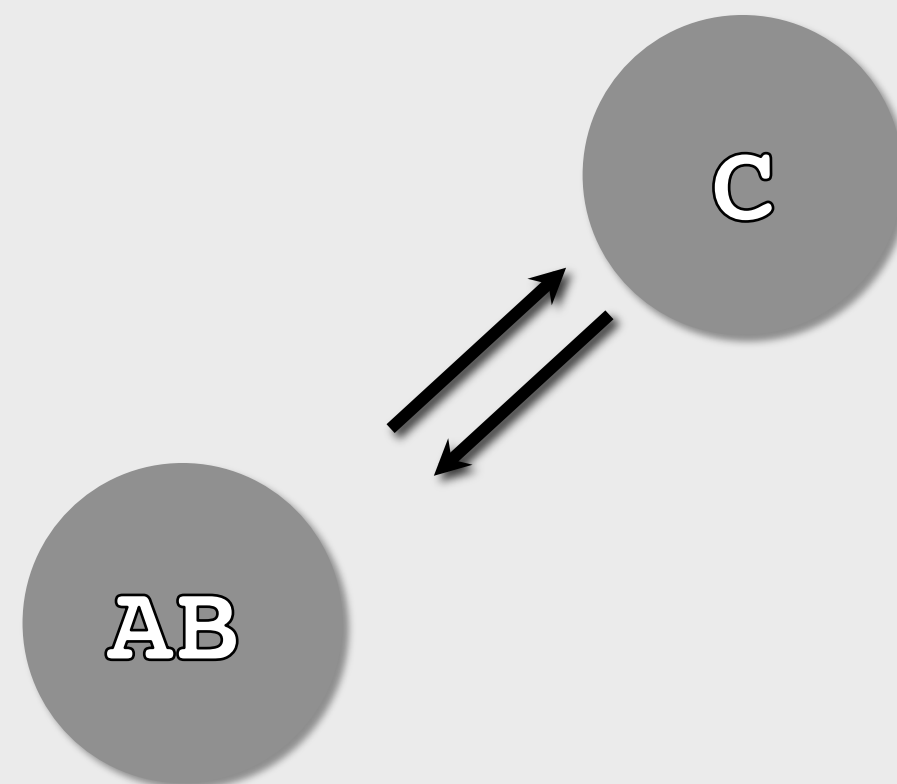
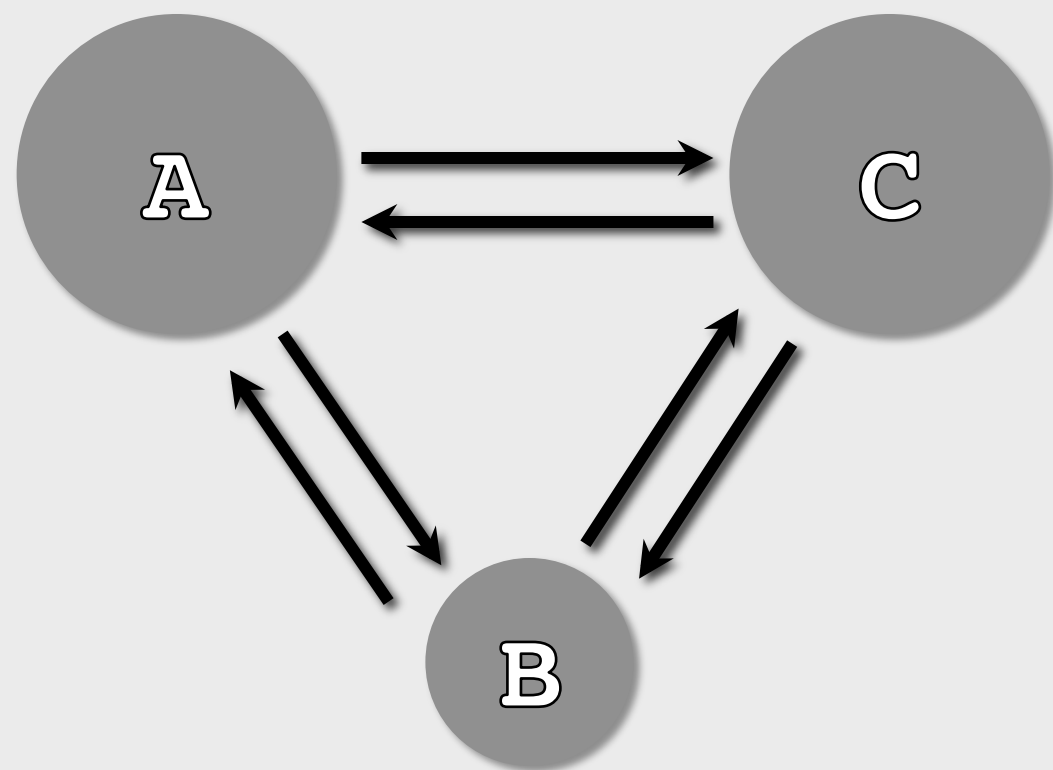
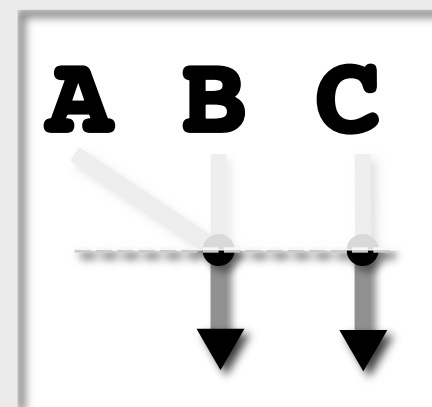
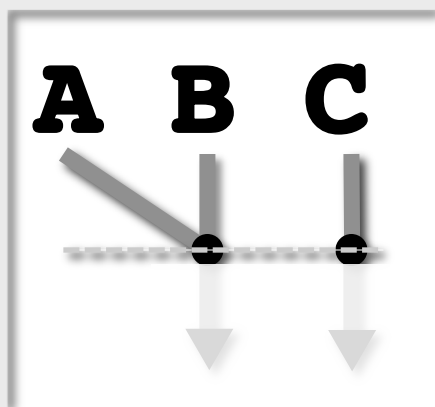


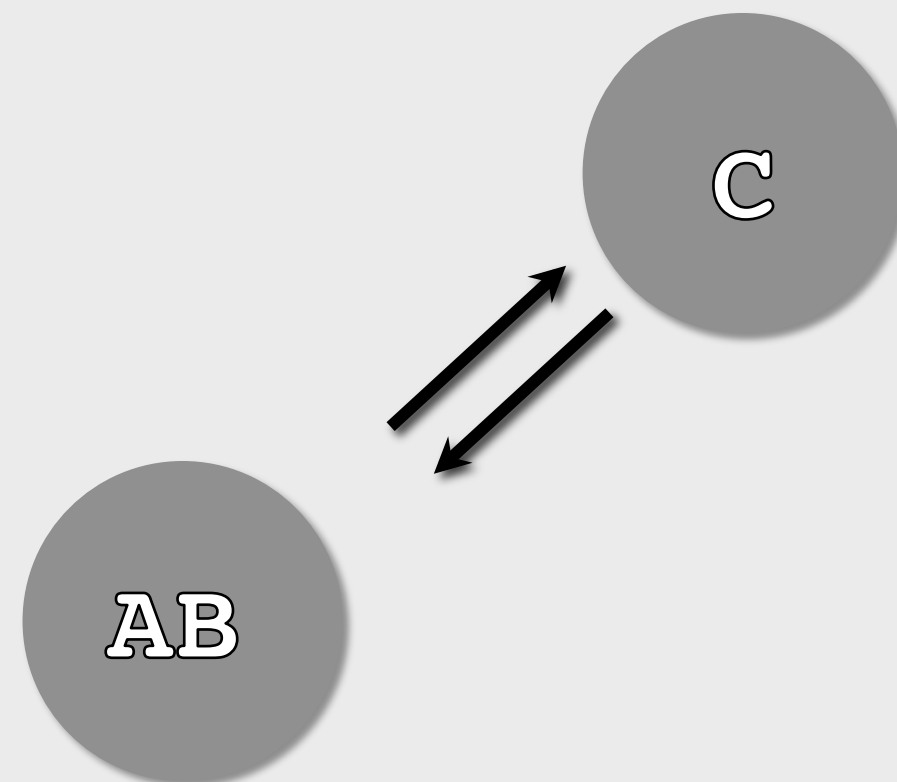
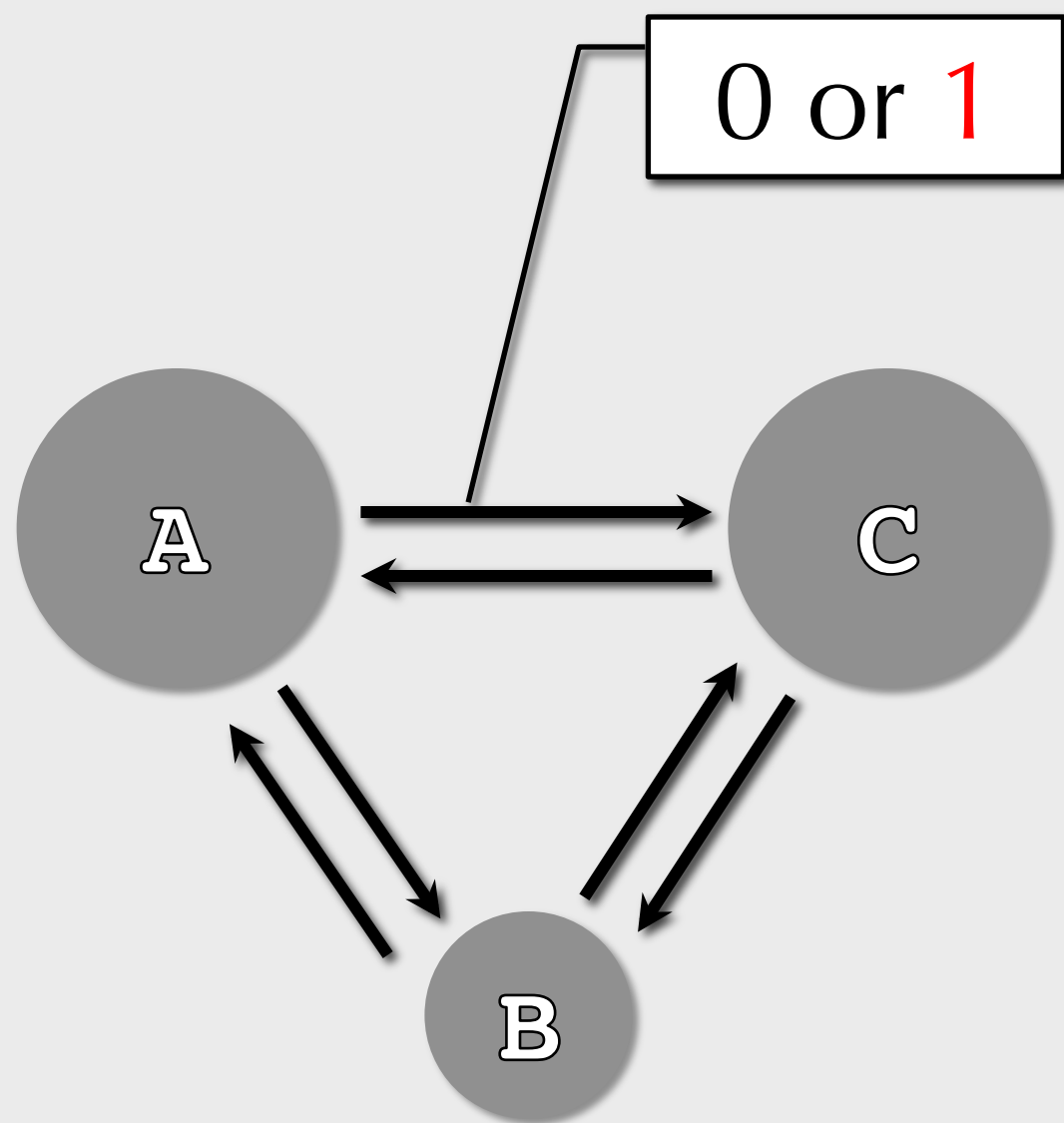
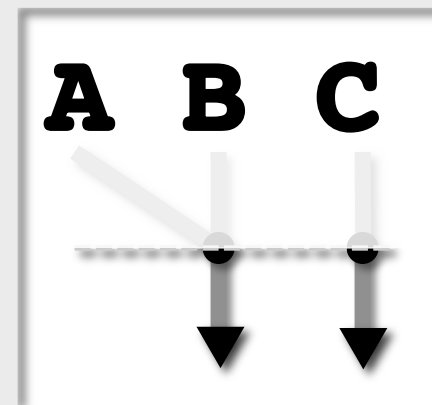
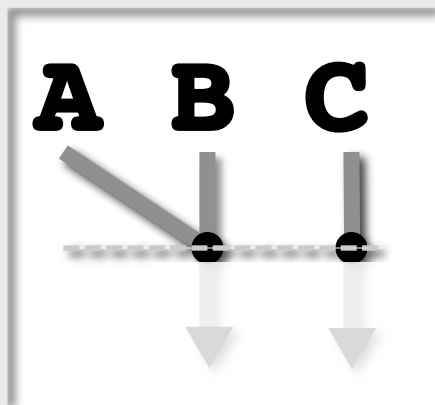


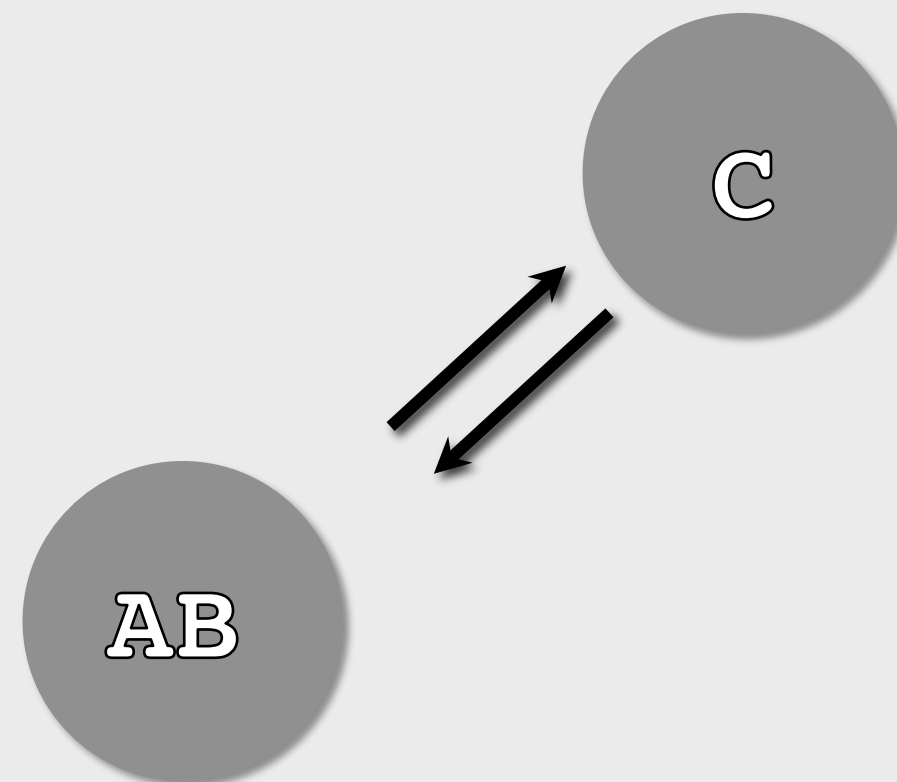
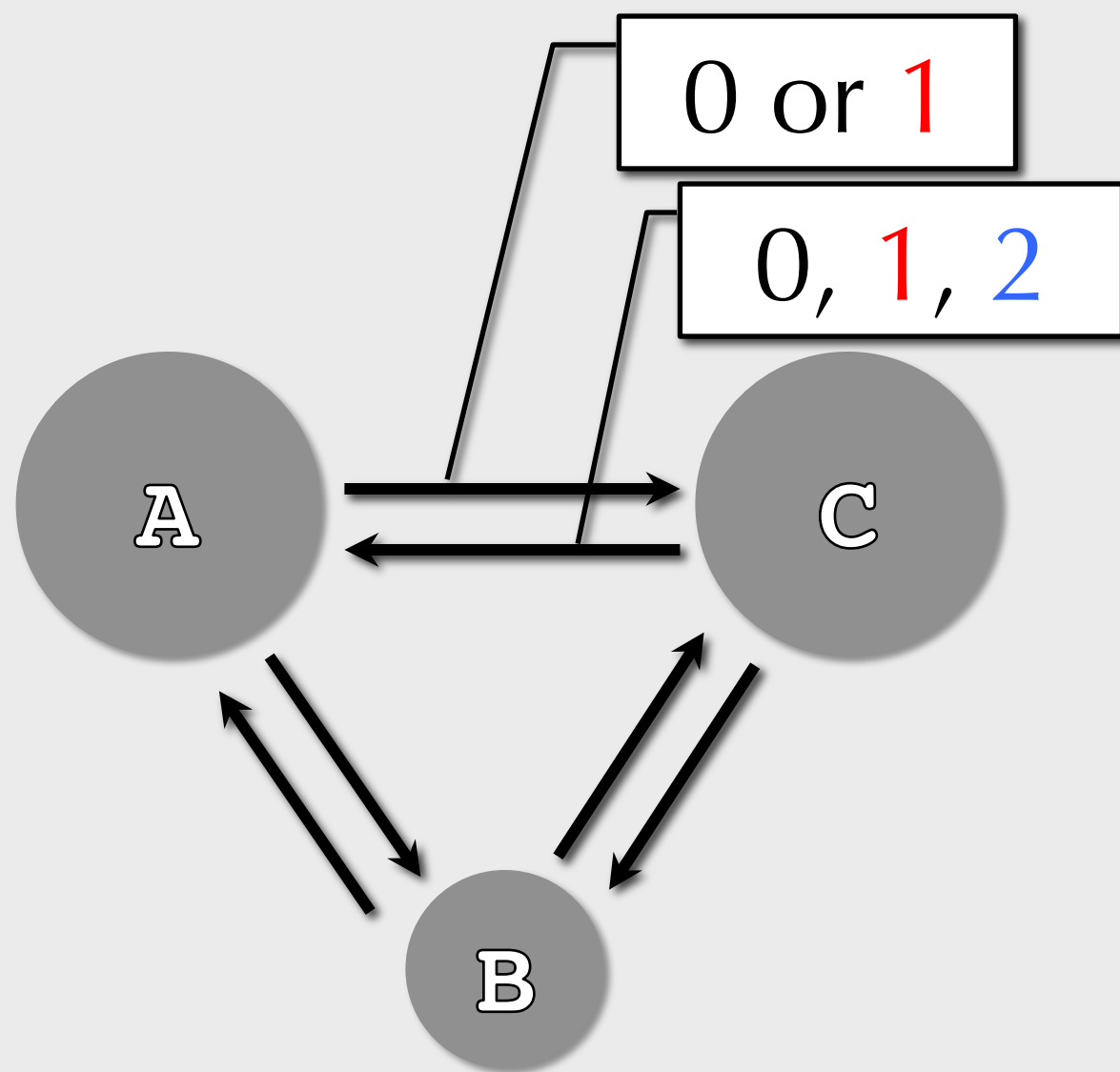
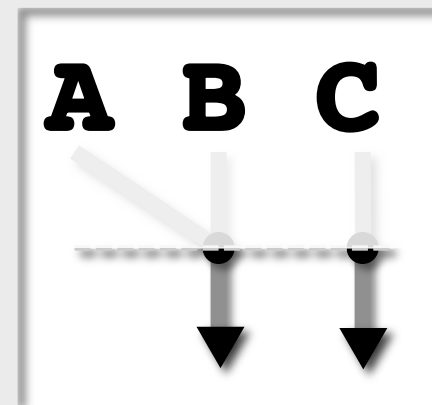
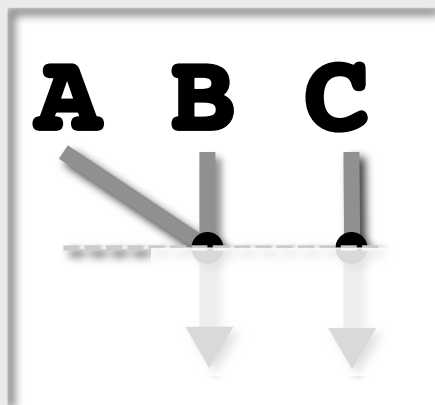


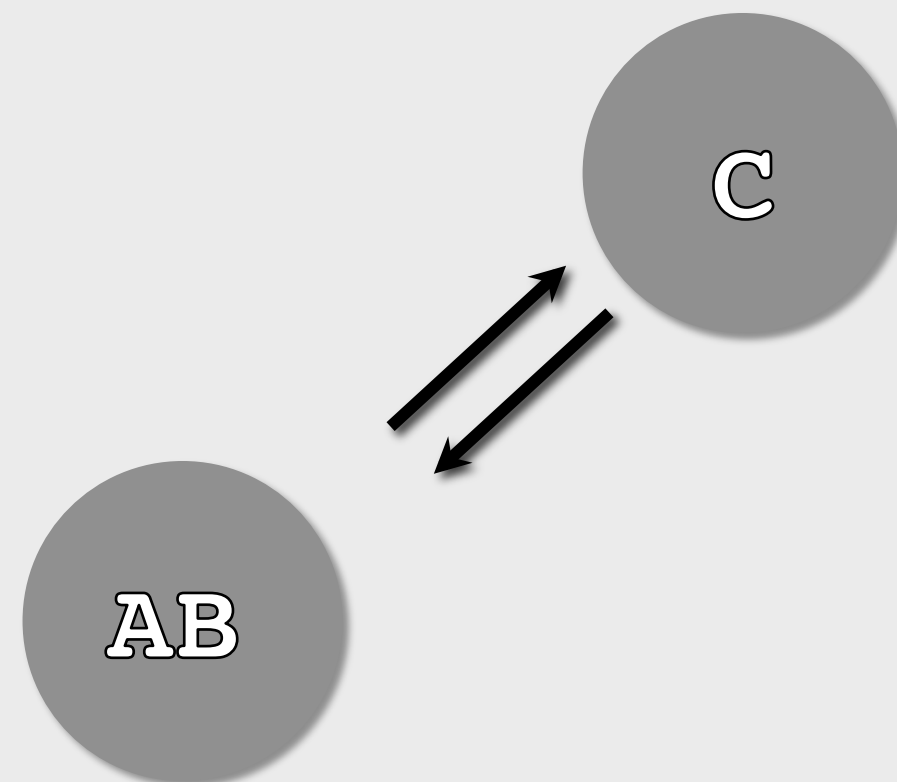
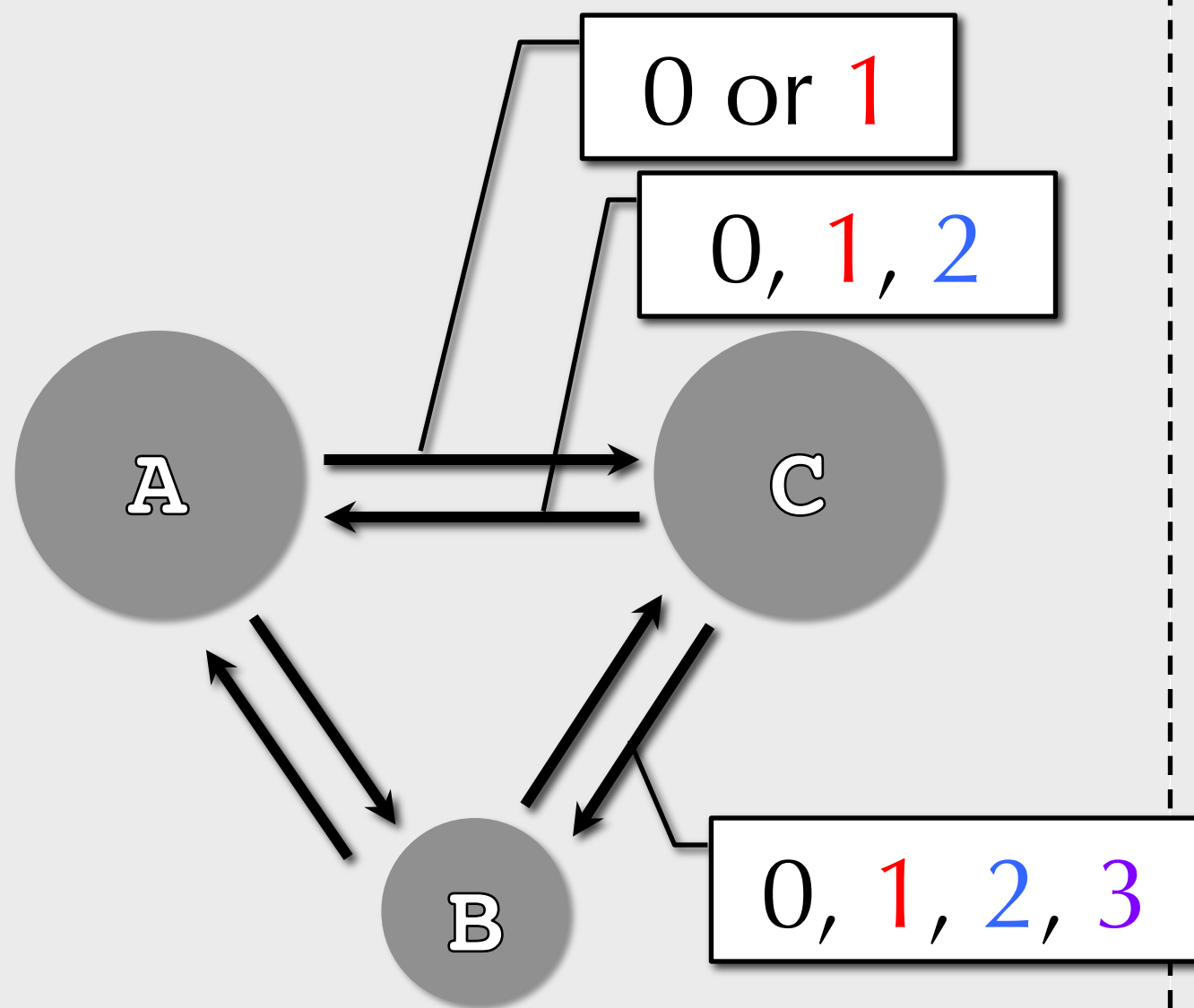
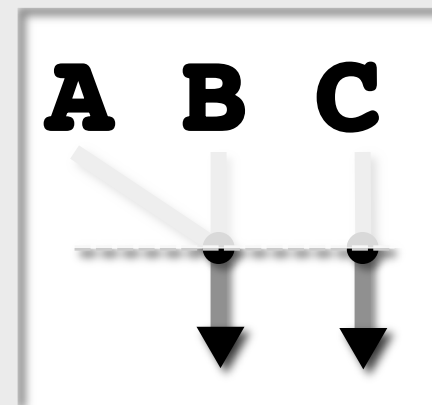
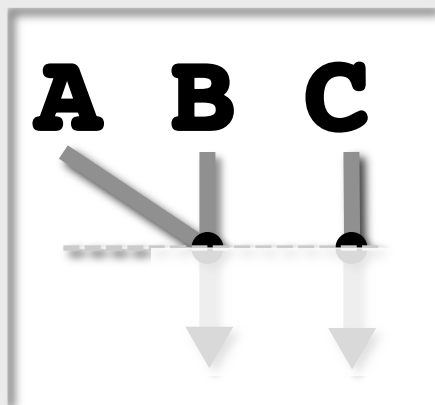
1 collapse parameter
2 population size parameters

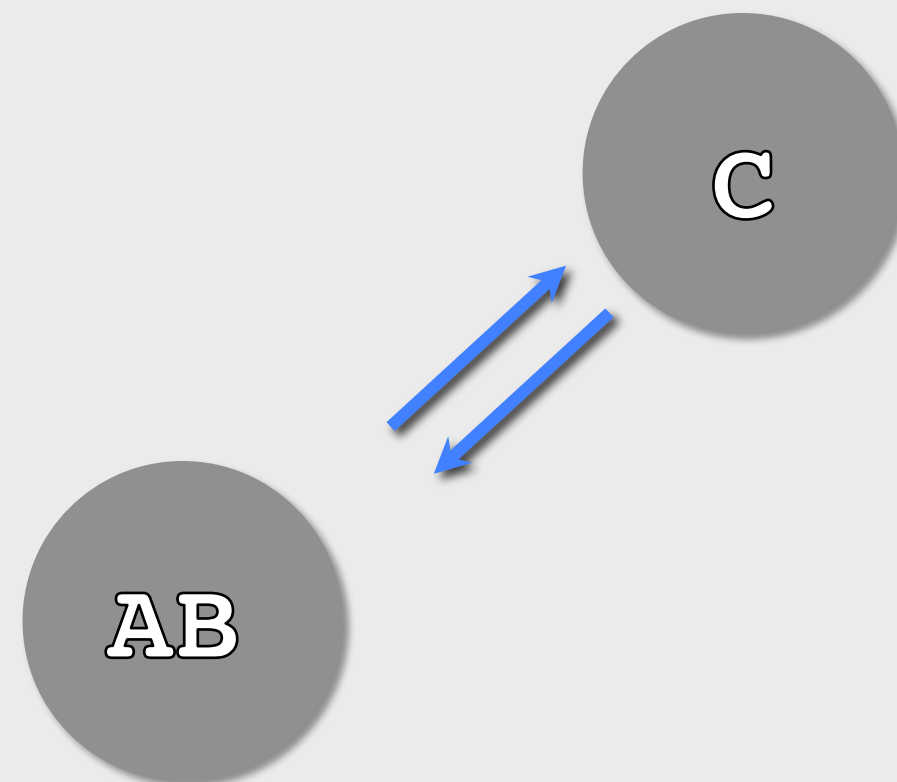
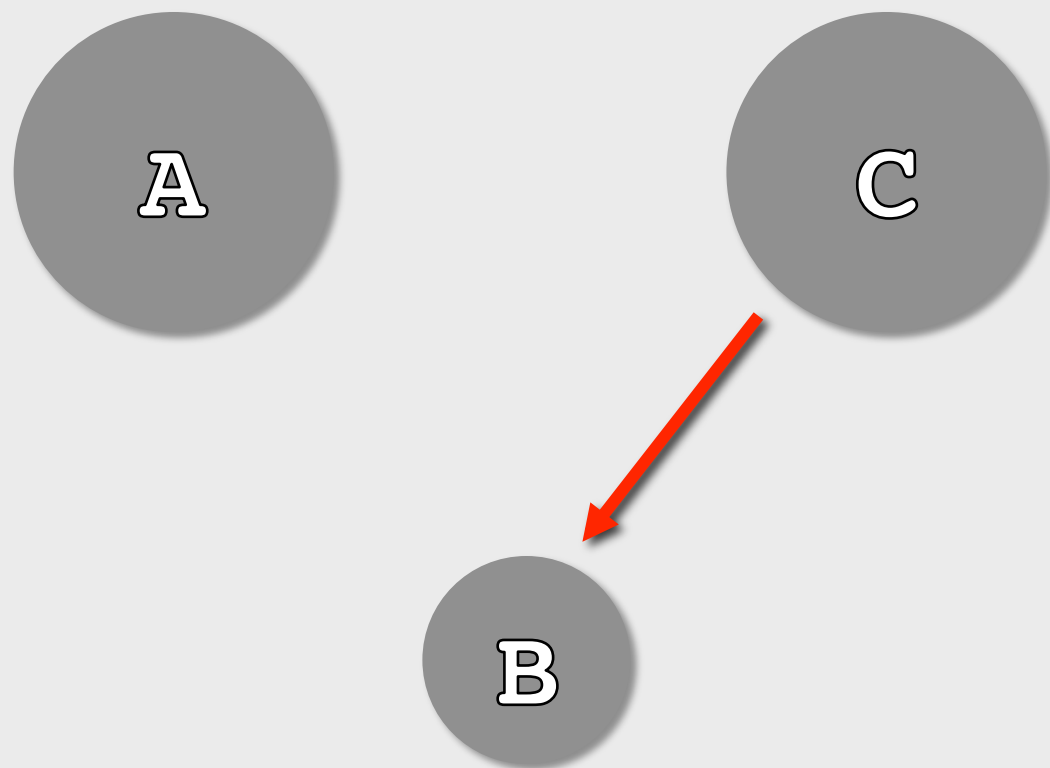
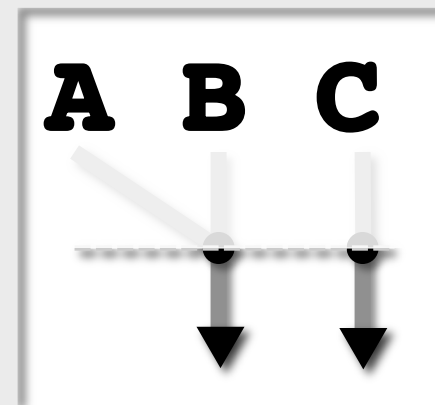
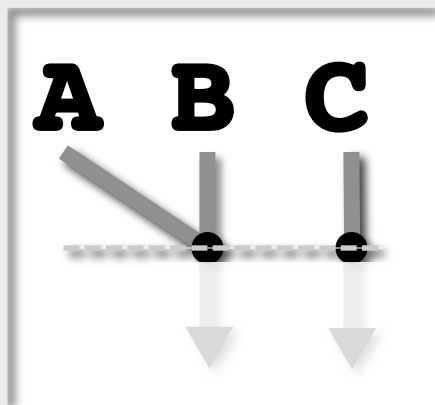




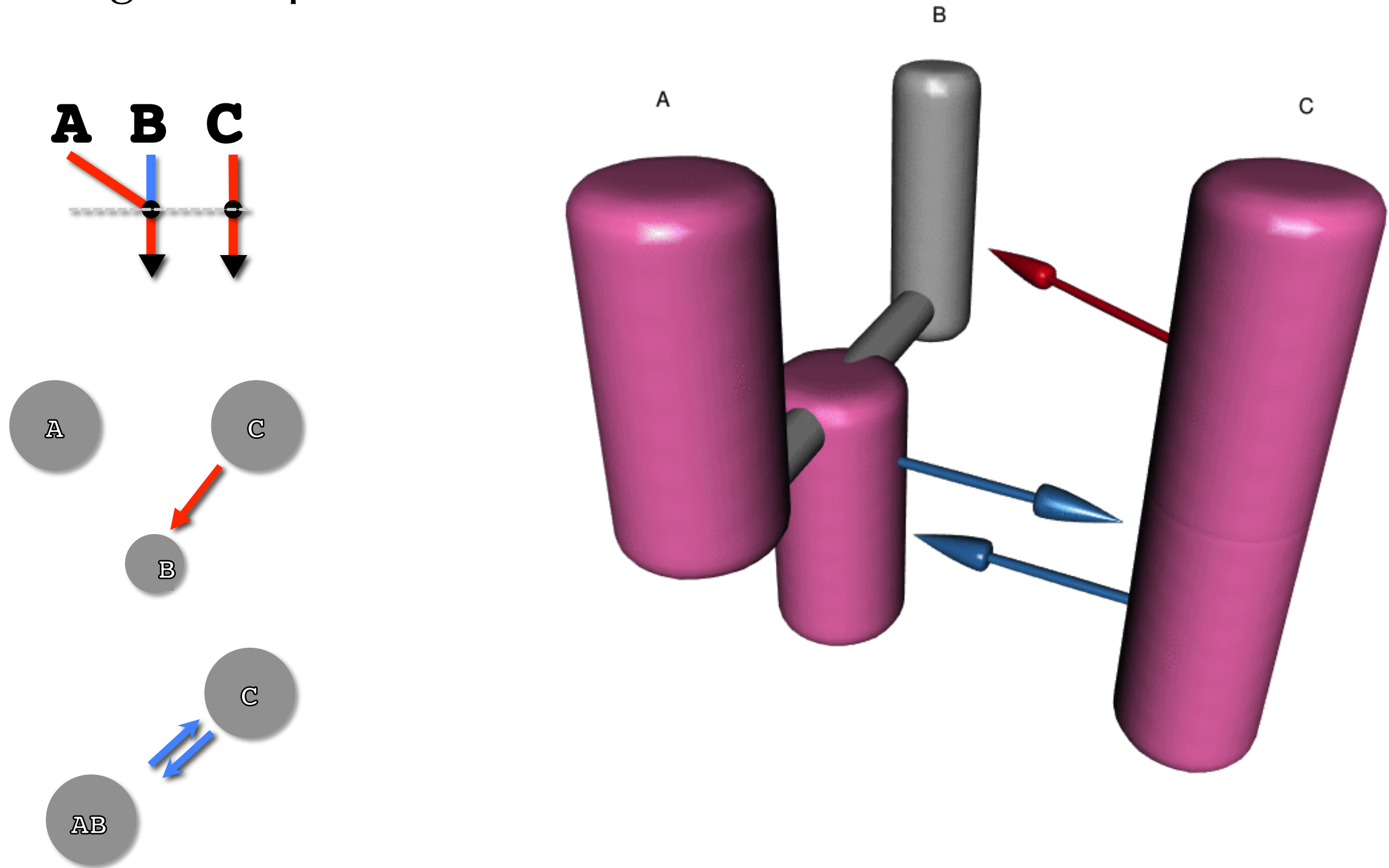


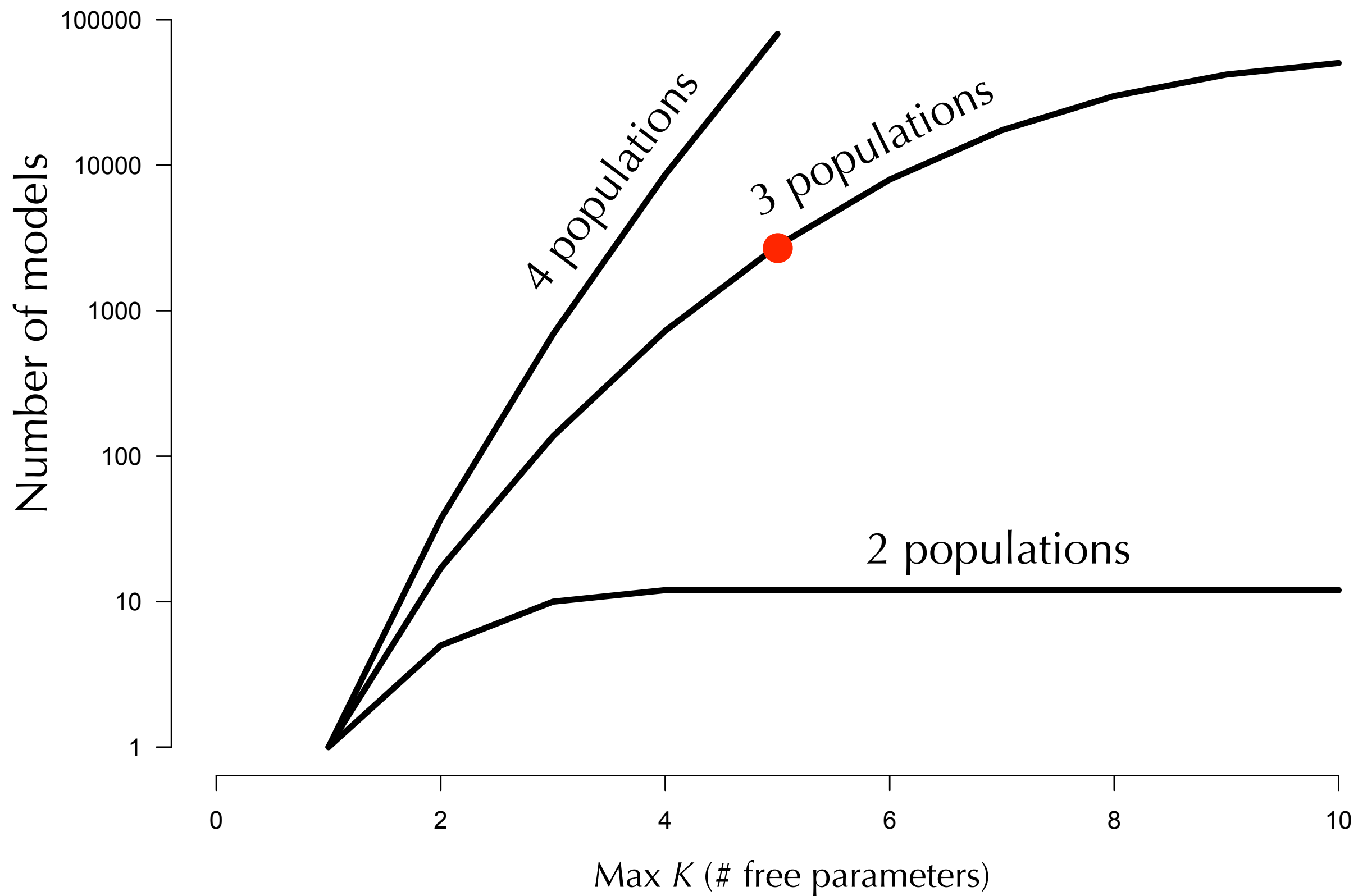




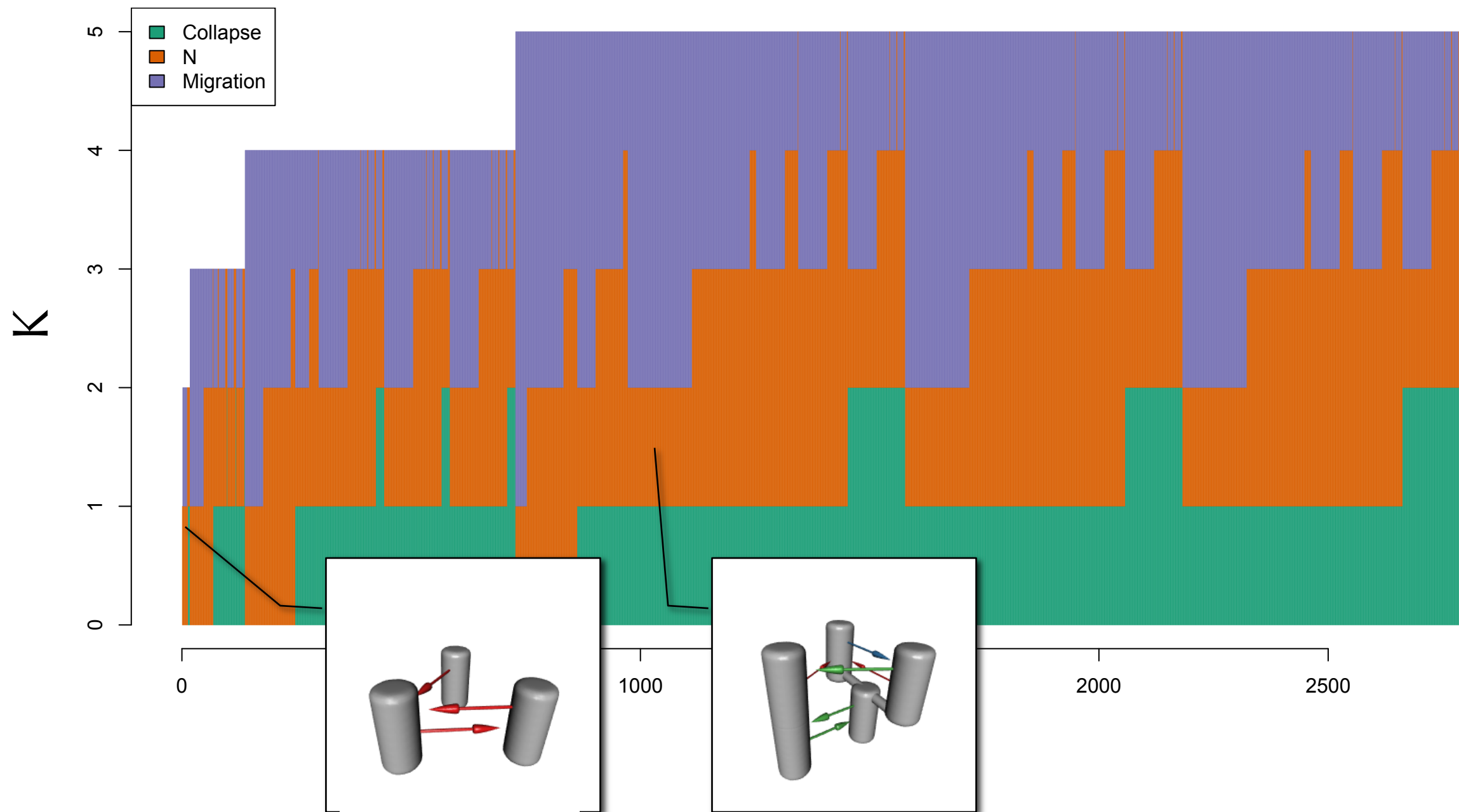


- 1 collapse parameter
- 2 population size parameters
- 2 migration parameters





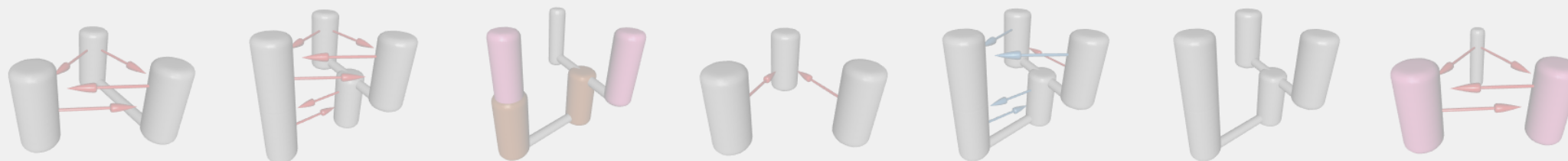
Model free parameters, three populations, max K=5



1

Generate all possible models

Given N populations, $\leq K$ free parameters



Filter

(only tree models, no more than two migration rates, etc.)

2



Analyze

(find best, find AIC for all)

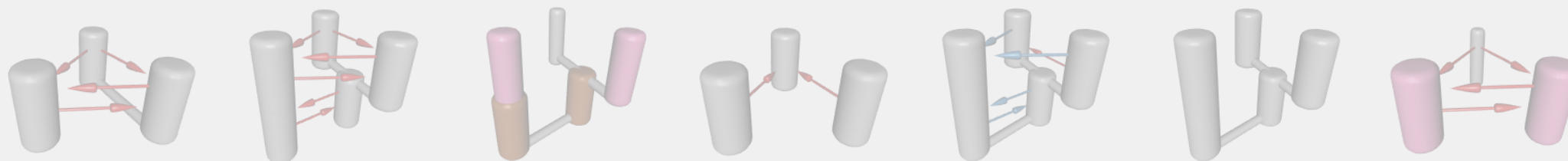
3



1

Generate all possible models

Given N populations, $\leq K$ free parameters



Filter

(only tree models, no more than two migration rates, etc.)

2

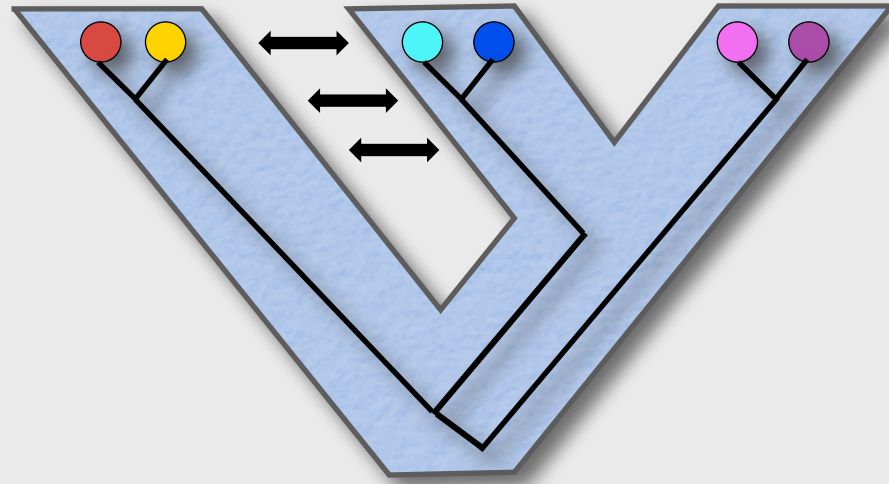


Analyze

(find best, find AIC for all)

3

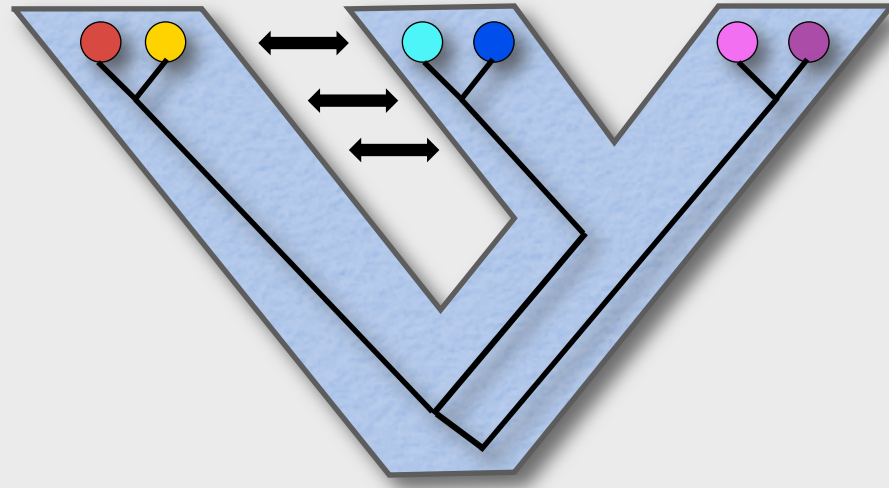




Exact likelihood: $P(\text{gene tree} \mid \text{population history})$

COAL (Degnan and Salter 2005), STEM (Kubatko et al. 2009), BPP (Yang and Rannala 2010), IM (Nielsen and Wakely 2001), etc.

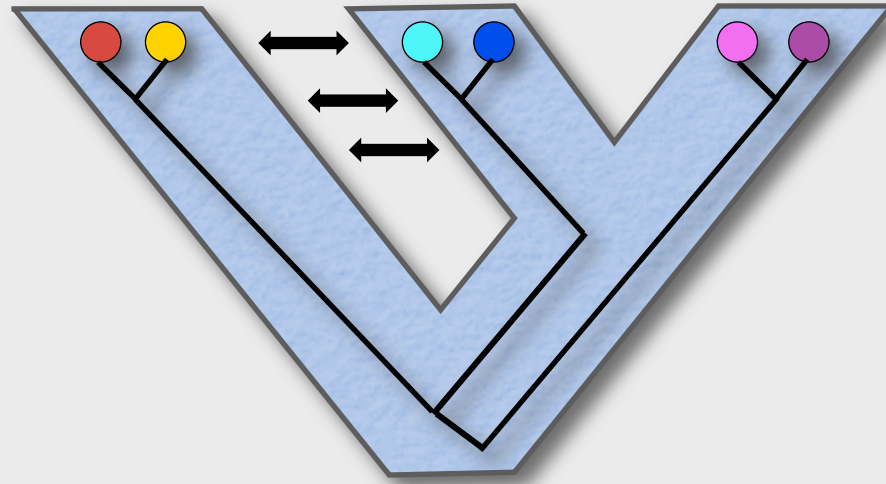
Issue: Not all models (currently) have analytical solutions, and it's always simpler to simulate (think of selection). But if it does work, can use it.



Approximate Bayesian Computation

Carmago et al. 2012, DIYABC (Cornuet et al. 2010), PopABC (Lopes et al. 2009)

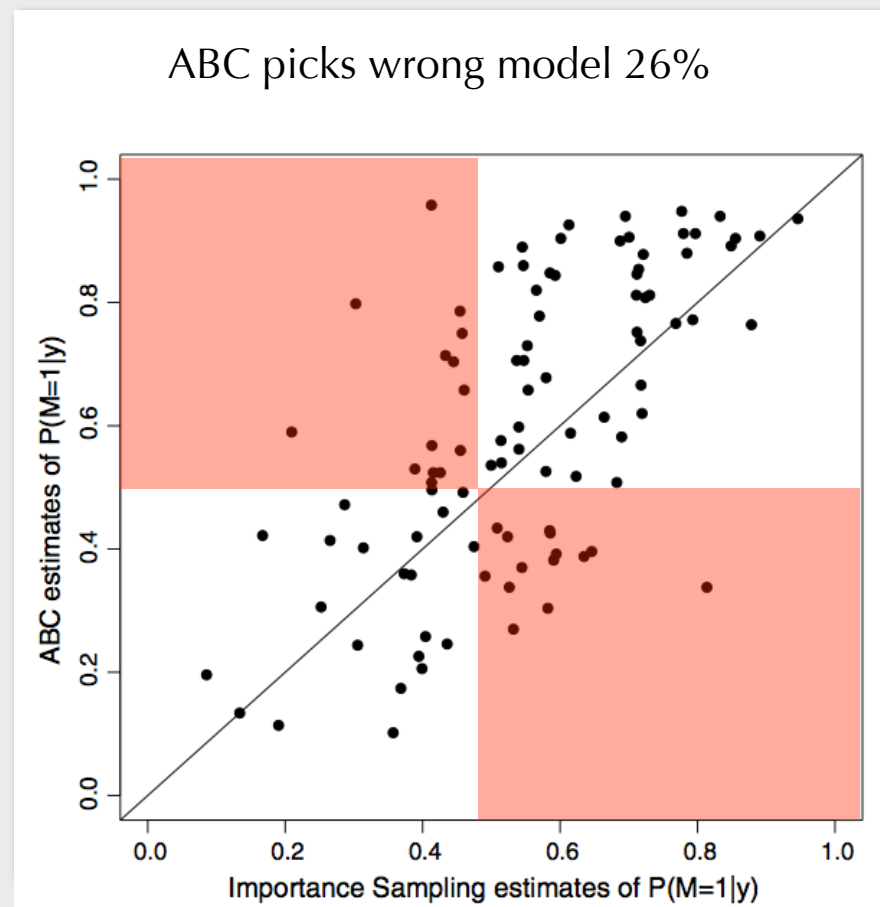
```
for (many reps) {  
  sample parameters (or models) from prior  
  simulate dataset using these parameters  
  if ( |(summary(dataset) - summary(observed)| <  $\epsilon$  ) {  
    count as match, save these parameters  
  }  
}  
  
plot(saved parameters) #or other way to summarize
```



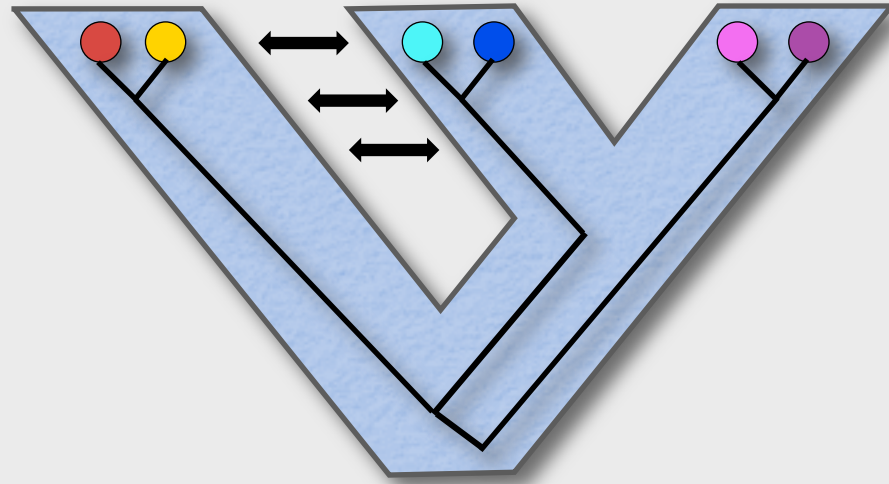
Approximate Bayesian Computation

Model choice hard to do with ABC: can depend on summary stats, priors, ϵ

Use of priors itself

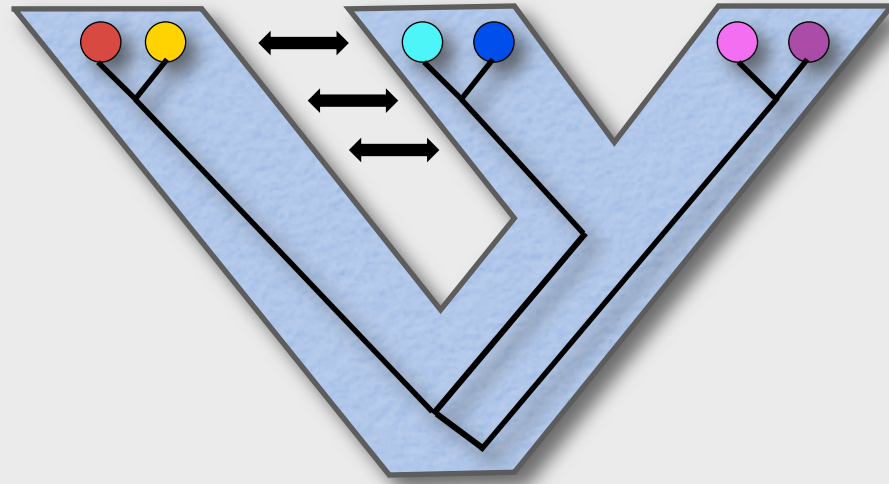


<https://www.facebook.com/LizardProtest>



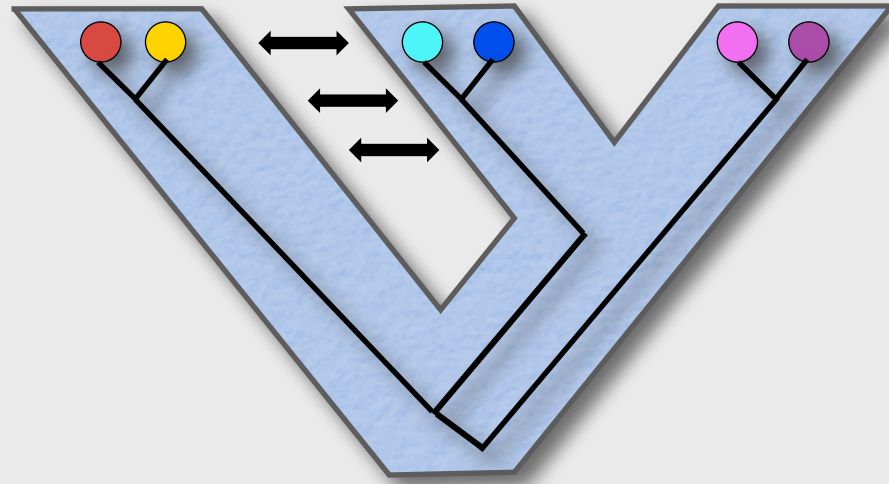
Approximate Likelihood

Estimate probability of data
counting exact matches



Approximate Likelihood

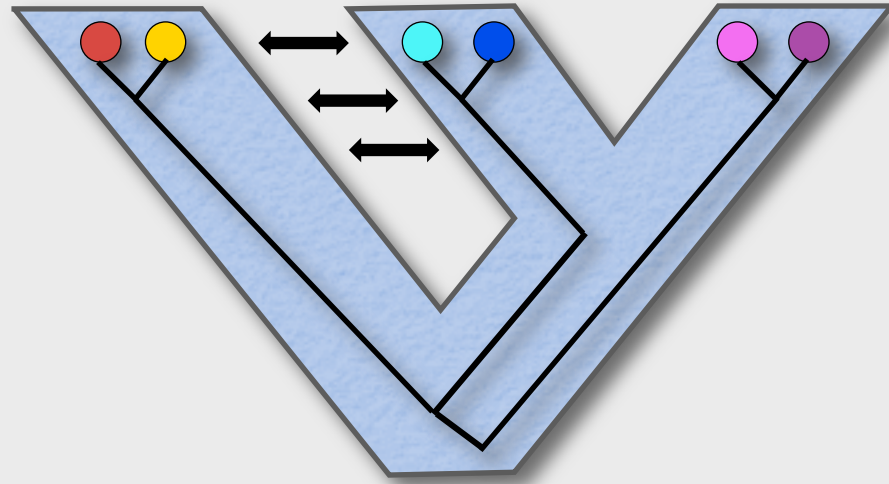
```
for (many reps) {  
  sample parameters (or models) from prior  
  simulate dataset using these parameters  
  if ( |(summary(dataset) - summary(observed)| <  $\epsilon$  ) {  
    count as match, save these parameters  
  }  
}  
  
plot(saved parameters) #or other way to summarize
```



Approximate Likelihood

```
sample parameters (or models) from prior guesses
for (many reps) {
  sample parameters (or models) from prior
  simulate dataset using these parameters
  if ( |(summary(dataset) - summary(observed)| <  $\epsilon$  ) {
    count as match, save these parameters
  }
}

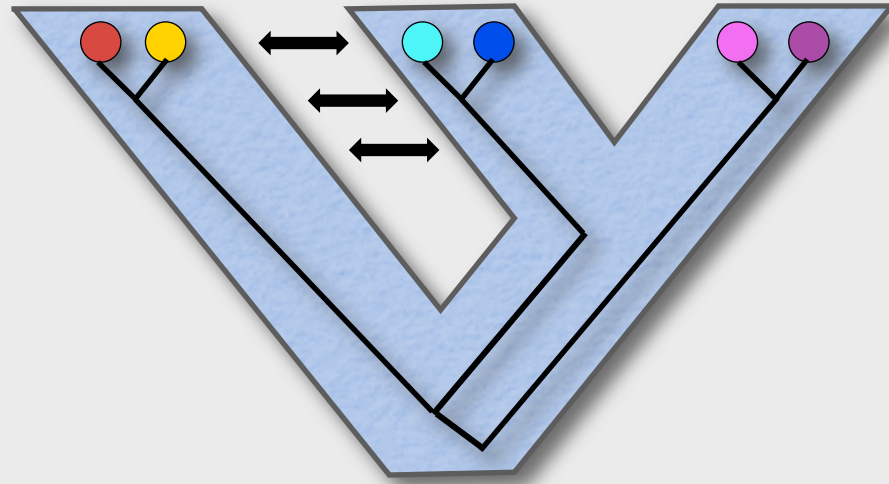
plot(saved parameters) #or other way to summarize
```



Approximate Likelihood

```
sample parameters (or models) from prior guesses
for (many reps) {
  sample parameters (or models) from prior
  simulate dataset using these parameters
  if ( topology(dataset) == topology(observed) ){
    count as match, save these parameters
  }
}

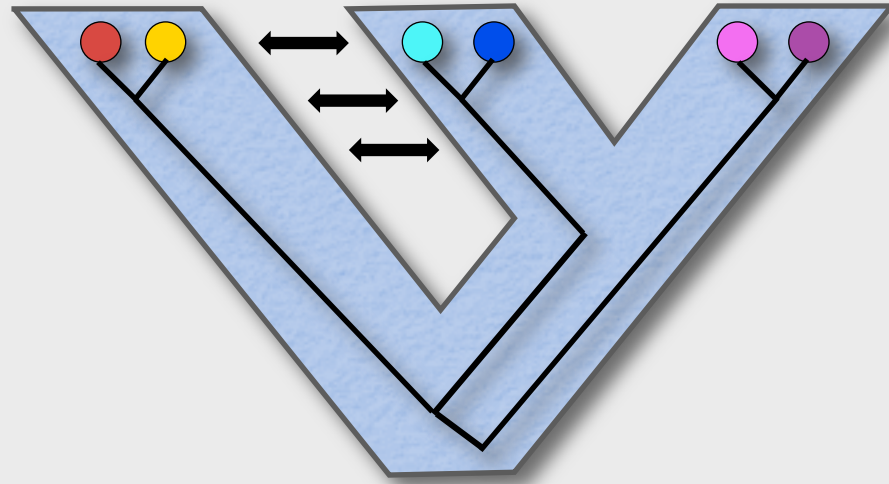
plot(saved parameters) #or other way to summarize
```



Approximate Likelihood

```
sample parameters (or models) from prior guesses
for (many reps) {
  sample parameters (or models) from prior
  simulate dataset using these parameters
  if ( topology(dataset) == topology(observed) ){
    count as match, save these parameters
  }
}

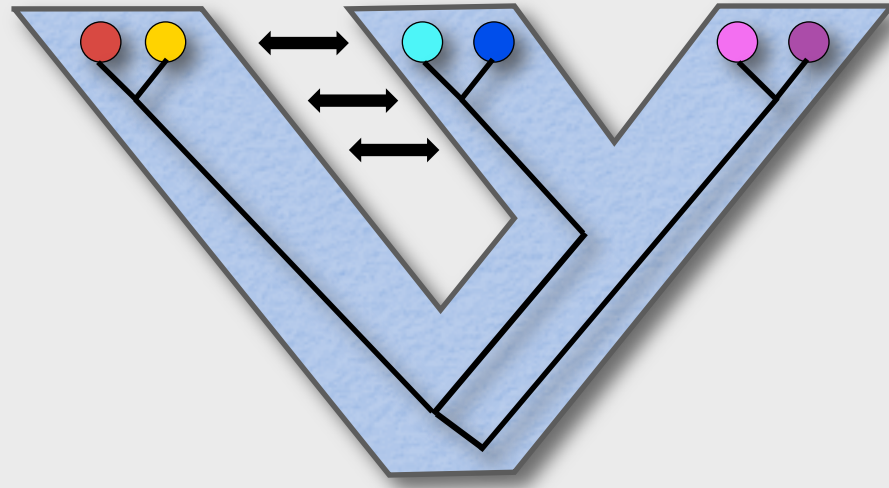
plot(saved parameters) #or other way to summarize
```

Approximate Likelihood

```
sample parameters (or models) from prior guesses
for (many reps) {
  sample parameters (or models) from prior
  simulate dataset using these parameters
  if ( topology(dataset) == topology(observed) ){
    count as match, save these parameters
  }
}
```

Likelihood = $\text{prob}(\text{topology}(\text{observed}) \mid \text{parameters})$



Approximate Likelihood

```
sample parameters (or models) from prior guesses
for (many reps) {
  sample parameters (or models) from prior
  simulate dataset using these parameters
  if ( topology(dataset) == topology(observed) ){
    count as match, save these parameters
  }
}
```

$\text{Likelihood} = \text{prob}(\text{topology}(\text{observed}) \mid \text{parameters}) \approx \text{\#matches} / \text{\#reps}$

Observed



Observed



Summary

4/7

Match

NA

Sim 1



3/7

N

Sim 2



5/7

N

Sim 3



1/7

N

Sim 4



4/7

N

Sim 5



2/7

N

Sim 6



4/7

Y

Sim 7



5/7

N

Observed



Summary

4/7

Match

NA

Sim 1



3/7 ✓

N

Sim 2



5/7 ✓

N

Sim 3



1/7

N

Sim 4



4/7 ✓

N

Sim 5



2/7

N

Sim 6



4/7 ✓

Y ✓

Sim 7



5/7 ✓

N

Approximate Bayesian Computation

Summary stats as functions of the data

Simulation

Approximate match

Priors

Integration

Approximate Likelihood

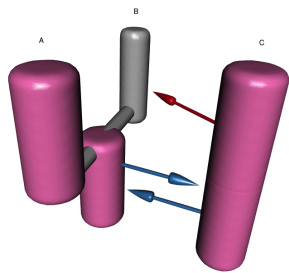
Topologies (= the data)

Simulation

Exact match

Starting values, but optimize

Maximum



Collapse 1

4

Pop size 1

1

Pop size 2

3

Migration 1

0.2

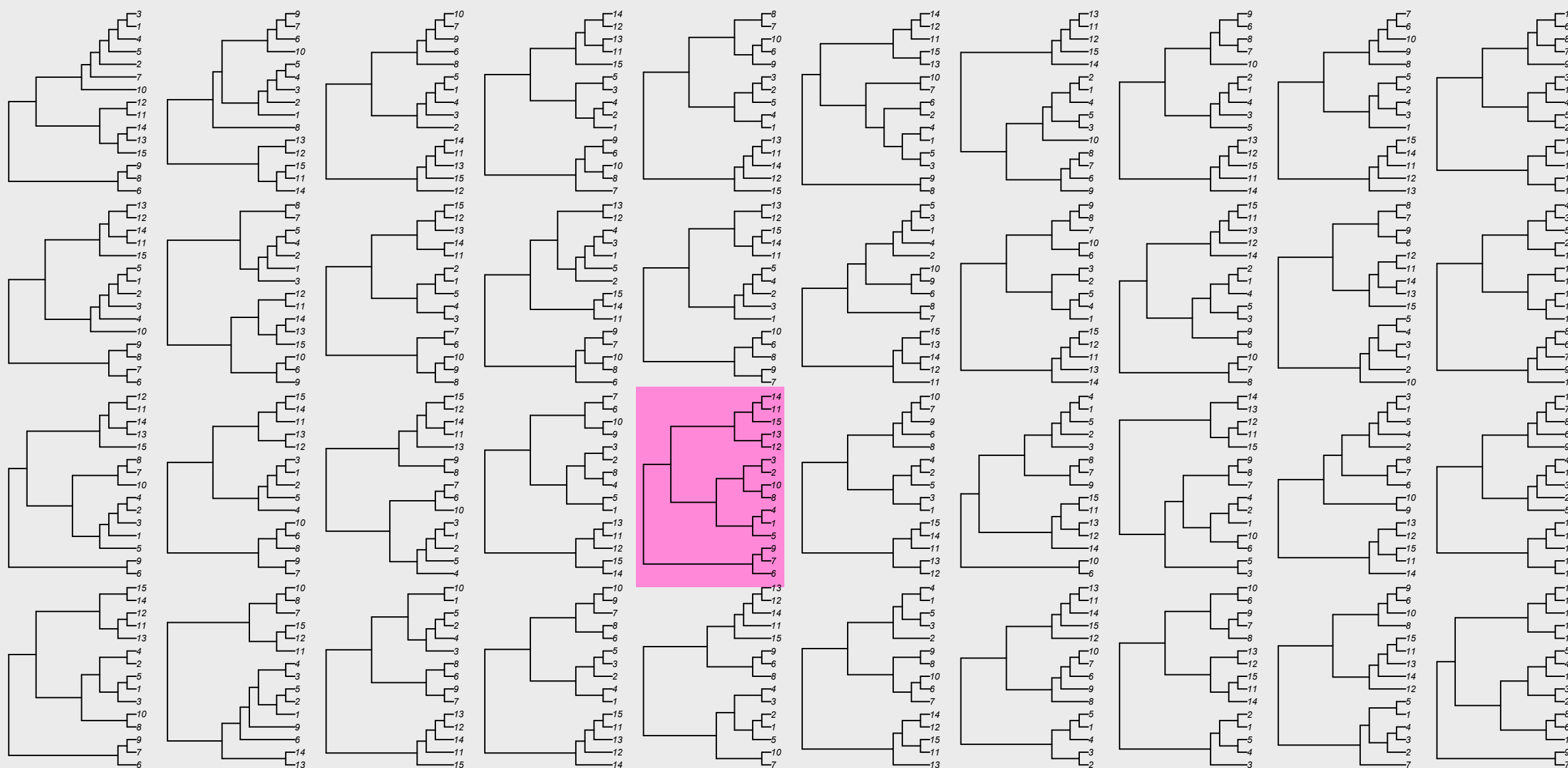
Migration 2

0.6

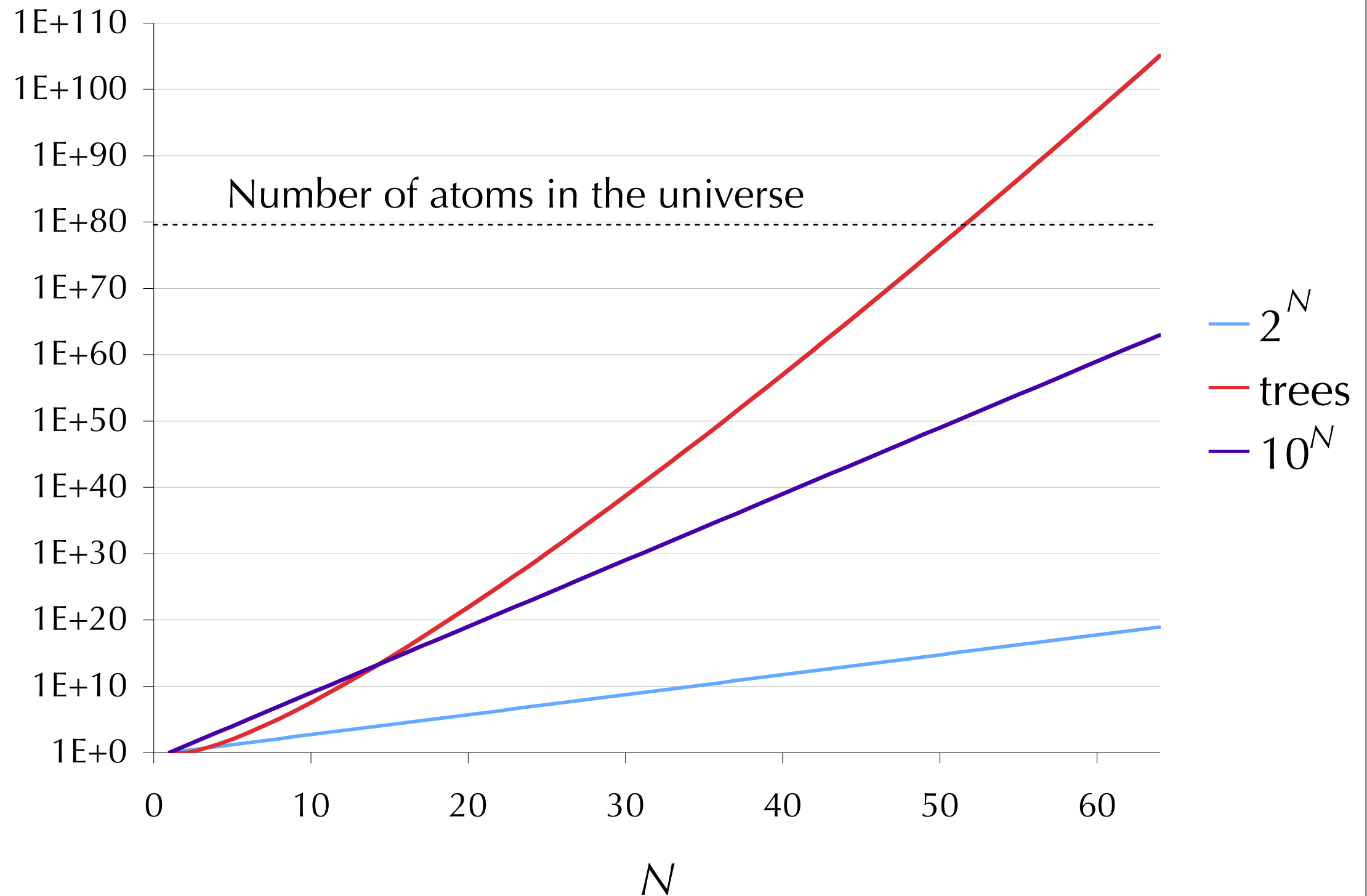
Generate 100,000 trees, three populations, five samples per population (~5 sec):

```
ms 15 100000 -T -I 3 5 5 5 -n 1 1.0 -n 2 3.0 -n 3 3.0 -ma x 0.0 0.0 0.2 x 0.0 0.0 0.0 x
-ej 4.0 3 1 -en 4.0 1 3.0 -en 4.0 2 3.0 -em 4.0 1 2 0.6 -em 4.0 2 1 0.6
```

ms (Hudson 2002)

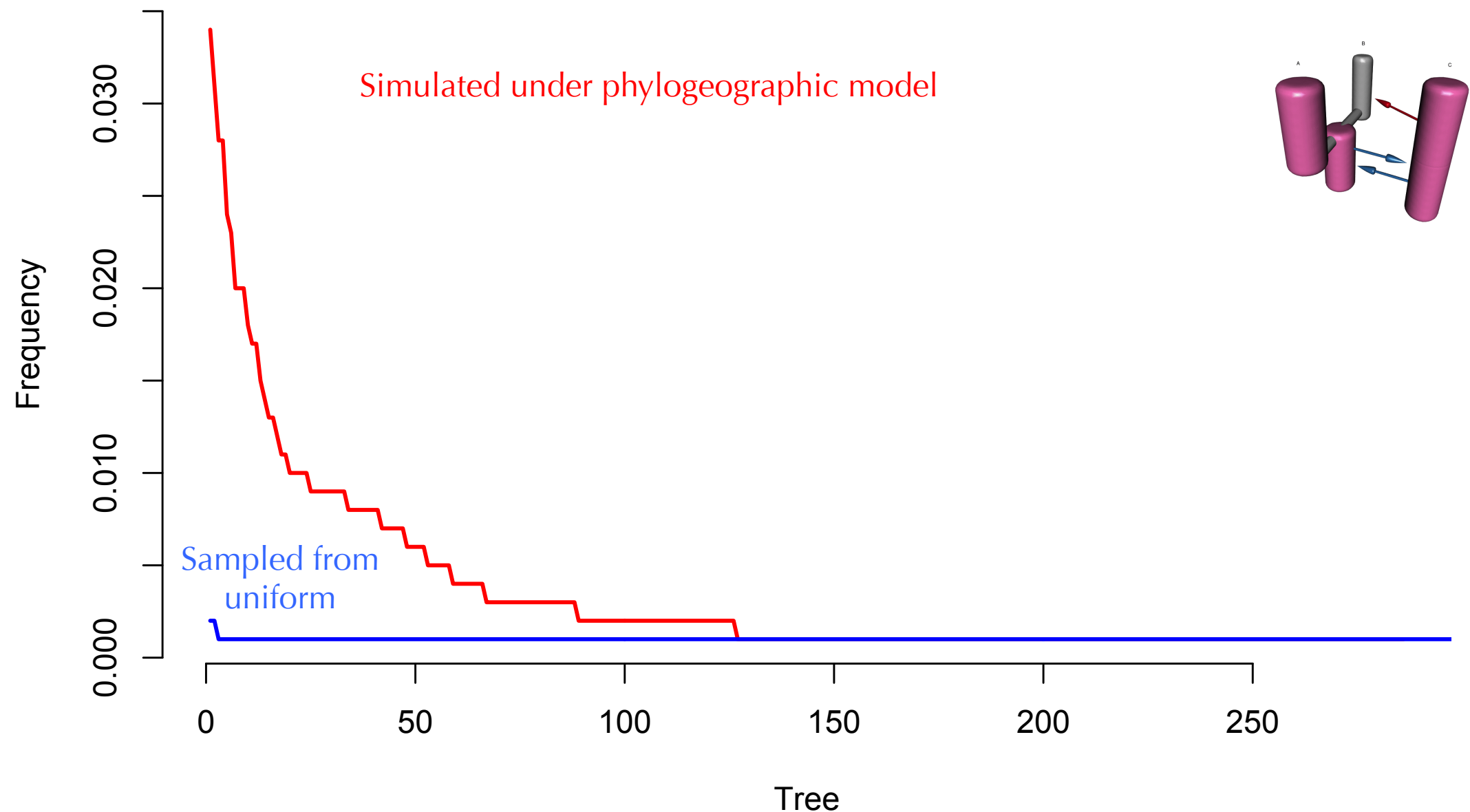


A dumb idea?



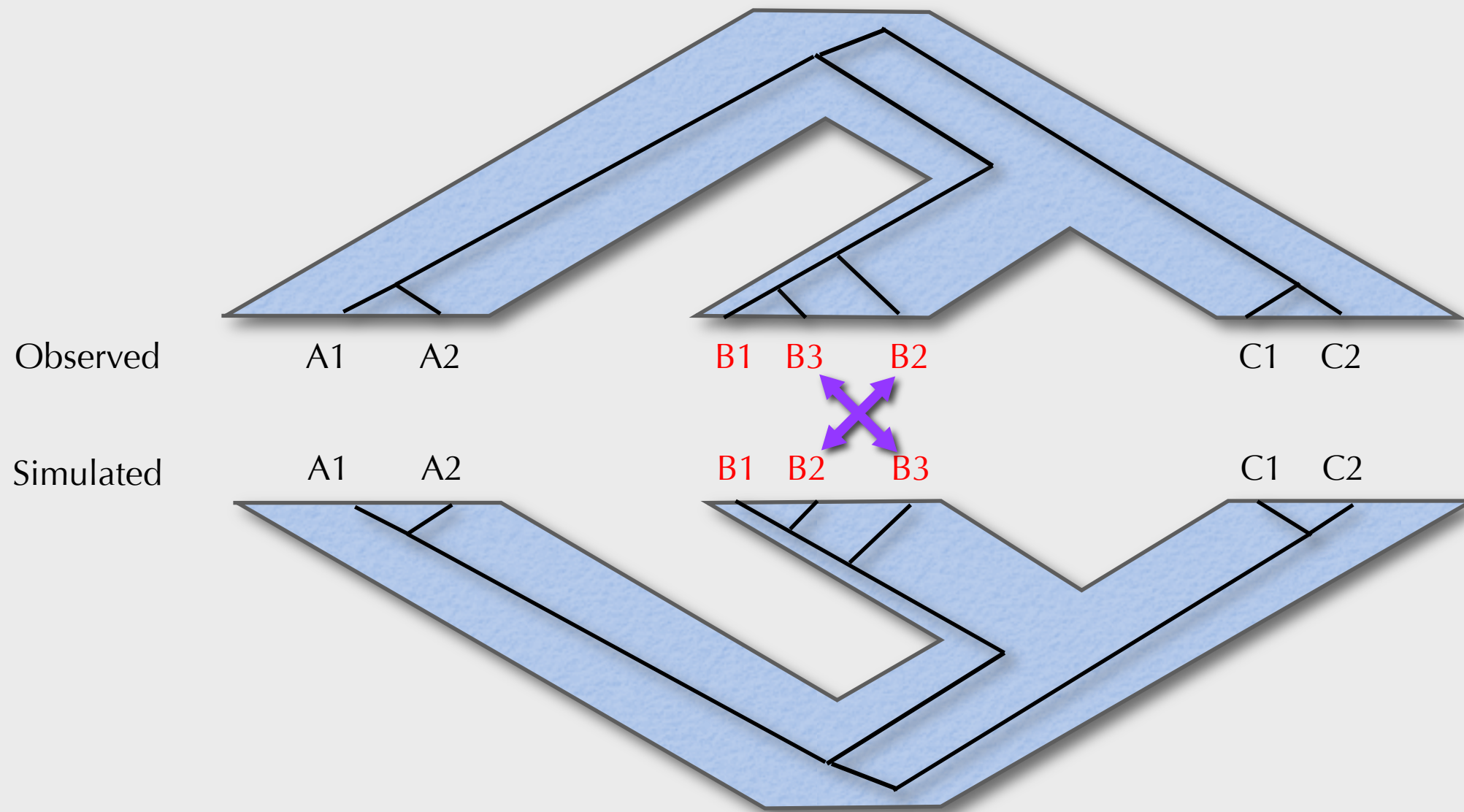
A dumb idea?

Tree probabilities not uniform (some trees *much* more likely than others)



A dumb idea?

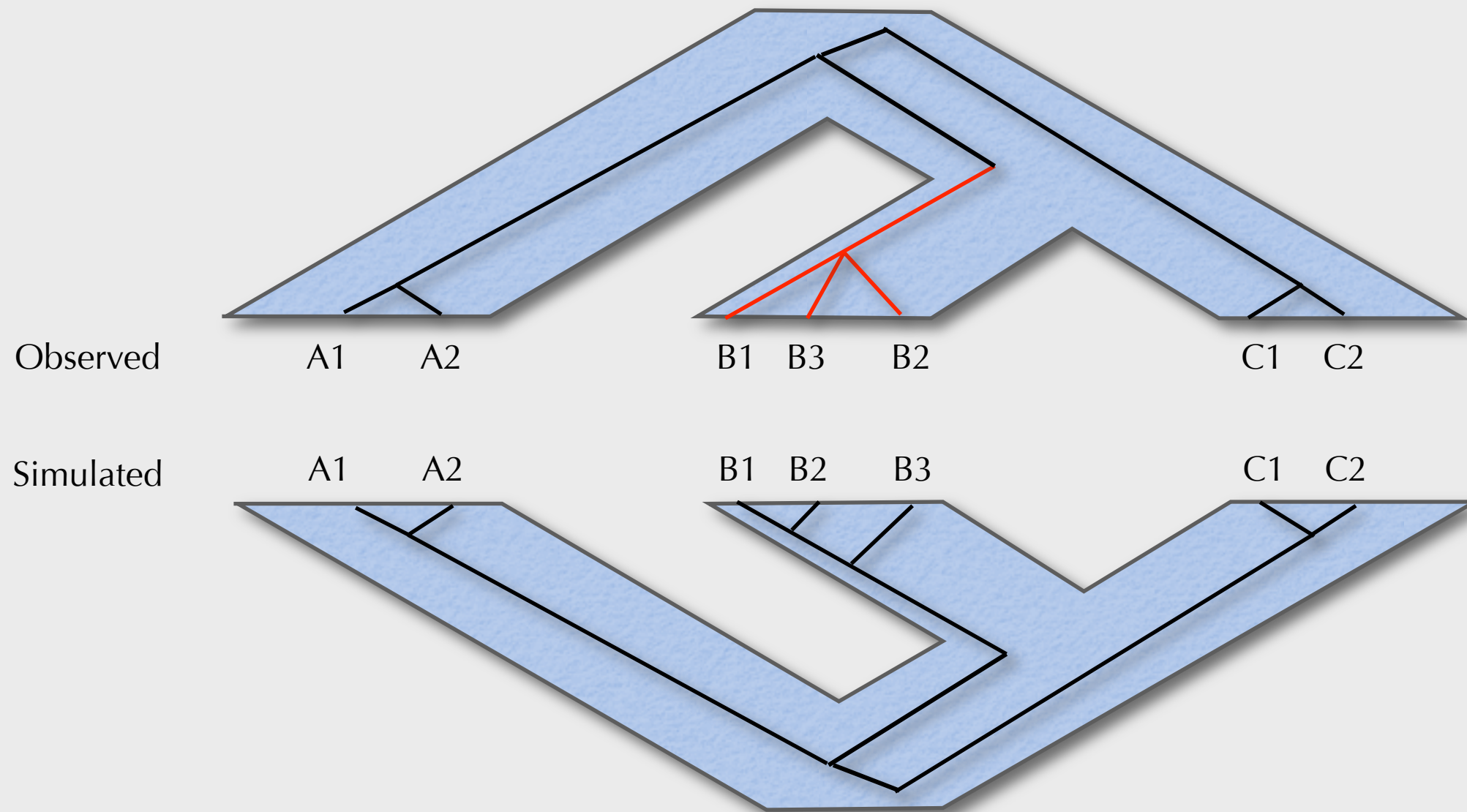
Clever idea 1: Sample labels within populations arbitrary



Match based on all possible labeling, then correct for this
i.e., three possible permutations, so if there is a match divide by 3 to get probability

A dumb idea?

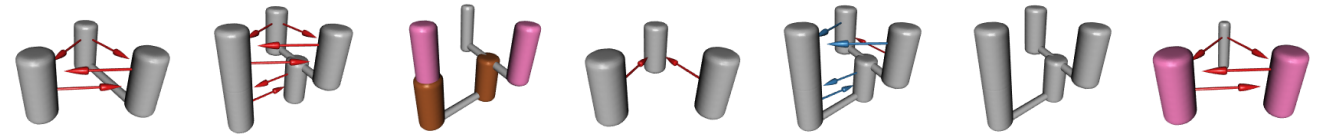
Clever idea 2: Polytomies are soft in gene trees (optional)



Match based on all possible resolutions, then correct for this

Grows (quickly) with
number of populations

Generate all possible models



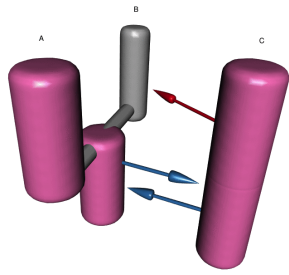
Filter



Analyze



Grows quickly with
number of samples per
population, but only
linearly with number
of genes



Collapse 1

4

Pop size 1

1

Pop size 2

3

Migration 1

0.2

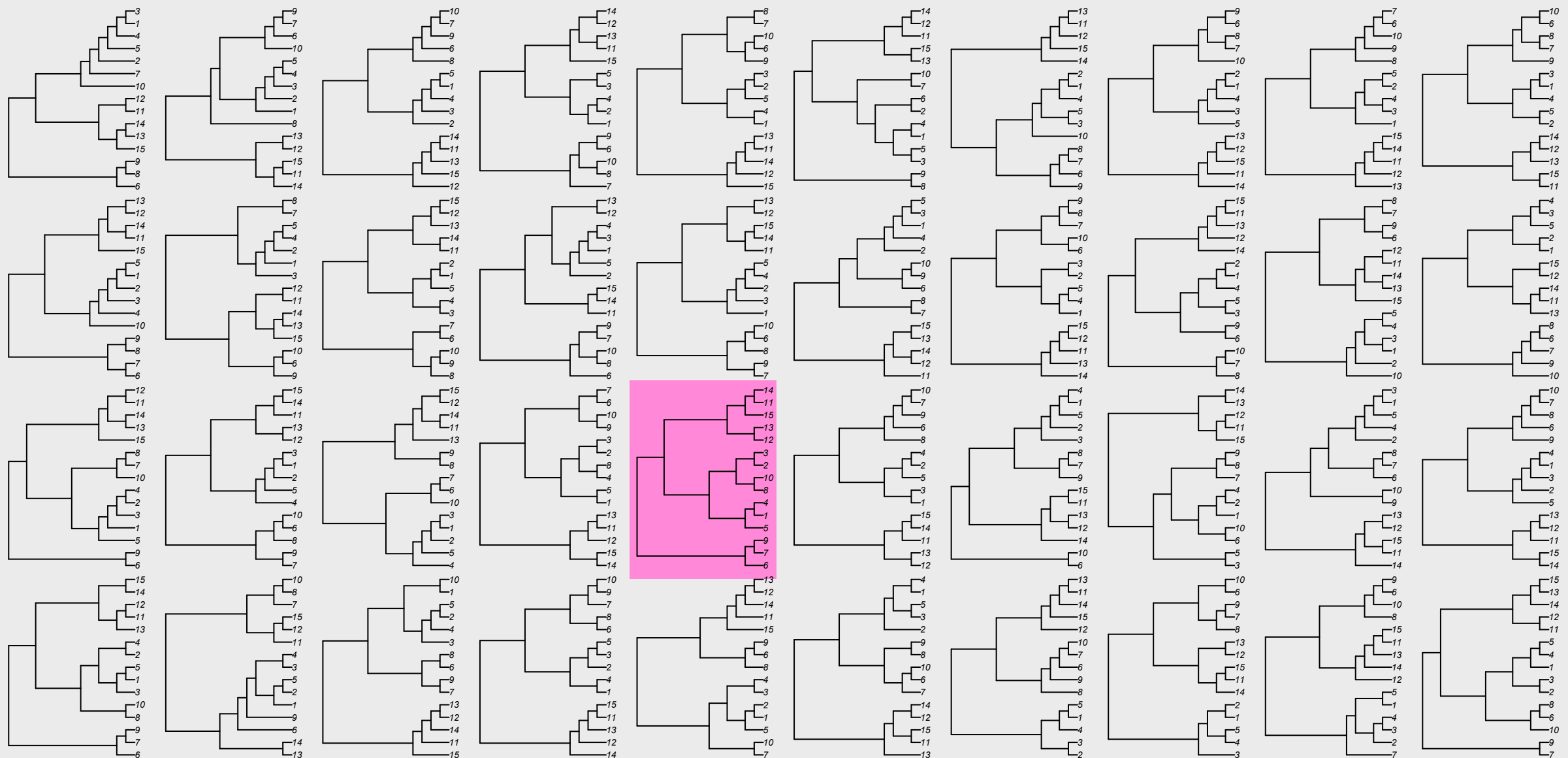
Migration 2

0.6

Generate 100,000 trees, three populations, five samples per population (~5 sec):

```
ms 15 100000 -T -I 3 5 5 5 -n 1 1.0 -n 2 3.0 -n 3 3.0 -ma x 0.0 0.0 0.2 x 0.0 0.0 0.0 x
-ej 4.0 3 1 -en 4.0 1 3.0 -en 4.0 2 3.0 -em 4.0 1 2 0.6 -em 4.0 2 1 0.6
```

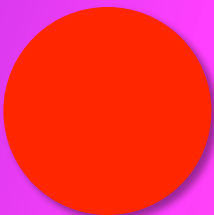
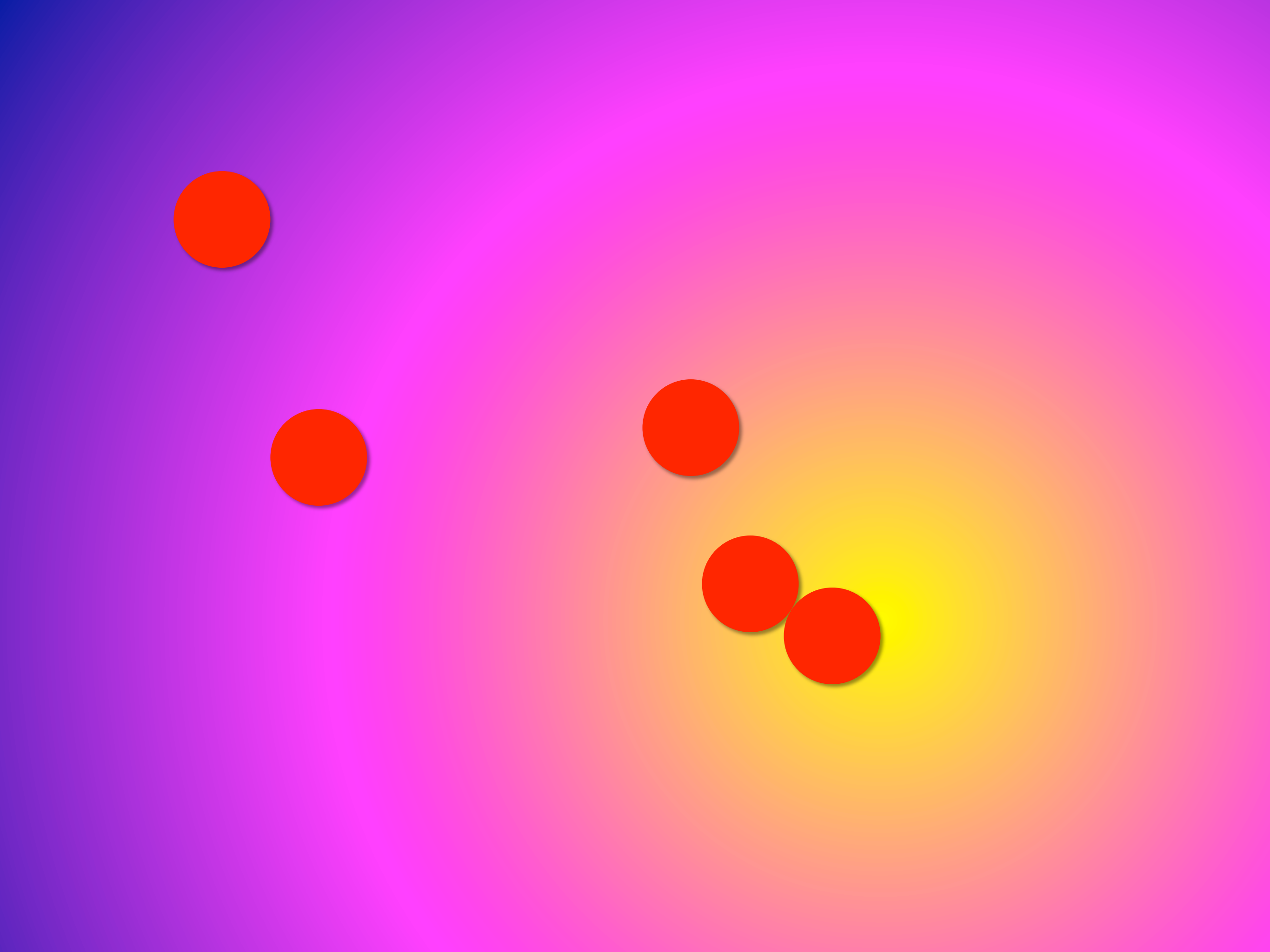
ms (Hudson 2002)



Optimization

Numerical parameters

Collapse 1	Pop size 1	Pop size 2	Migration 1	Migration 2
4	1	3	0.2	0.6

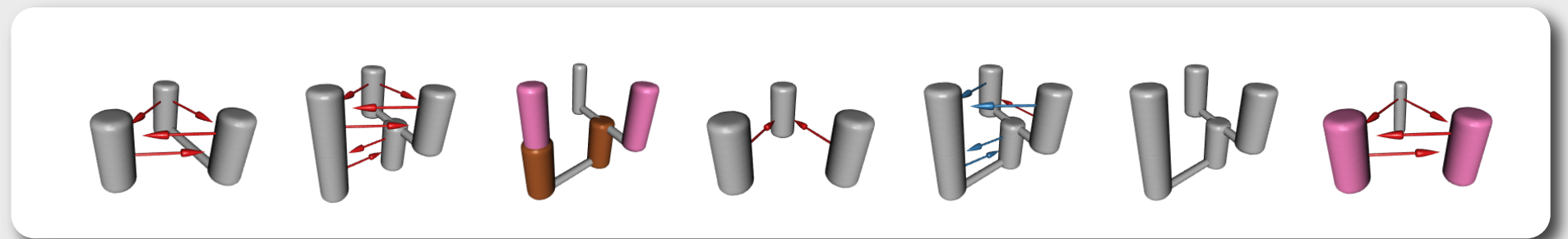


Optimization

Numerical parameters

Model

Exhaustive

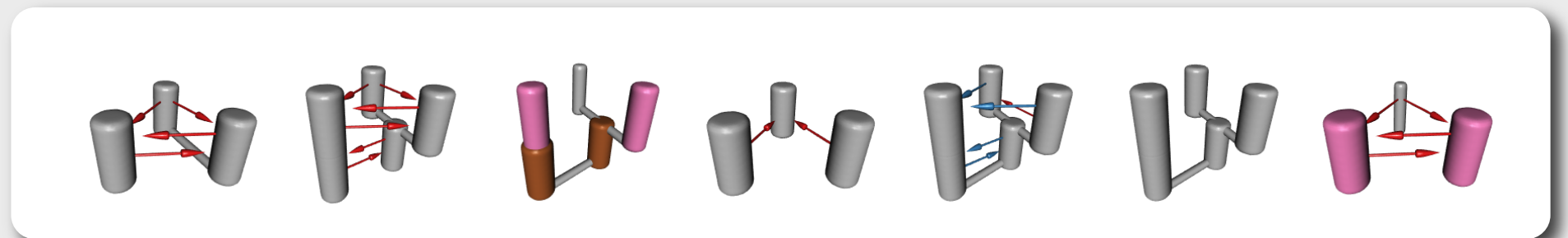


Optimization

Numerical parameters

Model

Exhaustive



Genetic

Break model into genes (one gene for collapse structure, one for migration structure, etc.)
Let them evolve, AIC is fitness

1.2

Collapse Model	Pop size Model	Migration Model	Continuous parameters
1	1	1	

0.0

Collapse Model	Pop size Model	Migration Model	Continuous parameters
2	1	3	

5.0

Collapse Model	Pop size Model	Migration Model	Continuous parameters
3	4	2	

20.5

Collapse Model	Pop size Model	Migration Model	Continuous parameters
1	2	1	

2.0

Collapse Model	Pop size Model	Migration Model	Continuous parameters
2	2	2	



1.2

Collapse Model	Pop size Model	Migration Model	Continuous parameters
1	1	1	...

0.0

Collapse Model	Pop size Model	Migration Model	Continuous parameters
2	1	3	...

5.0

Collapse Model	Pop size Model	Migration Model	Continuous parameters
3	4	2	...

20.5

Collapse Model	Pop size Model	Migration Model	Continuous parameters
1	2	1	...

2.0

Collapse Model	Pop size Model	Migration Model	Continuous parameters
2	2	2	...

•
•
•

1.2

Collapse Model	Pop size Model	Migration Model	Continuous parameters
1	1	1	...

0.0

Collapse Model	Pop size Model	Migration Model	Continuous parameters
2	1	3	...

Collapse Model	Pop size Model	Migration Model	Continuous parameters
			...

Collapse Model	Pop size Model	Migration Model	Continuous parameters
			...

2.0

Collapse Model	Pop size Model	Migration Model	Continuous parameters
2	2	2	...



1.2

Collapse Model	Pop size Model	Migration Model	Continuous parameters
1	1	1	...

0.0

Collapse Model	Pop size Model	Migration Model	Continuous parameters
2	1	3	...


Collapse Model	Pop size Model	Migration Model	Continuous parameters
2	1	1	...

Collapse Model	Pop size Model	Migration Model	Continuous parameters
			...

2.0

Collapse Model	Pop size Model	Migration Model	Continuous parameters
2	2	2	...





```

createModelSpaceNloptr<-function(migrationArrayMap, migrationArray, popVector, print.ms.string=FALSE, badAIC=1000000000000000,
parameterValue=100, nTrees=1, msLocation="/usr/local/bin/
compareLocation="comparecladespipe.pl", assign="assign.txt", observed="observed.txt", unresolvedTest=TRUE, debug=FALSE, method="nlminb",
itnmax=NULL, pop.size=50, print.results=FALSE, maxtime=0, maxeval=0, ...) {
915 Domains<-matrix(ncol=2, nrow=3)
916 Domains[1,]<-range(migrationArrayMap$collapseMatrix.number)
917 Domains[2,]<-range(migrationArrayMap$n0multiplierMap.number)
918 Domains[3,]<-range(migrationArrayMap$migrationArray.number)
919
920 results<-genoud(searchContinuousModelSpaceNloptr, nvars=3, max=FALSE, starting.values=c(1,1,1), MemoryMatrix=TRUE, boundary.enforcement=2,
data.type.int=TRUE, Domains=Domains, migrationArrayMap=migrationArrayMap, migrationArray=migrationArray, popVector=popVector,
print.ms.string=print.ms.string, badAIC=badAIC, maxParameterValue=maxParameterValue,
nTrees=nTrees, msLocation=msLocation, compareLocation=compareLocation, assign=assign, observed=observed, unresolvedTest=unresolvedTest,
debug=debug, method=method, itnmax=itnmax, pop.size=pop.size, print.results=print.results, maxtime=maxtime, maxeval=maxeval,
return.all=FALSE, ...)
921 return(results)
922 }

```



```

sub returnclades {
  my ($intree, %localassignments) = @_;
  my @cladearray=();
  $intree=~s/:\d+\.\.?d*e?\-\?d*//ig; #remove branch length (regex from Olaf Bininda-Emonds'
partitionmetric.pl script)
  while ($intree=~m/(\([\w\,]+\))/) {
    my $matchstring=$1;
    my $innermatch="";
    if ($matchstring=~m/\([\w\,]+\)/) {
      $innermatch=$1;
    }
    my @splitinner=split(/\,/,$innermatch);
    my @renamedinner=();
    foreach my $inner (@splitinner) {

```

When the option -T is used the trees representing the history of the sampled chromosomes are output. For example, the command line `ms 5 2 -T` results in the following output:

```

ms 5 2 -T
3579 27011 59243

//
((2:0.074,5:0.074):0.296,(1:0.311,(3:0.123,4:0.123):0.187):0.060);

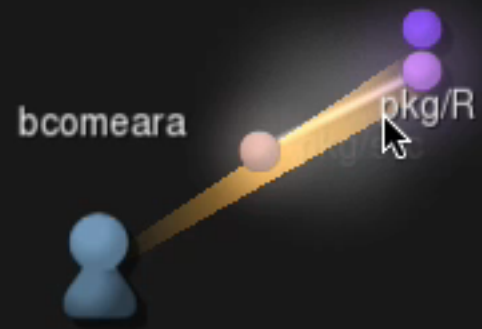
//
(2:1.766,(4:0.505,(3:0.222,(1:0.163,5:0.163):0.059):0.283):1.261);

```

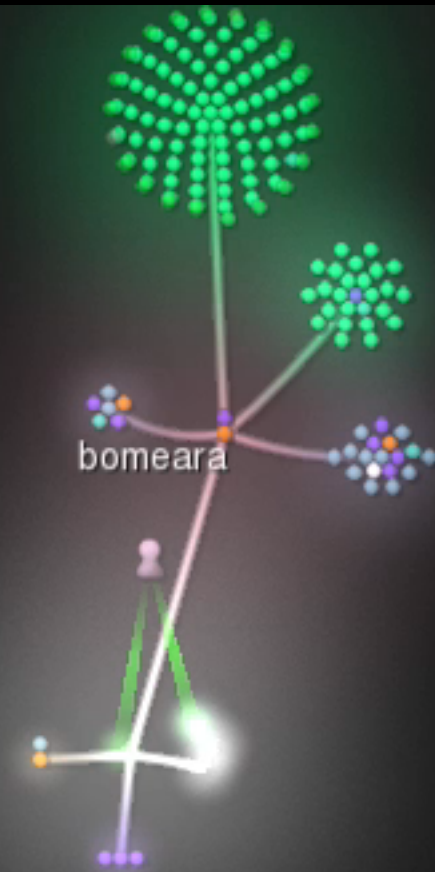
ms
Hudson
2002

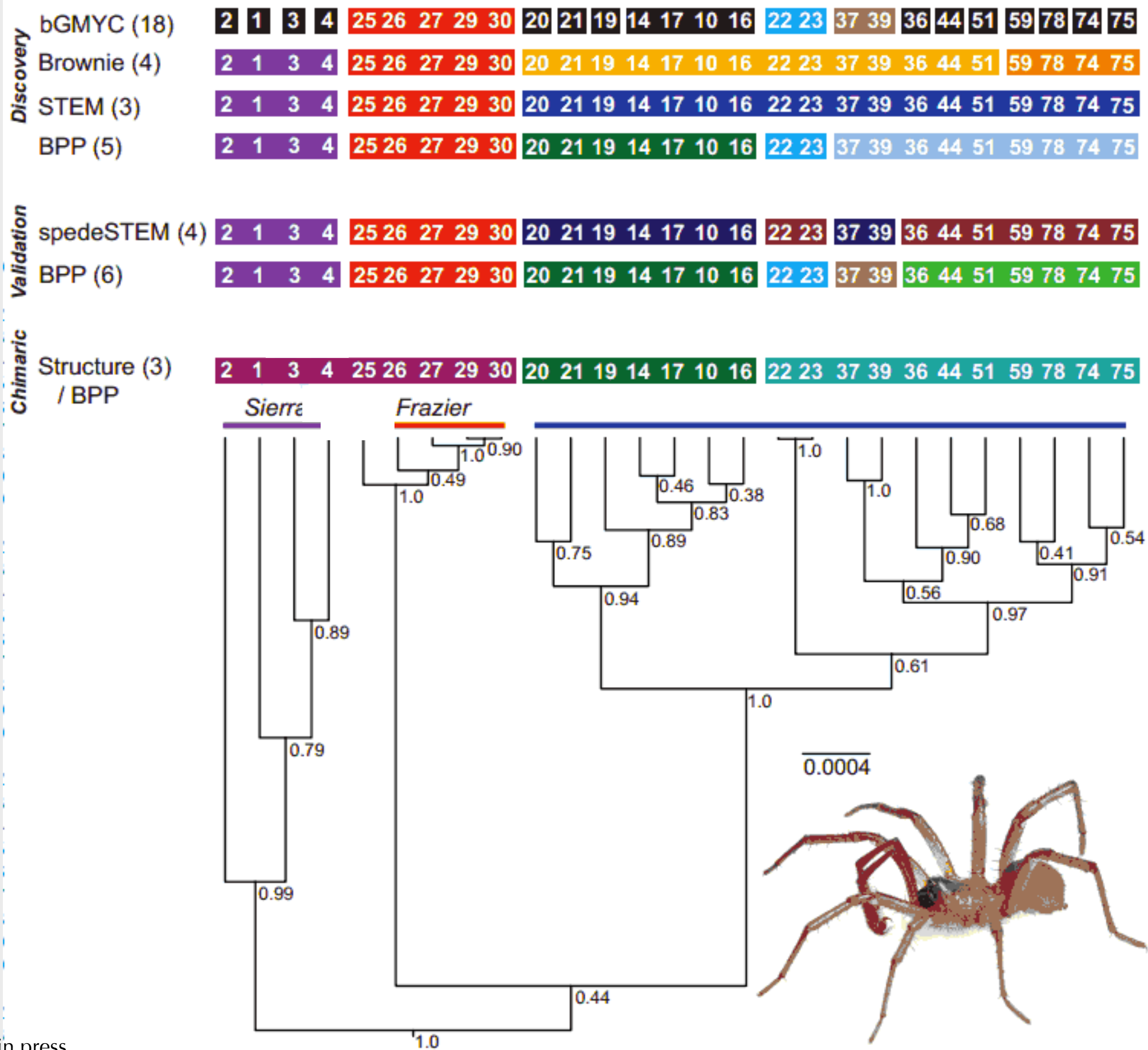
Code

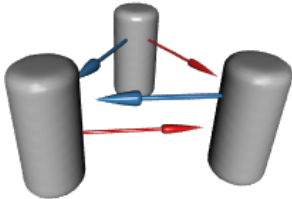
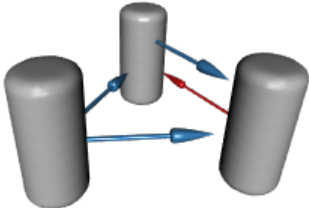
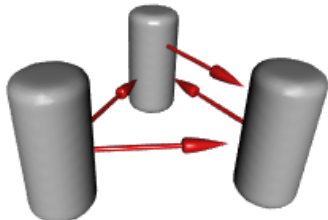
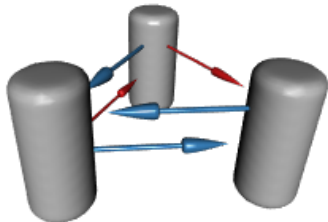
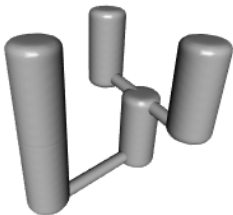
11/27/11



Analyses

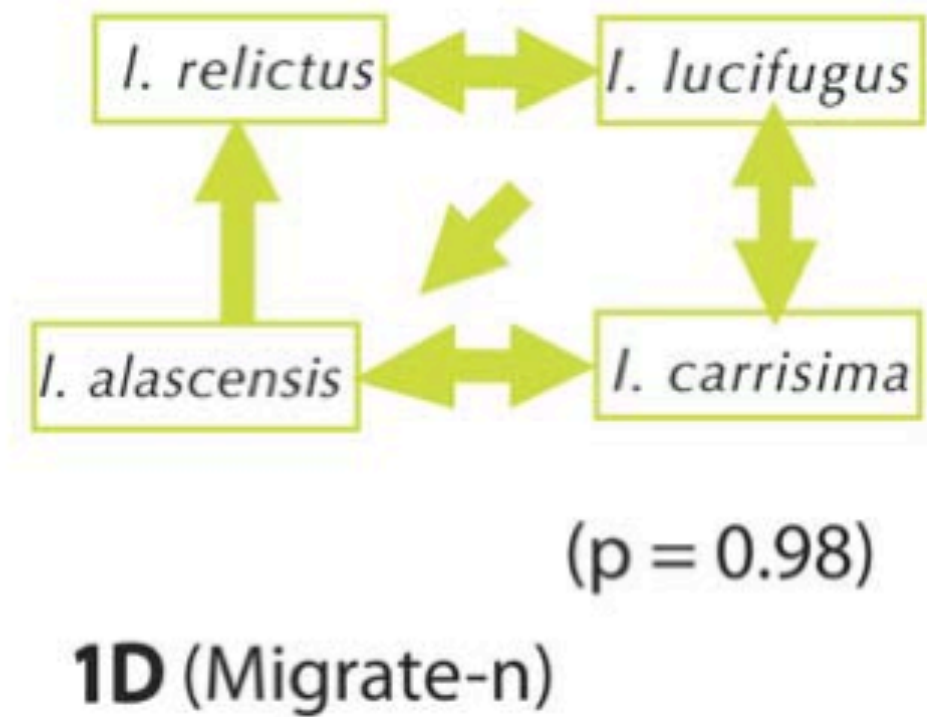
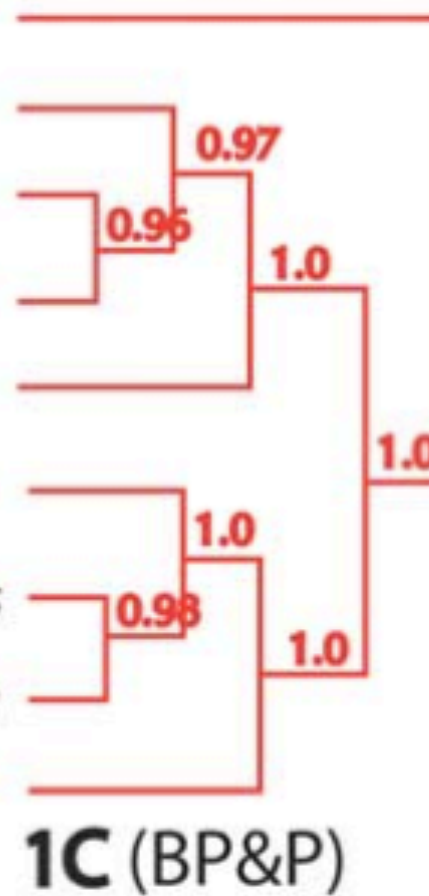
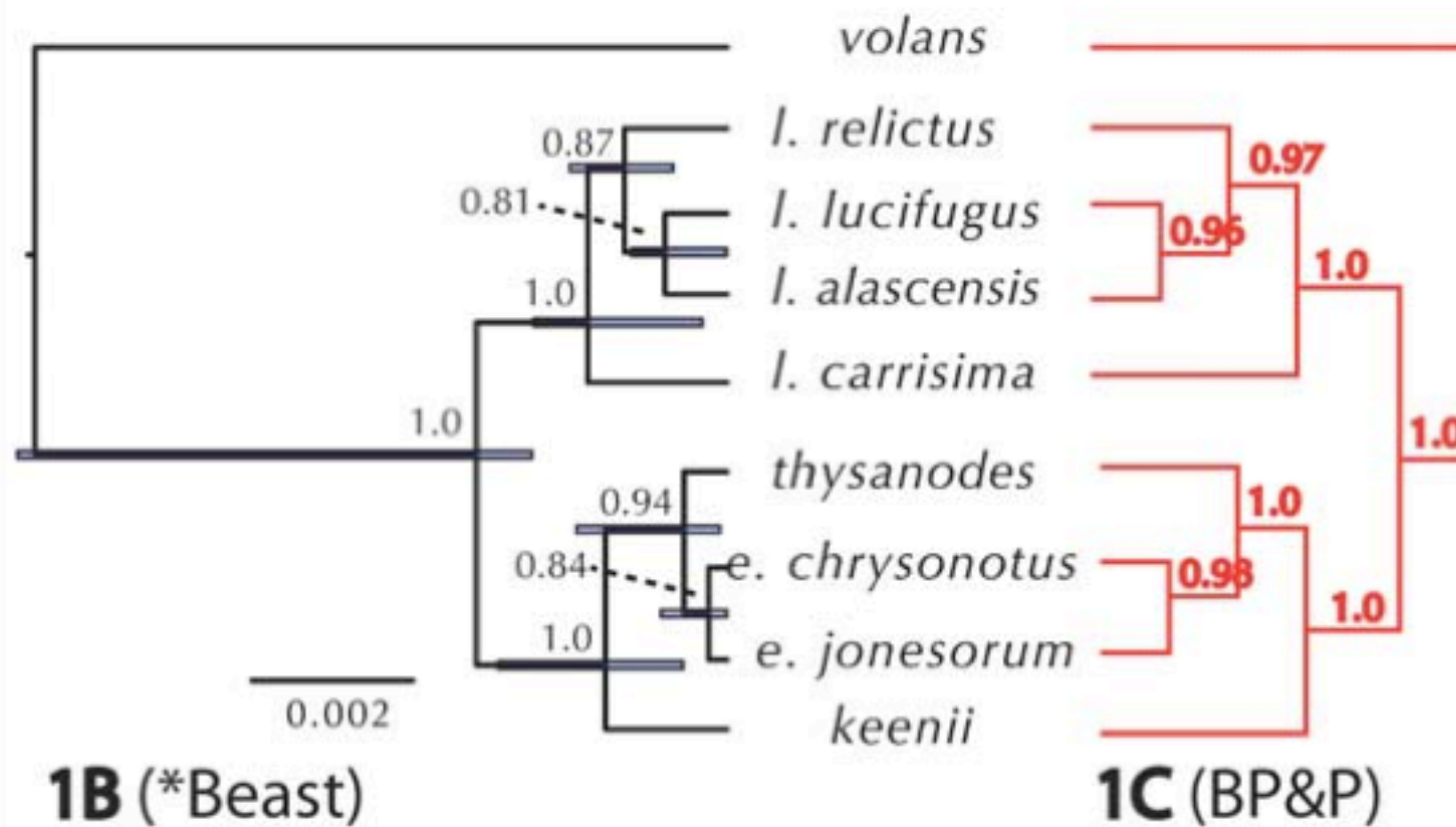
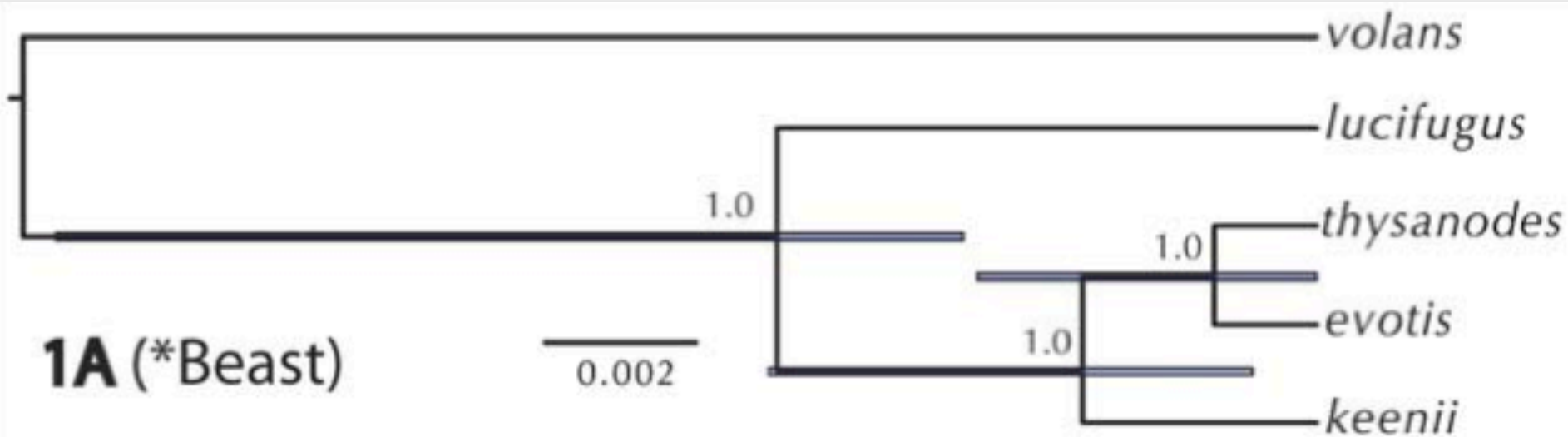




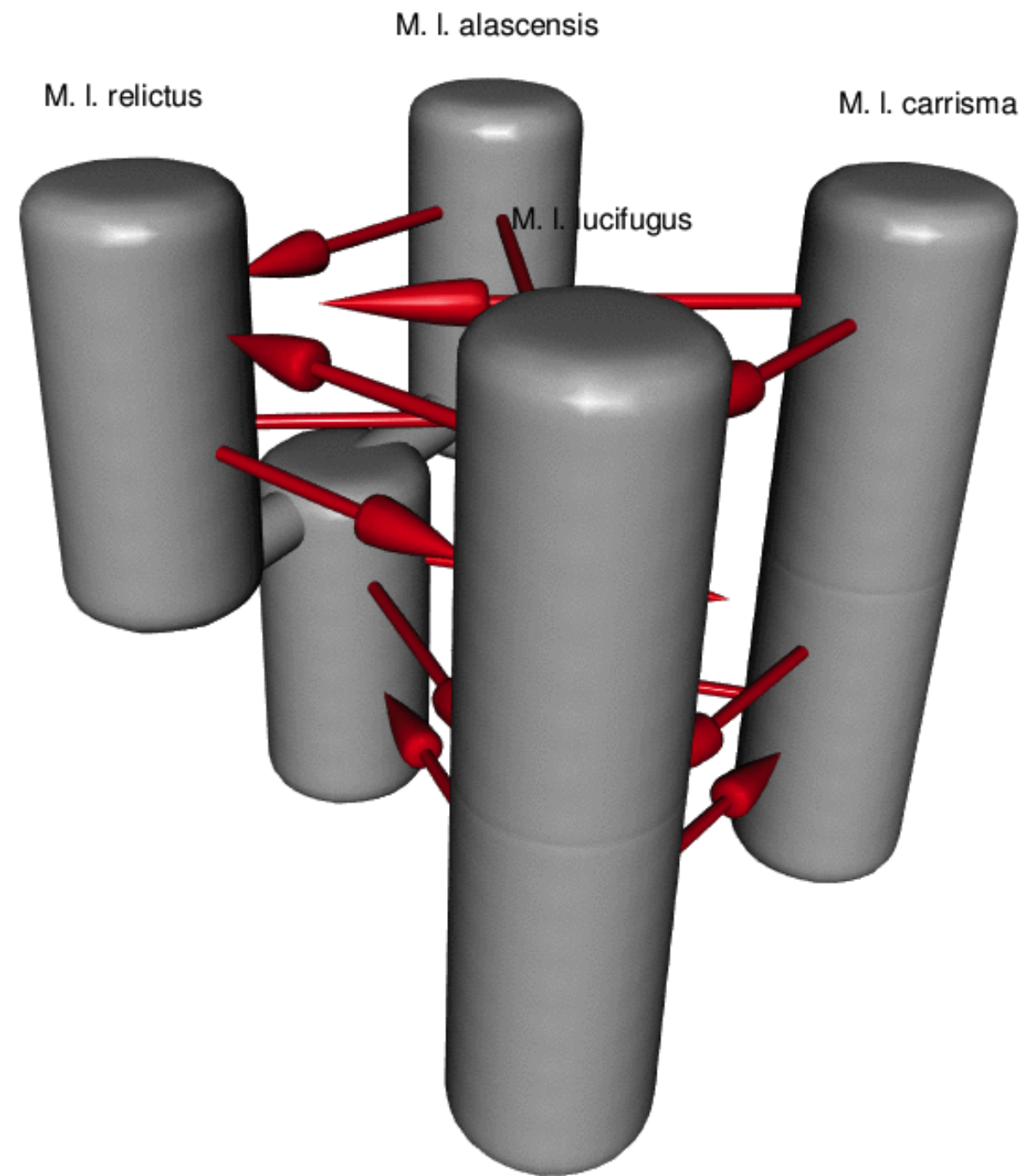
ΔAIC	AIC weight	Model
0.0	0.954	
9.5	0.008	
9.6	0.008	
10.6	0.005	
135 other models		
261.8	0.000	



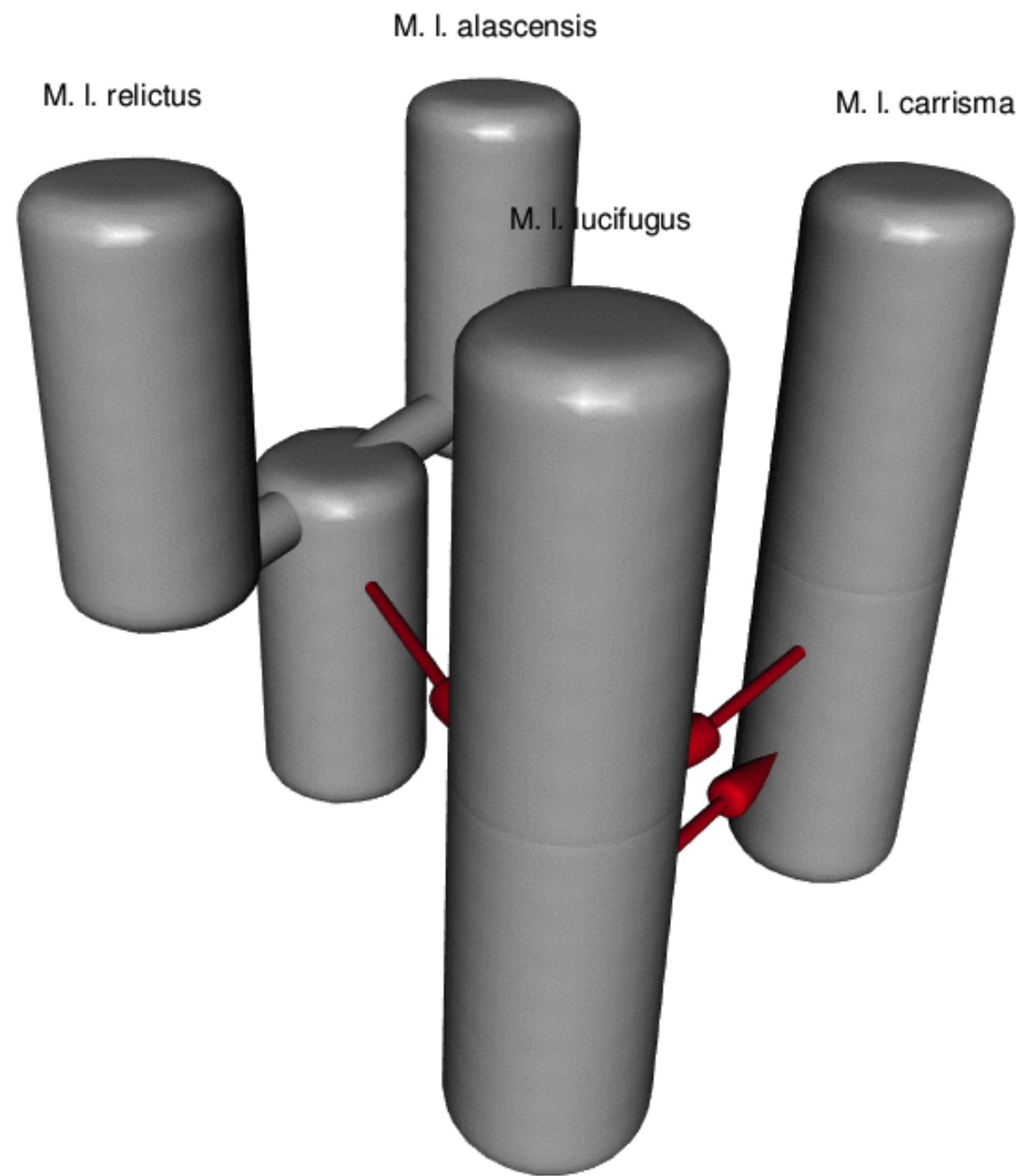
Kent McFarland



$$\Delta AIC = 0$$



$$\Delta AIC = 5.6$$

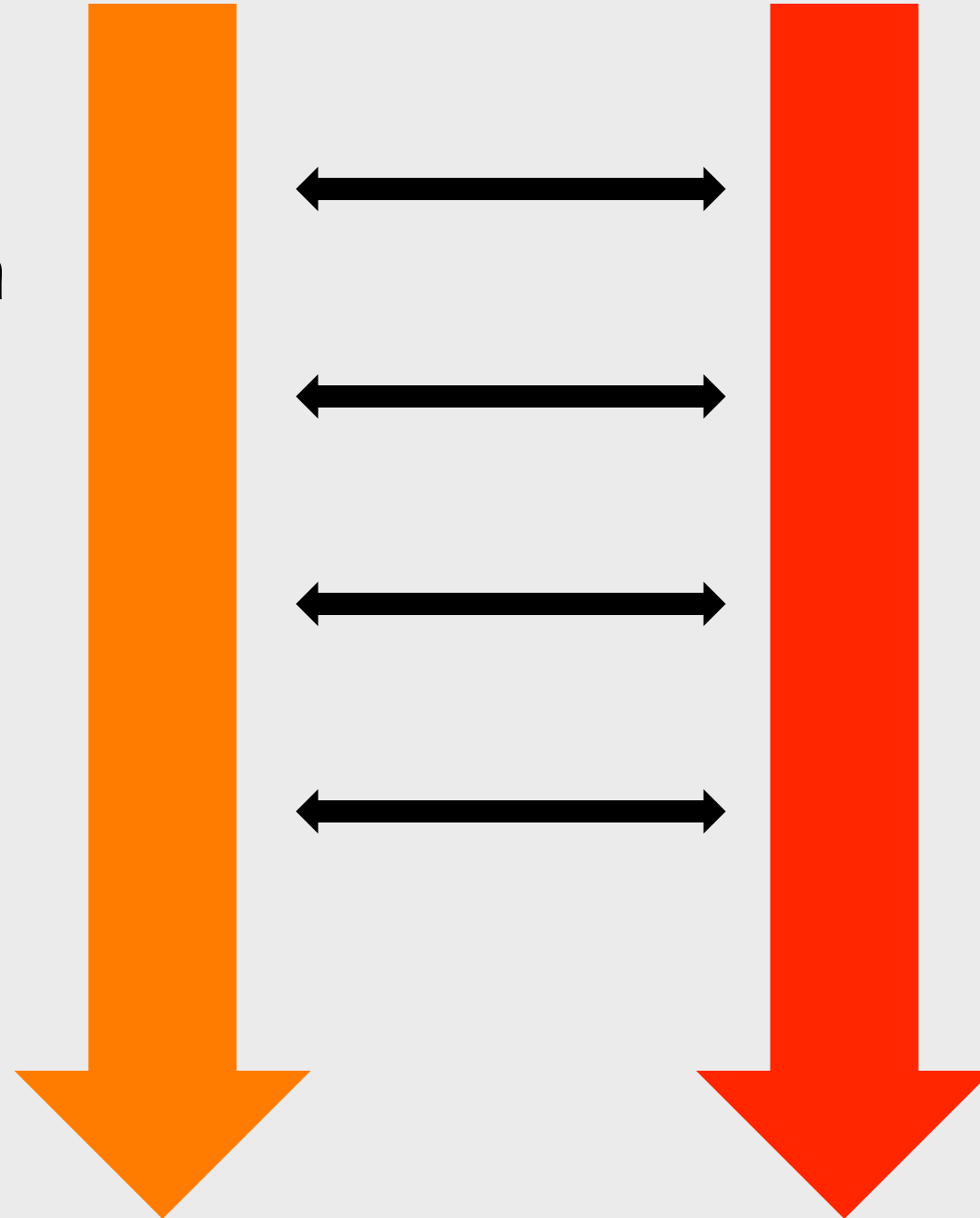


Summary

- Way to create many models
- Evaluate many models
- Prior free, but complementary with other approaches
- Should scale linearly with gene number

Present/Future

Coding
Brian O'Meara
Nathan Jackson



Empirical example,
use, documentation
Bryan Carstens
Ariana Morales

Present/Future

1. Better code (changing pop size, better algorithms), incl. more testing and comparison.
2. ~~Parameter estimation~~ (done)
3. Additional *Myotis* samples
4. Shortcuts (equivalent to ModelTest using a single NJ tree rather than tree search)
5. Sensitivity to tree uncertainty
6. Statistical phylogeography workshops (2014 & 2015)

Acknowledgements



1257784
1257699



Note I am **looking for grad students**. Phylogeography, character evolution, diversification, etc. Five years guaranteed support.

End