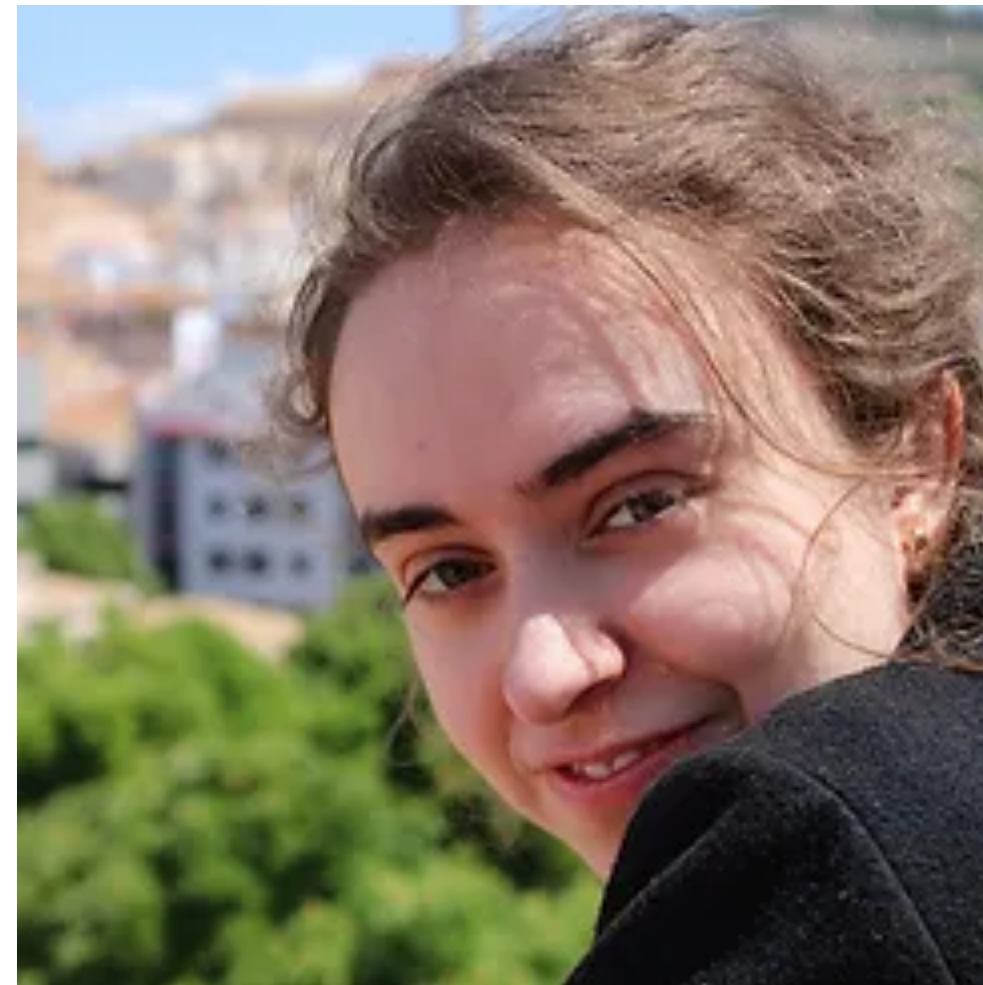
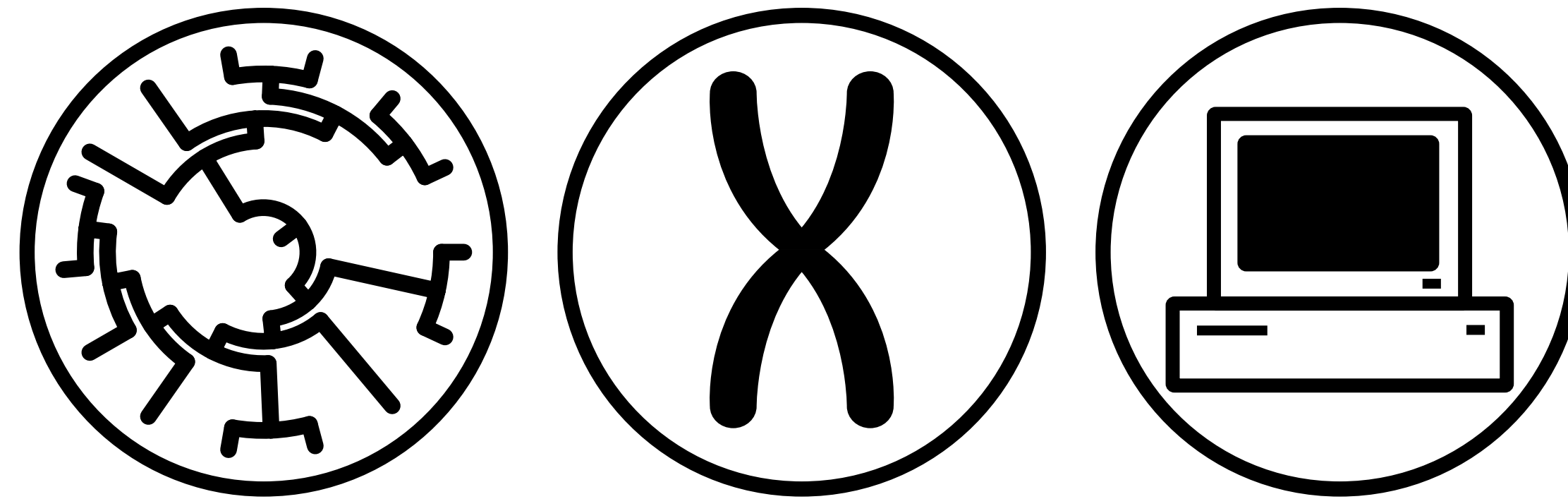


Many thanks to Gemma!



Outline



- Introduction
- Inferring genetic networks from phylogenies
- Phylogenomic subsampling
- Misc. notes before the tutorial



Jacob

Howard

Mercy

Emily

Josh

Nina

Gary

Art for Earth



Purchase with purpose: 100% of profits go toward global conservation efforts

Shop by product type, conservation status, or buy a sticker of [Sciart logo!](#)



Vinyl Stickers



Poster Prints



Camper Mugs



Endangered Animals

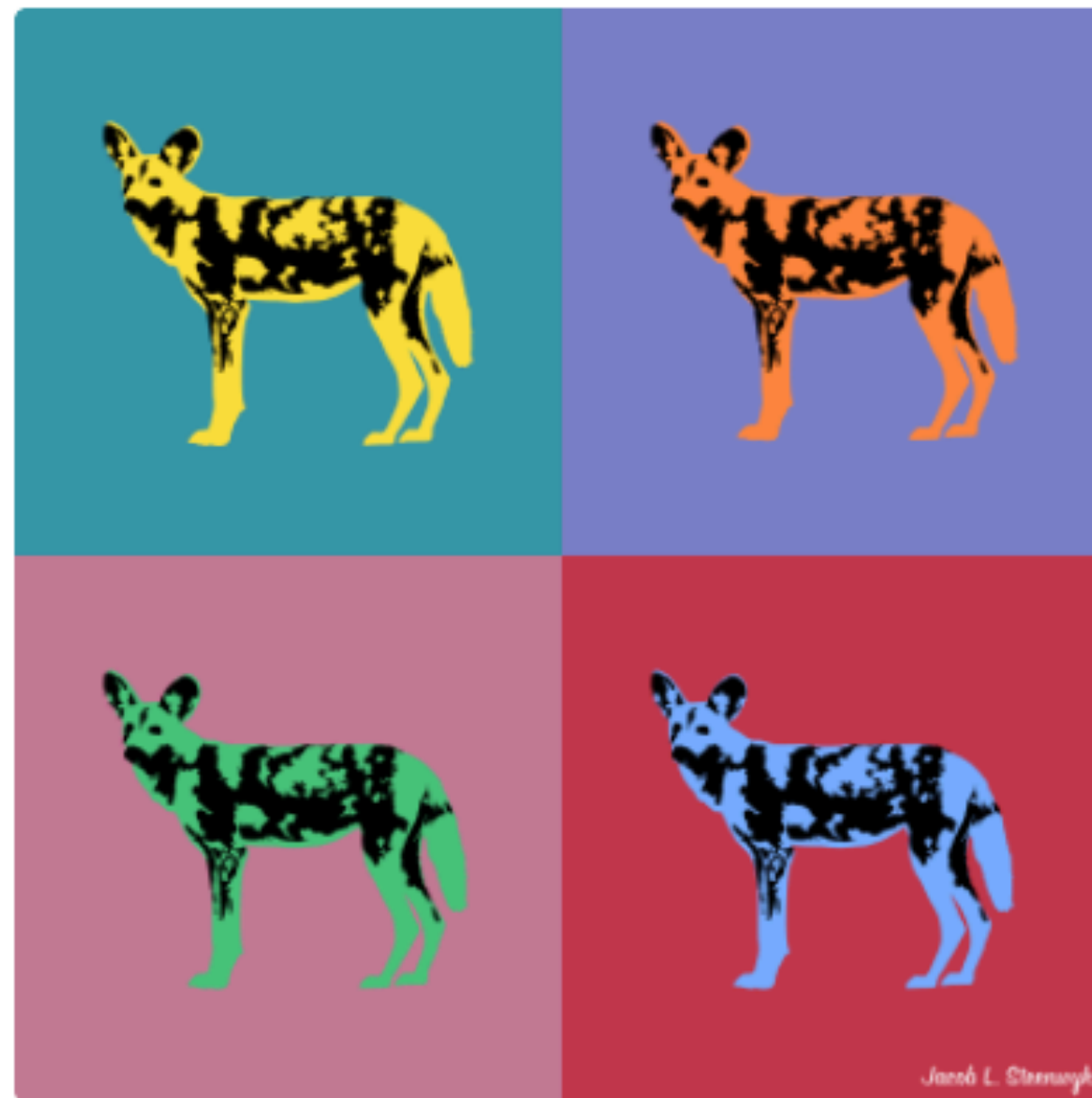


No current concern

Have a question? Check out the [Frequently Asked Questions \(FAQ\)](#) section or get in touch via [twitter!](#)

Art for Earth

Using art to raise awareness and immortalize endangered species



African wild dog (*Lycaon pictus*)

- Status: Endangered
- Population: 1,409



Blue whale (*Balaenoptera musculus*)

- Status: Endangered
- Population: 10,000 - 25,000



Galápagos penguin (*Spheniscus mendiculus*)

- Status: Endangered
- Population: fewer than 2,000

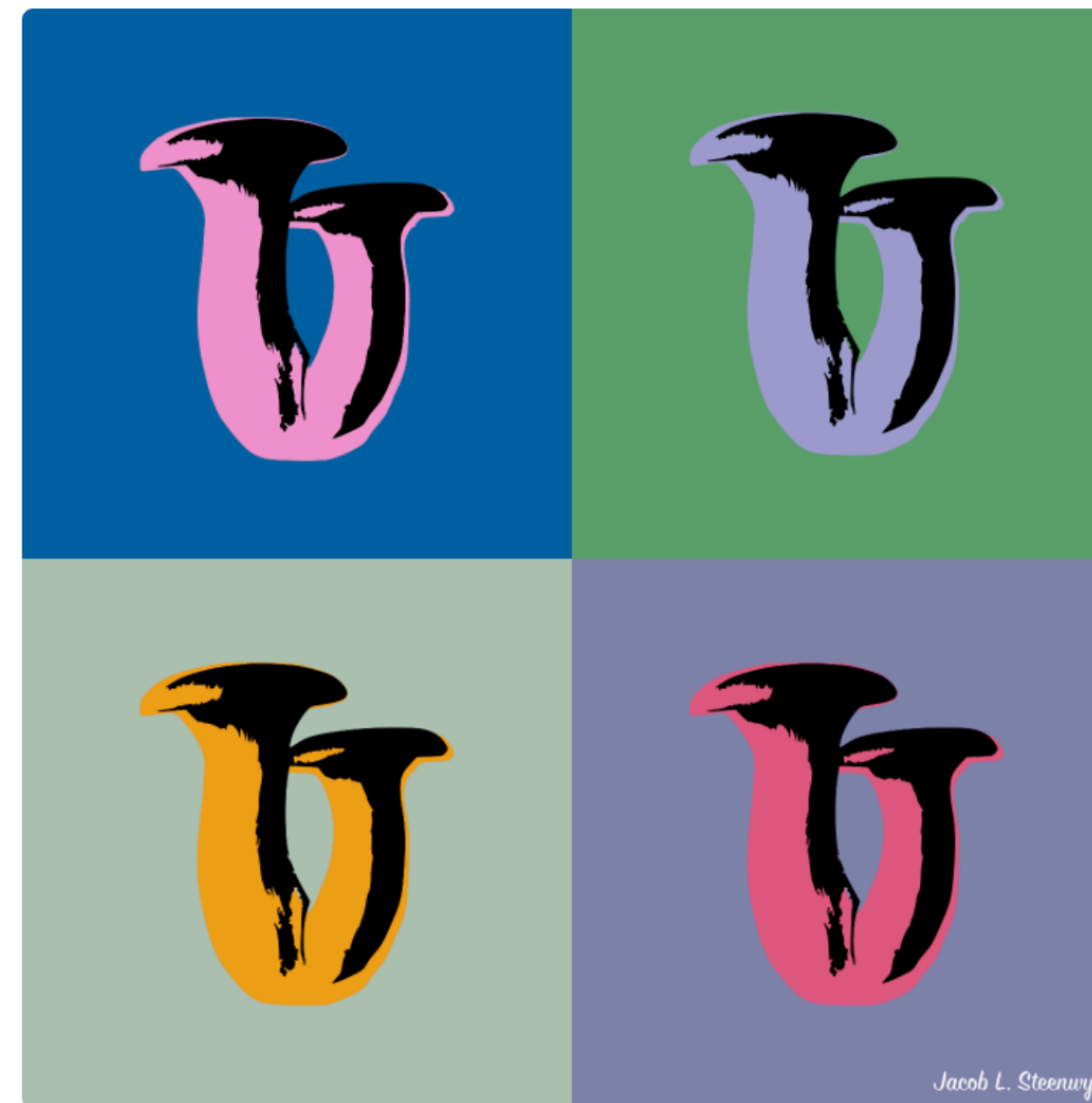
Art for Earth

Using art to raise awareness and immortalize endangered species or species I love



Fly agaric (*Amanita muscaria*)

- \$8.99 vinyl sticker (FREE shipping)
- \$22.99+ poster print (FREE shipping)
- \$28.99 camper mug (FREE shipping)



Oyster mushroom (*Pleurotus ostreatus*)

- \$8.99 vinyl sticker (FREE shipping)
- \$22.99+ poster print (FREE shipping)
- \$28.99 camper mug (FREE shipping)



Morel mushroom (*Morchella esculenta*)

- \$8.99 vinyl sticker (FREE shipping)
- \$22.99+ poster print (FREE shipping)
- \$28.99 camper mug (FREE shipping)

Featured on *Yeast* magazine

Received: 29 July 2020 | Accepted: 20 August 2020
DOI: 10.1002/yea.3518

SPECIAL ISSUE ARTICLE

Yeast WILEY

A portrait of budding yeasts: A symbol of the arts, sciences and a whole greater than the sum of its parts

Jacob L. Steenwyk^{1,2} 

¹Department of Biological Sciences, Vanderbilt University, Nashville, TN, USA

²Early Career Leadership Program Communication and Outreach Subcommittee, Genetics Society of America, Rockville, MD, USA

Correspondence

Jacob L. Steenwyk, Department of Biological Sciences, Vanderbilt University, Nashville, TN 37235, USA.

Email: jacob.steenwyk@vanderbilt.edu

Funding information

Vanderbilt University; Howard Hughes Medical Institute

KEYWORDS: art, budding yeast, cell cycle, Merian, naturalist, non-conventional yeasts, sciart, science, STEAM, Warhol

1 | INTRODUCTION

In the year 1660, 13-year-old Maria Sibylla Merian roamed the gardens and countryside of Germany taking detailed notes about caterpillars, moths, butterflies and their interactions with host plants, accompanying her notes were elaborate multimedia depictions of insect and plant life cycles (Figure 1). Merian's efforts in documenting interspecies relationships are regarded as early contributions to modern natural history and ecology, although the term 'ecology' was coined approximately two centuries later (Etheridge, 2011a, 2011b; Pieters & Winthagen, 1999). Her influence can be seen in the work of naturalists such as John James Audubon (Etheridge, 2015; Palmeri, 2017). Merian's success in part stems from her ability to use art to bolster her science and vice versa.

Merian is one of many scientists and artists who blended the arts and sciences over the centuries. In fact, scientist-artist polymaths like Aristotle and Leonardo da Vinci were more commonplace in part because of the common goal science and art share: interpreting and representing the natural world. The 'great divide' of the arts and sciences in Western cultures is thought to have started in the 19th century, coinciding with the term 'scientist' being coined (Braund & Reiss, 2019; Sumner, 1959; Zhu & Goyal, 2019). The division became reinforced. Schools for arts and sciences were separated as unfounded claims about brain differentiation formulated (Zhu & Goyal, 2019). For example, the right and left brain hemispheres were thought to be individually responsible for arts and science learning, respectively (Sperry, 1968). However, evidence from cognitive scientists favours a holistic view of the brain wherein a wide range of stimulation (e.g., arts and sciences) improves broad brain function and critical thinking skills (Braund & Reiss, 2019; Howes, Kaneva, Swanson, & Williams, 2013).

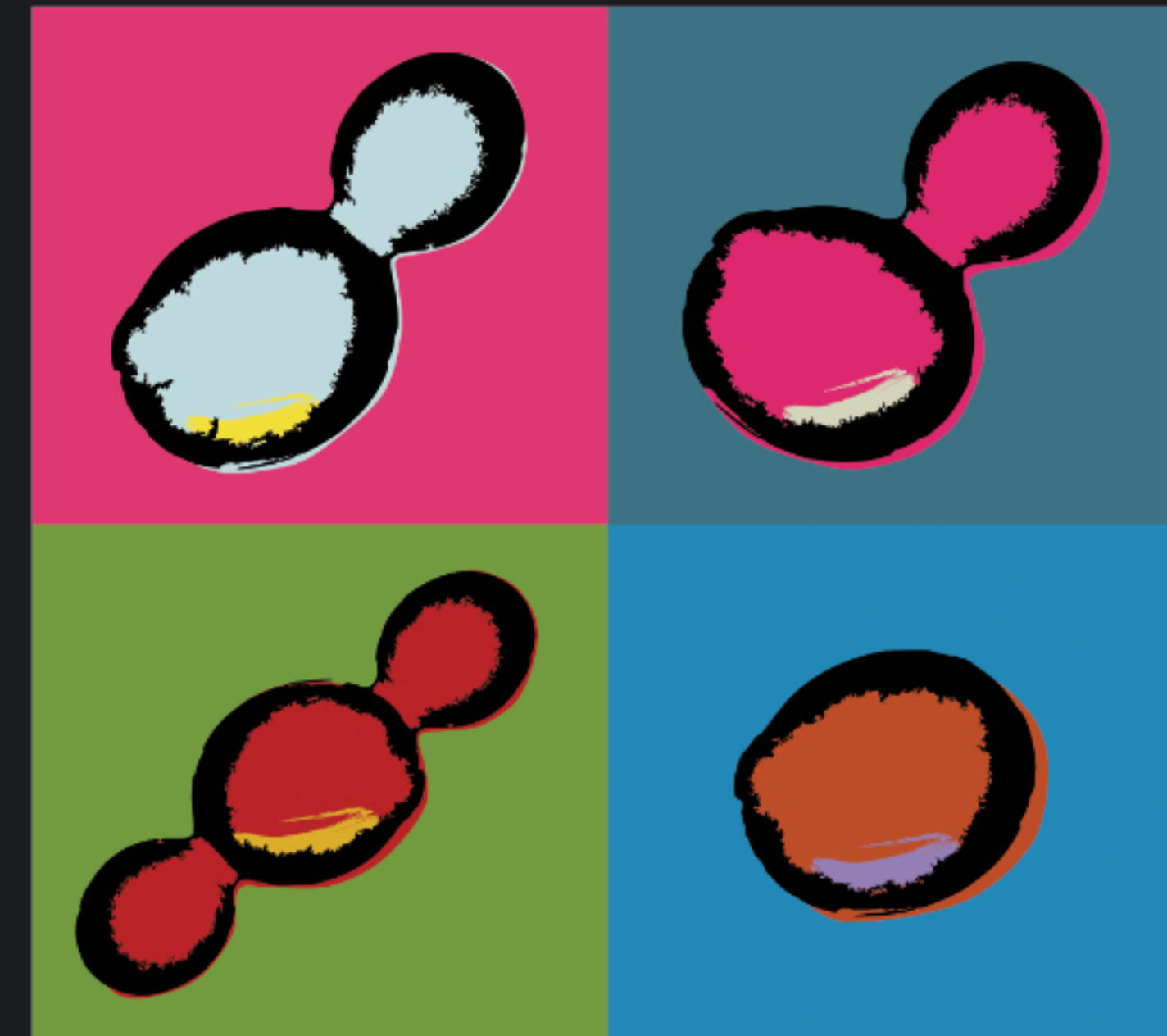
Today, the benefits of a holistic view of the arts and the sciences have been recognized by numerous institutions. For example, Science, Technology, Engineering, Arts and Mathematics (STEAM) inspired curriculum is used to help students build skills for broad problem solving in K-12 schools (Kim & Park, 2012; Pepler, 2013; Sochacka, Guyotte, & Walther, 2016). In higher education, artists, designers, researchers and inventors have formed forward-thinking coalitions such as the Center for Art, Science & Technology at Massachusetts Institute of Technology (<https://arts.mit.edu/cast/>) and ArtLab at Vanderbilt University (<https://artlabvanderbilt.com/>) to reunite the arts and sciences. These initiatives and many others have used the arts as an effective form of communication between scientists and the broader community (Illingworth, 2017), ultimately helping disseminate major scientific findings across society.

Perhaps one of the most important and recent scientific findings in the field of biological sciences is our understanding of the cellular life cycle. Seminal discoveries that unraveled the controls of the life cycle were made studying the model unipolar budding yeast *Saccharomyces cerevisiae* (Hartwell, Culotti, Pringle, & Reid, 1974). Comparative studies of *S. cerevisiae*, the fission yeast (*Schizosaccharomyces pombe*) and animals revealed striking similarities suggesting the life cycle is evolutionarily stable (Breedon & Nasmyth, 1987). Exploiting these similarities has enabled yeasts to be powerful models for cancer biology research and the development of anticancer therapeutics (Gao, Chen, & Huang, 2014; Guaragnella et al., 2014; Schwartz & Dickson, 2009). However, examination of non-conventional yeasts and their life cycles can provide novel insights important to the fields of cell biology, evolutionary biology and more. For example, species of the budding yeast genus *Hanseniaspora* have lost numerous cell cycle control genes, including *MAD1*, *MAD2* and *RAD9*, and components of the Anaphase Promoting Complex and display atypical bipolar budding patterns (Steenwyk

January 2021, Volume 38, Issue No. 1

Yeast

ISSN 0749-503X



Special Issue: Exploring the Yeast Life Cycles
Edited by Nishant KT, Nobile C, Wloch-Salamon D and Wolfe K
Dedicated to the memory of Angelika Amon

WILEY

Lineages of interest across my career



Lineages of interest across my career



~18,000
genomes



~15,000
genomes

Lineages of interest across my career



~18,000
genomes

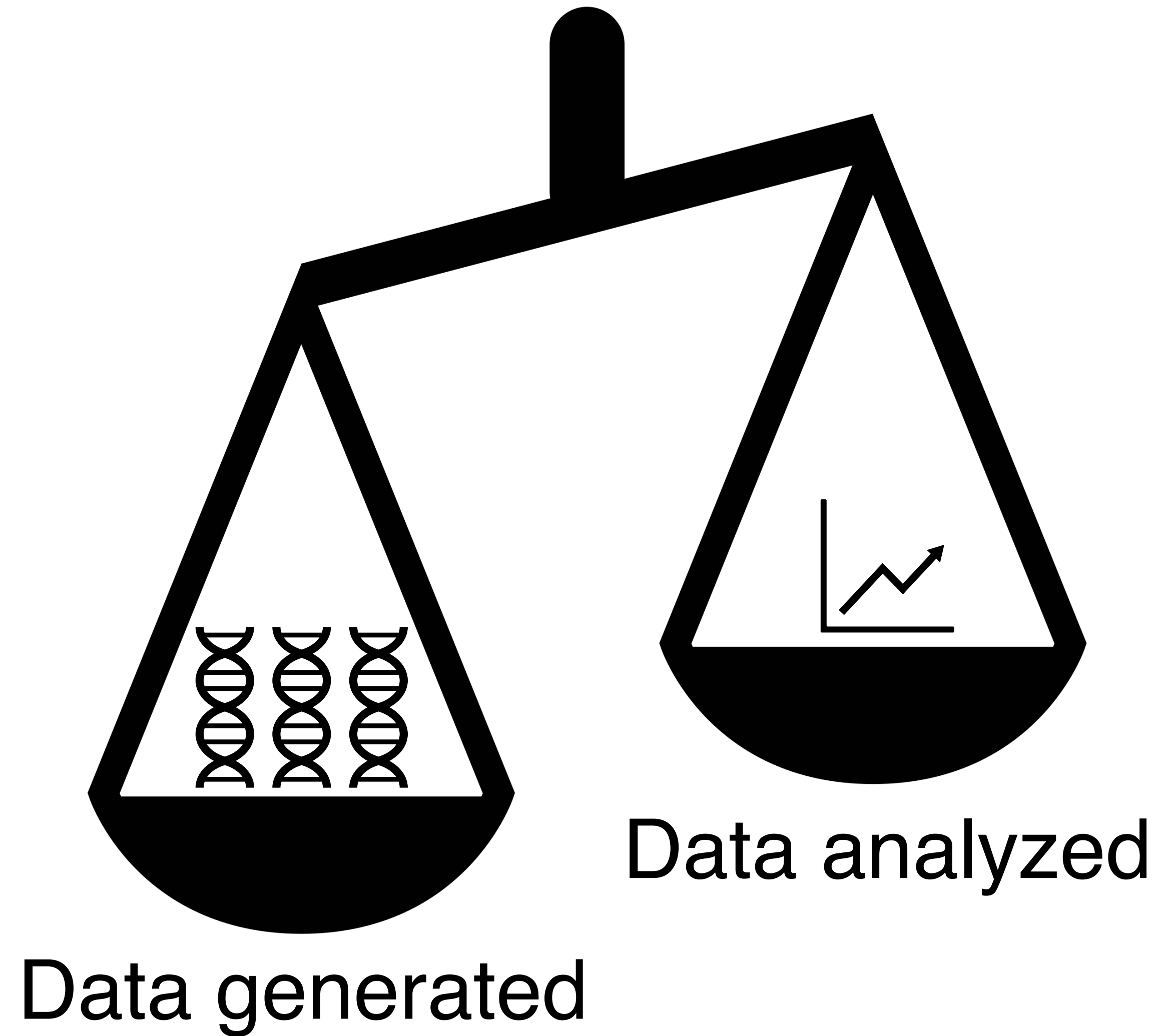


~15,000
genomes



~4,000
plant genomes

Data generation has outpaced data analysis



Engineering software for 'omic inquiry

Ortholog identification

Ortho**SNAP**

Steenwyk et al. (2022),
PLOS Biology

orthofisher

Steenwyk & Rokas (2021),
G3 Genes|Genomes|Genetics

Phylogenomics

Clip**KIT**

Steenwyk et al. (2020),
PLOS Biology

Phy**KIT**

Steenwyk et al. (2021),
Bioinformatics

Remar**KIT**

Steenwyk et al. (in prep.)

Genomics

Bio**KIT**

Steenwyk et al. (2022),
Genetics

LVBRS

Le and Steenwyk et al. (2022),
bioRxiv

Other

OVER
UNDER

Steenwyk et al. (in prep.)

ggpubfigs

Steenwyk & Rokas (2021),
Micro. Resource Announcements

treehouse

Steenwyk & Rokas (2019),
BMC Research Notes

 @JLSteenwyk

Engineering software for 'omic inquiry

Ortholog identification

Ortho**SNAP**

Steenwyk et al. (2022),
PLOS Biology

orthofisher

Steenwyk & Rokas (2021),
G3 Genes|Genomes|Genetics

Phylogenomics

Clip**KIT**

Steenwyk et al. (2020),
PLOS Biology

Phy**KIT**

Steenwyk et al. (2021),
Bioinformatics

Remar**KIT**

Steenwyk et al. (in prep.)

Genomics

Bio**KIT**

Steenwyk et al. (2022),
Genetics

LVBRS

Le and Steenwyk et al. (2022),
bioRxiv

Other

OVER
UNDER

Steenwyk et al. (in prep.)

ggpubfigs

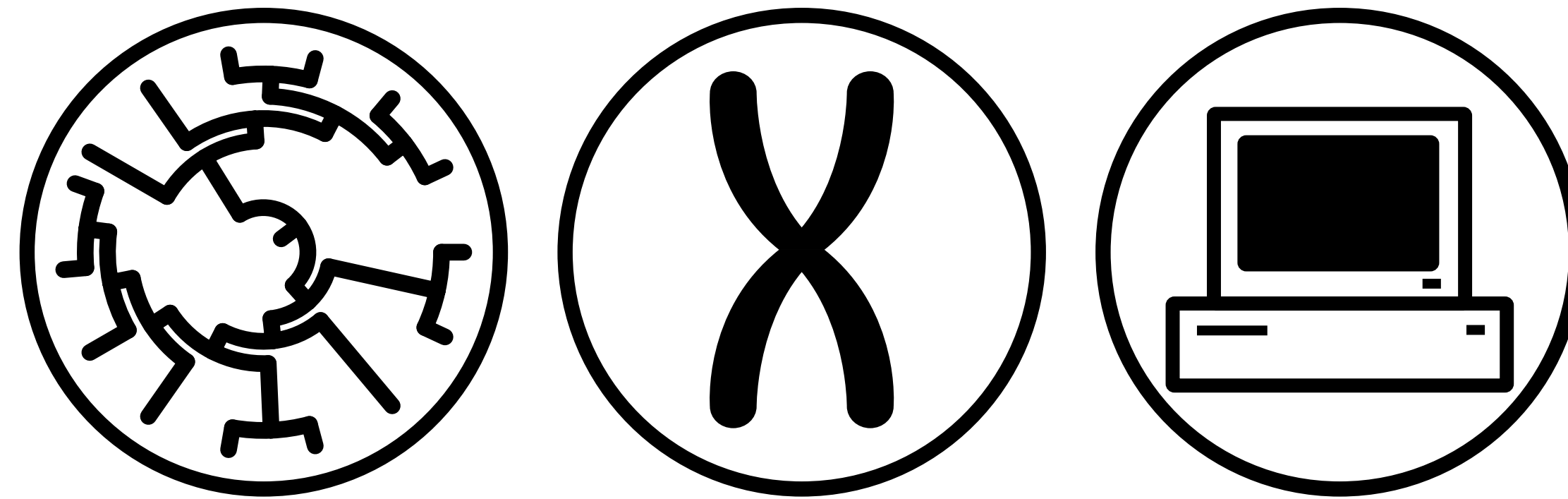
Steenwyk & Rokas (2021),
Micro. Resource Announcements

treehouse

Steenwyk & Rokas (2019),
BMC Research Notes

 @JLSteenwyk

Outline



- Introduction
- **Inferring genetic networks from phylogenies**
- Phylogenomic subsampling
- Misc. notes before the tutorial

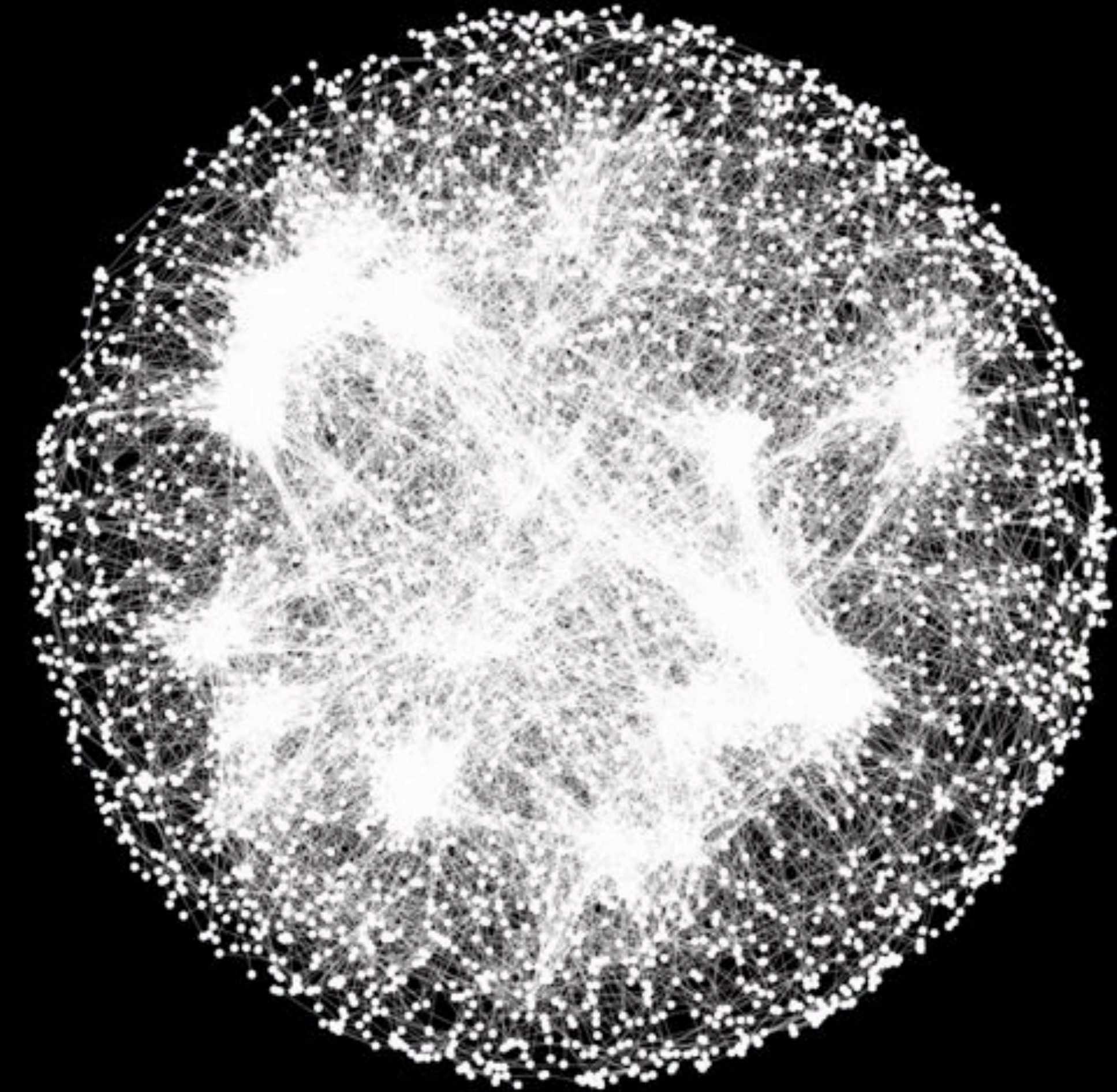
Inferring genetic networks from phylogenies

Jacob L. Steenwyk



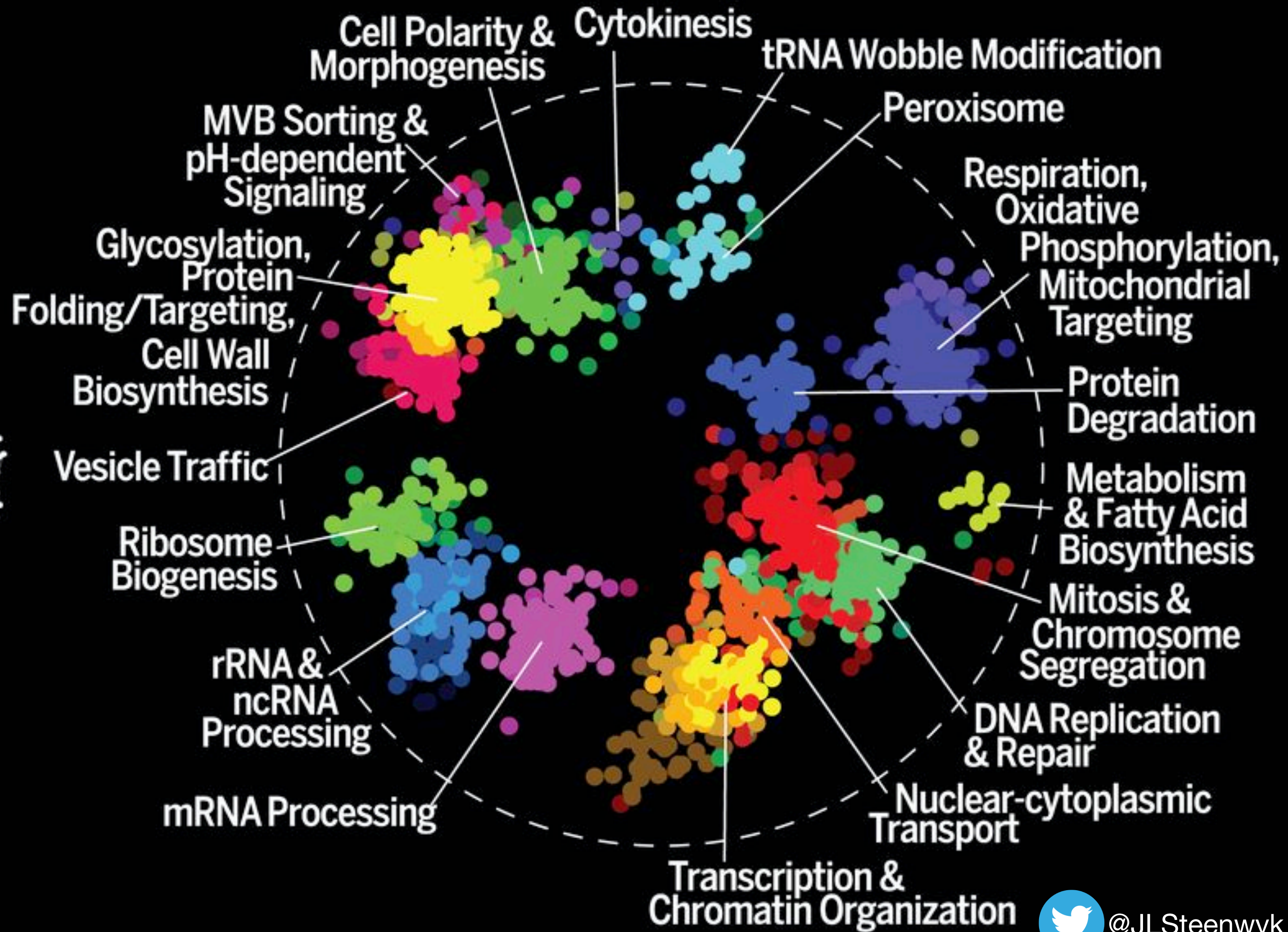
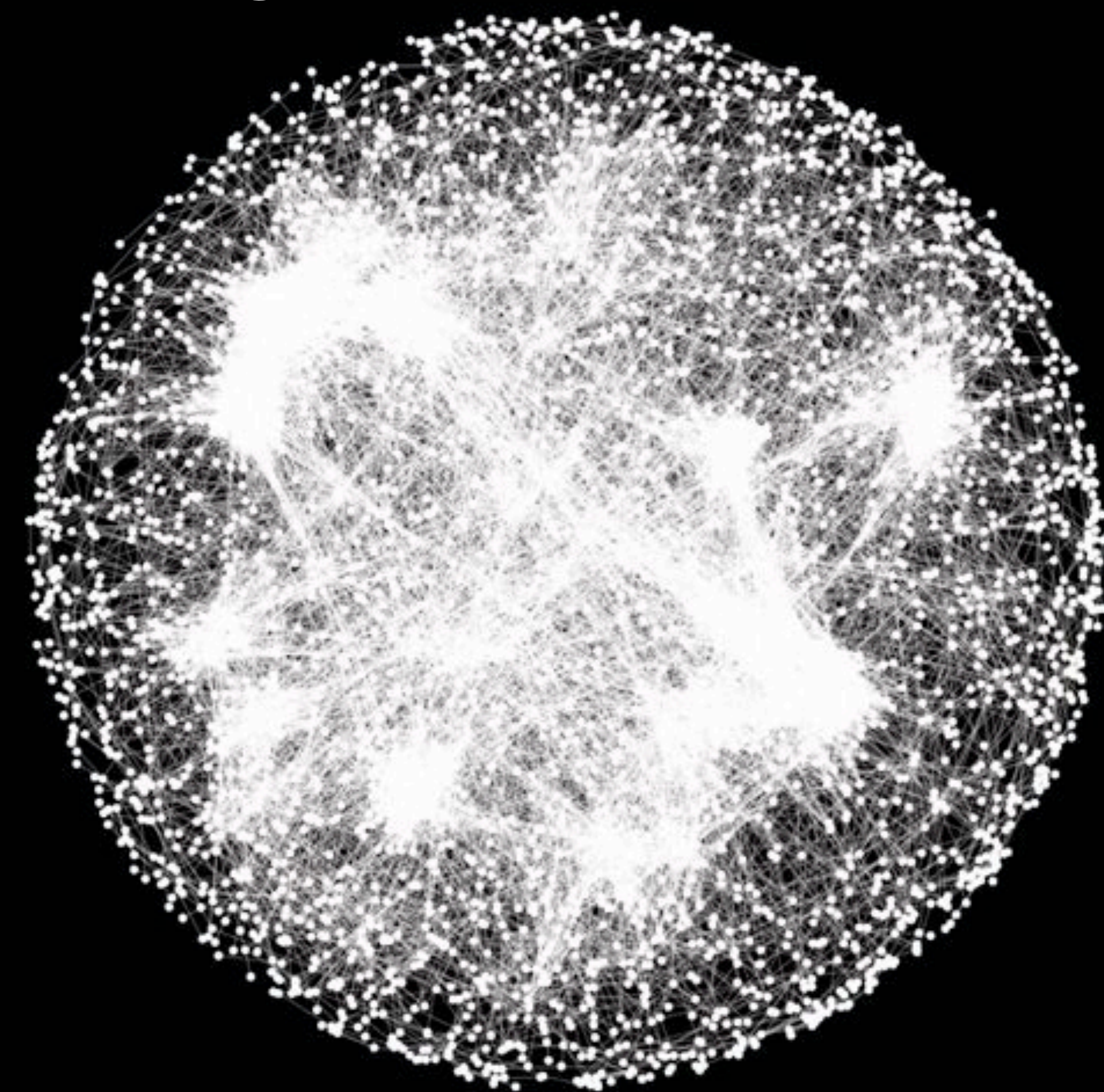
Networks capture the complexity of genomic function

 Baker's yeast



Networks capture the complexity of genomic function

 Baker's yeast



 @JLSteenwyk

PEX1 and PEX6 share function

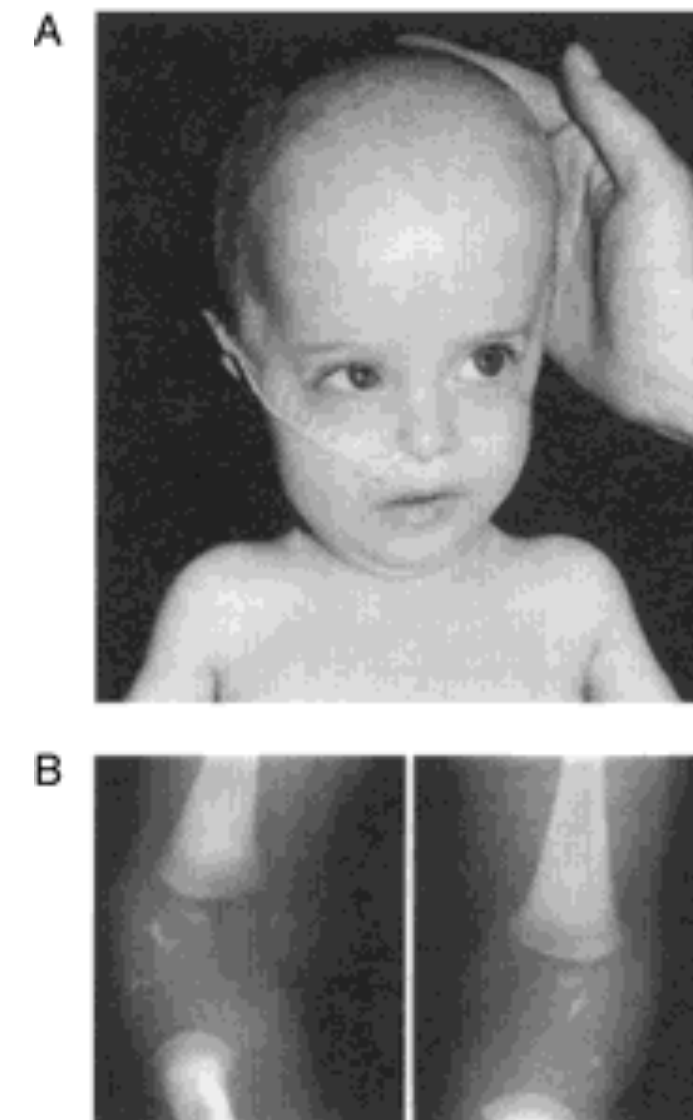
Pex1p & Pex6p: forms a heterodimer involved in recycling peroxisomal signal receptor Pex5p

PEX1 and PEX6 share function

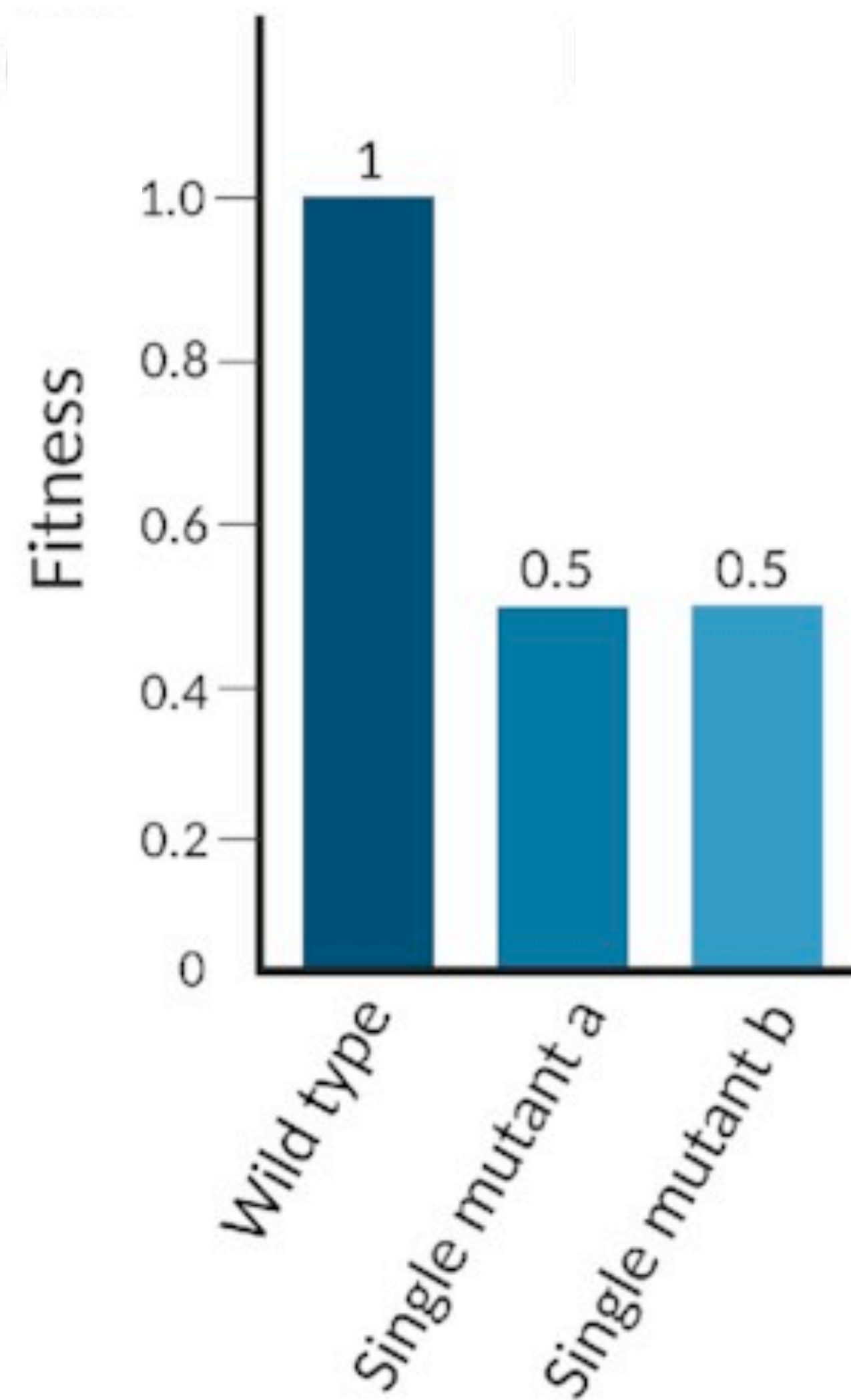
Pex1p & Pex6p: forms a heterodimer involved in recycling peroxisomal signal receptor Pex5p

Mutations that disrupt protein interactions cause neurologic disorders including:

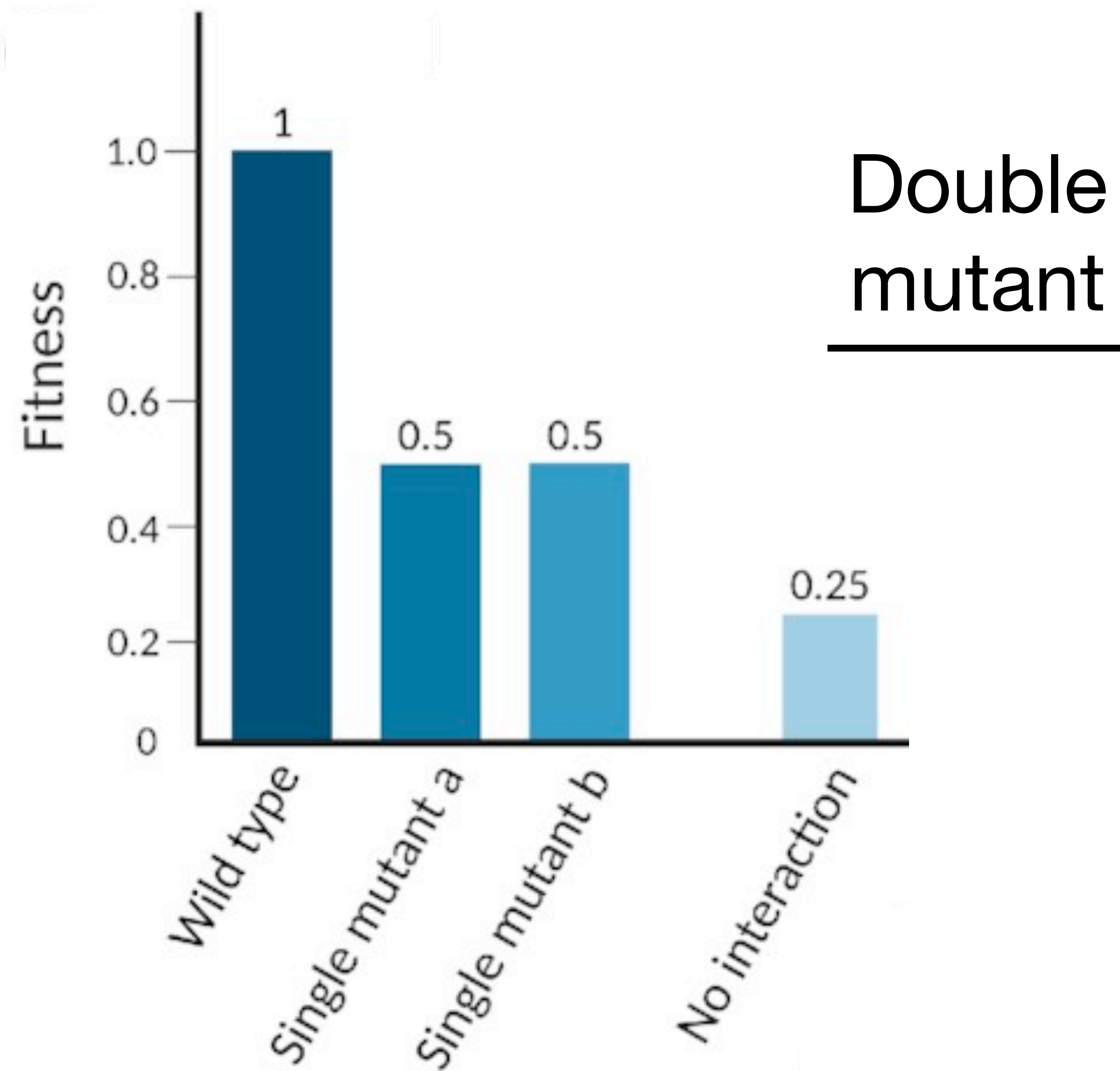
- Zellweger syndrome,
- neonatal adrenoleukodystrophy,
- infantile Refsum disease



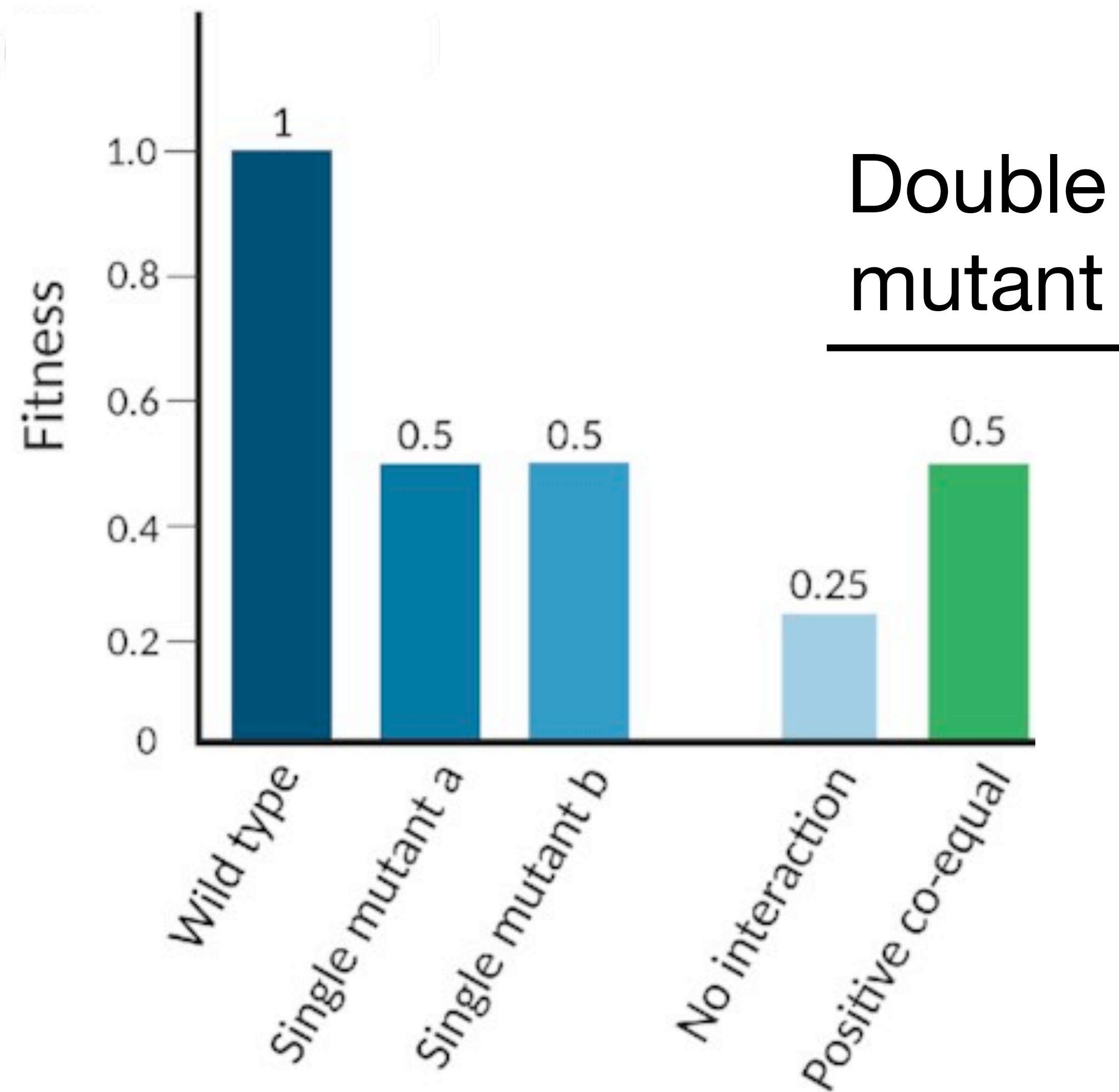
Genetic interactions are gene-gene associations



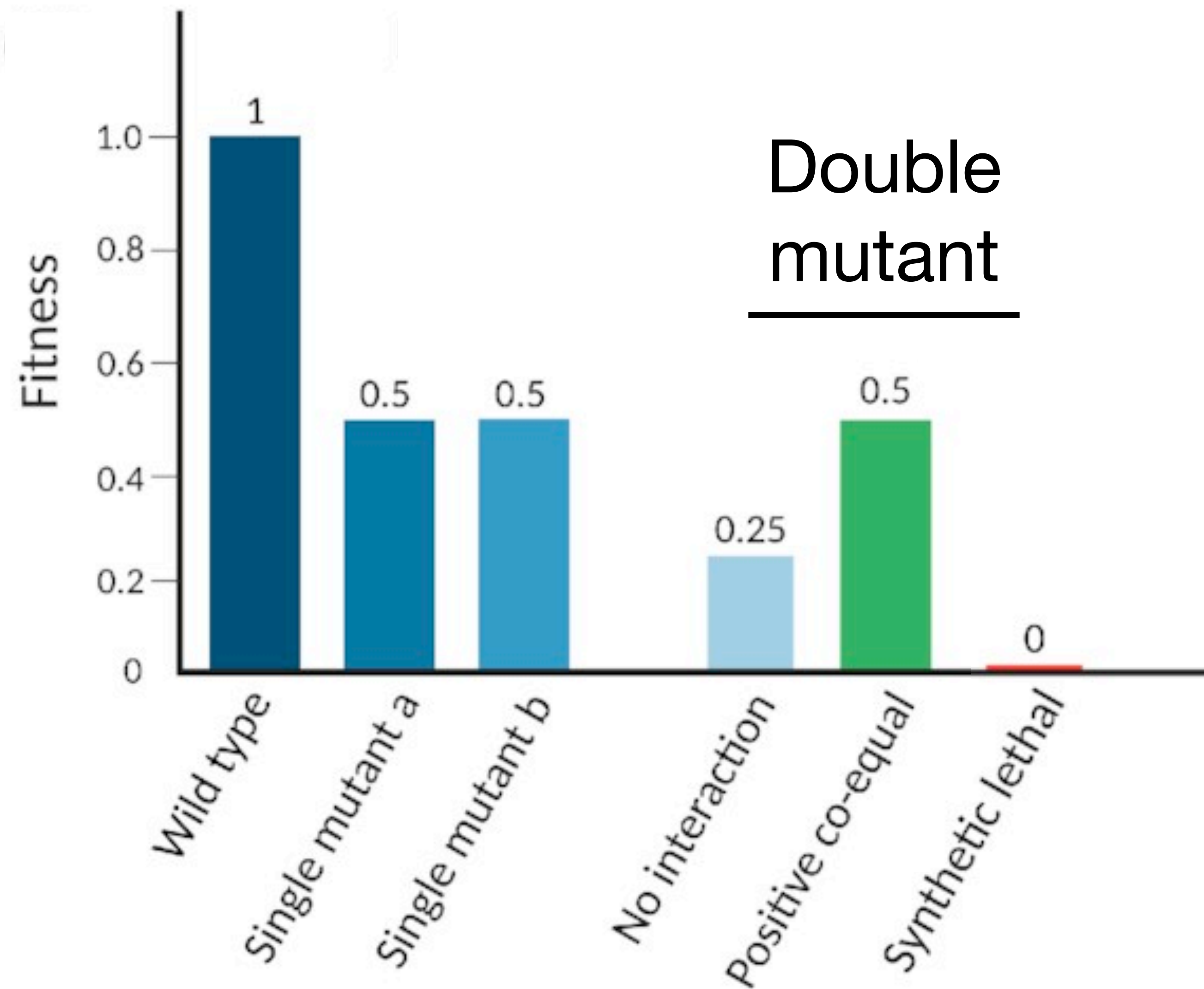
Genetic interactions are gene-gene associations



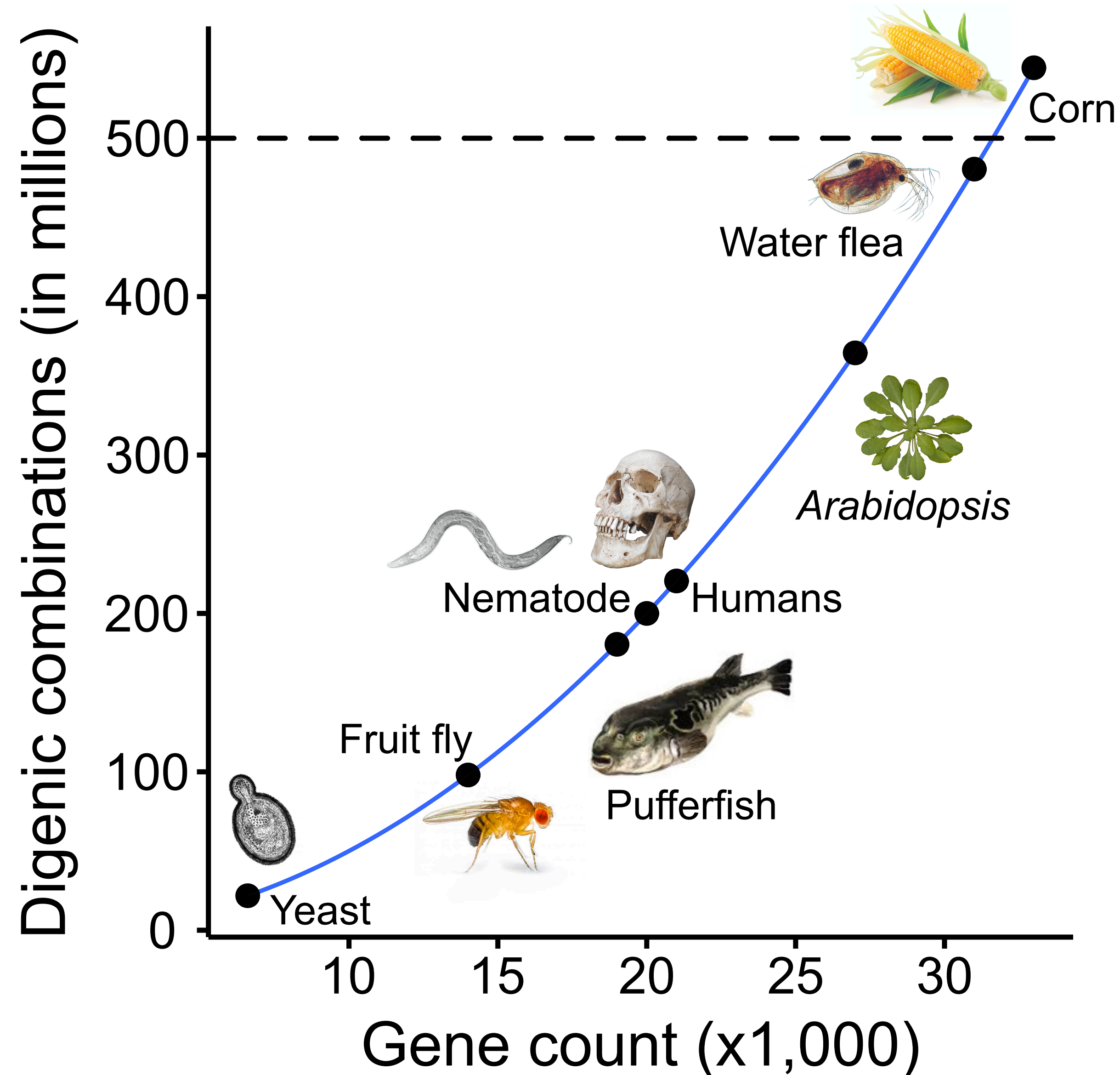
Genetic interactions are gene-gene associations



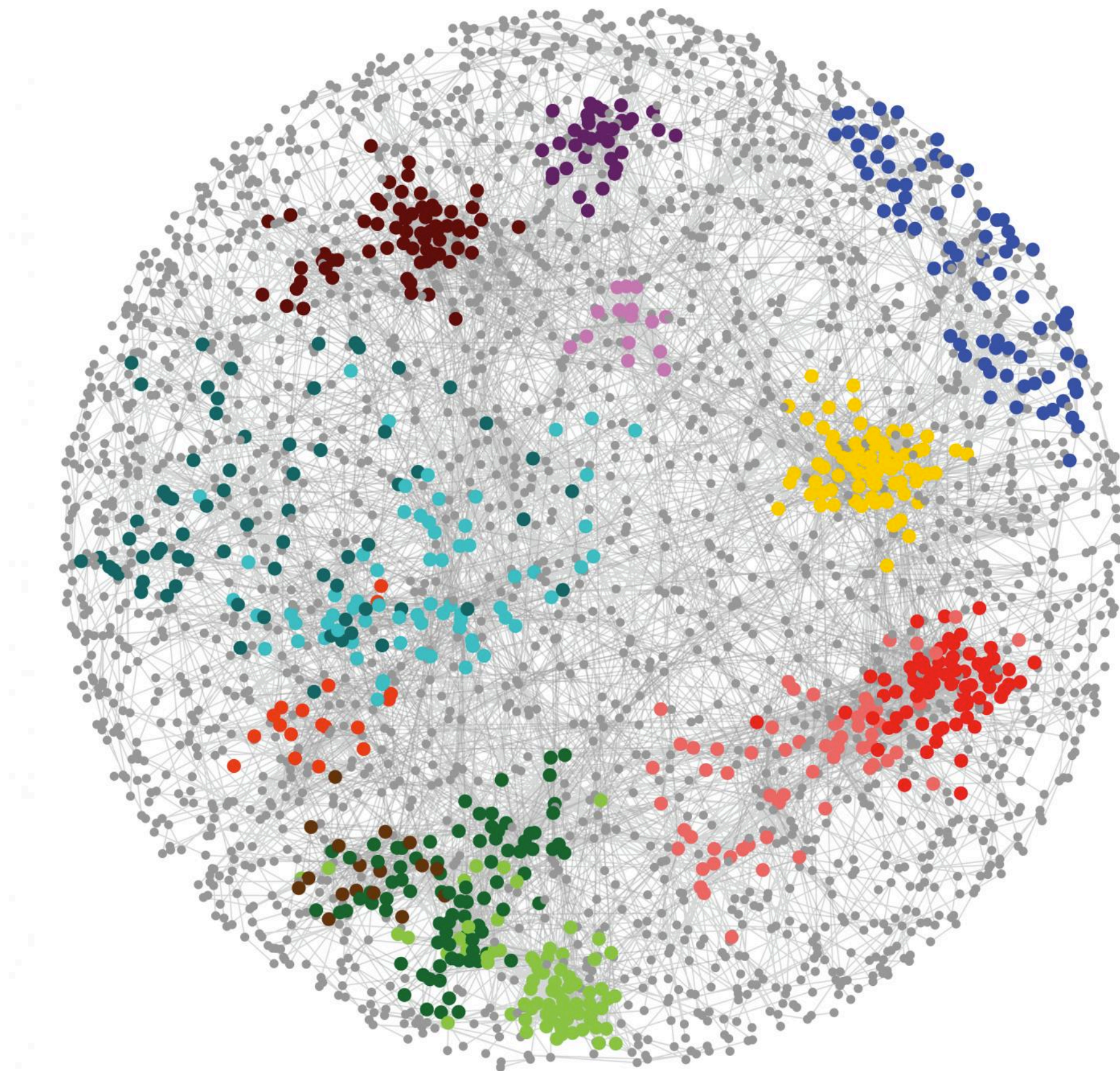
Genetic interactions are gene-gene associations



Too many pairwise combinations of genes

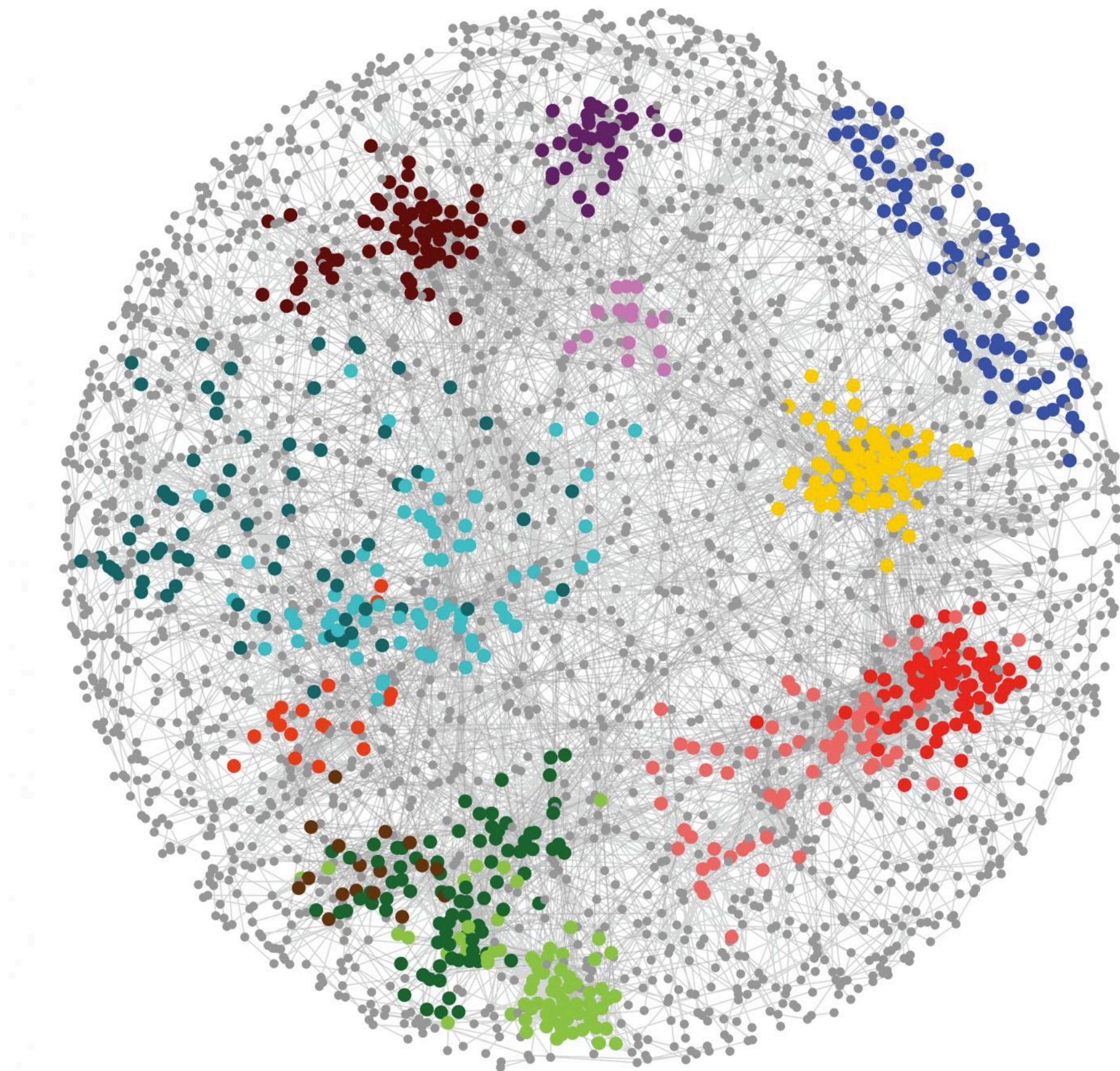


Methods to infer gene-gene associations



- Coexpression
- Gene presence/absence patterns

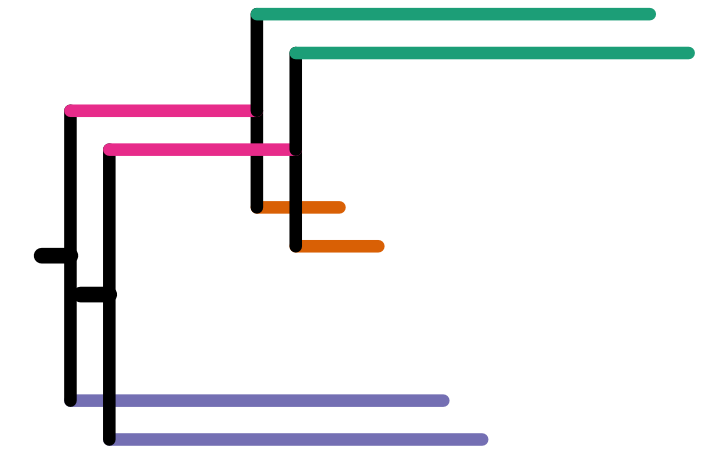
Methods to infer gene-gene associations



- Coexpression
- Gene presence/absence patterns
- **Gene coevolution**

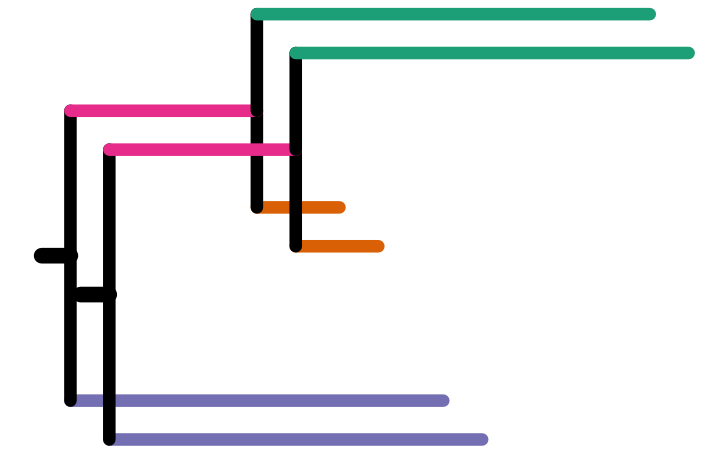
Gene-gene coevolution predicts shared function

- gene coevolution refers to:
 - two genes that covary in parallel across speciation events
 - often observed among genes that share function, are coexpressed, or are part of the same multi-meric complexes



Gene-gene coevolution predicts shared function

- gene coevolution refers to:
 - two genes that covary in parallel across speciation events
 - often observed among genes that share function, are coexpressed, or are part of the same multi-meric complexes



PhyKIT

a toolkit for examining multiple
sequence alignments and trees

PhyKIT: a broadly applicable UNIX shell
toolkit for processing and analyzing
phylogenomic data

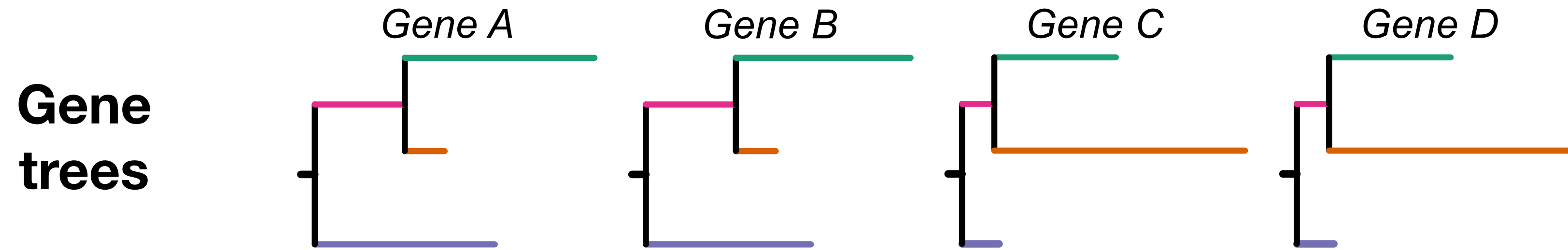
Jacob L Steenwyk ✉, Thomas J Buida, III, Abigail L Labella, Yuanning Li,
Xing-Xing Shen, Antonis Rokas ✉

Bioinformatics, btab096, <https://doi.org/10.1093/bioinformatics/btab096>

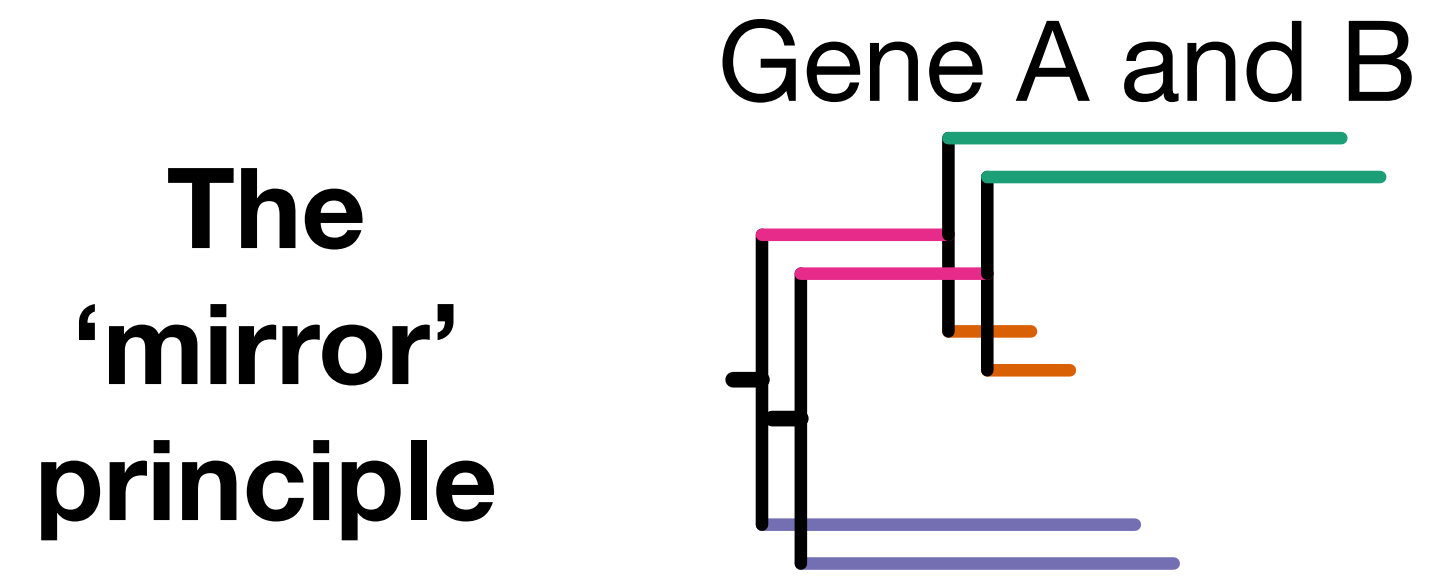
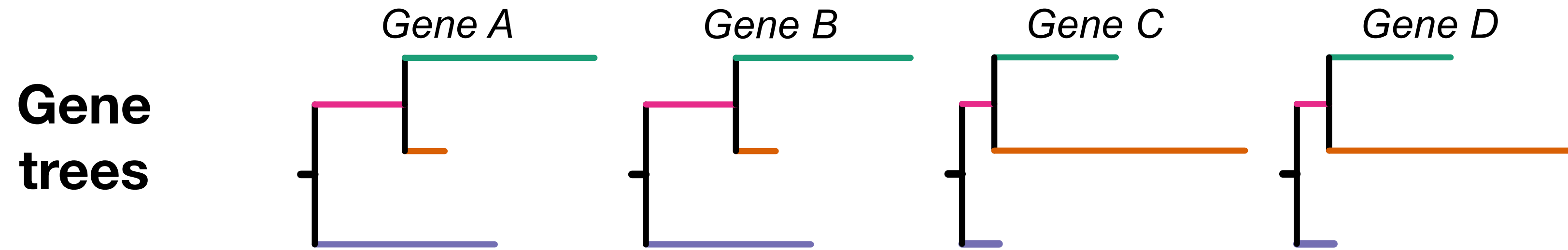
Published: 09 February 2021 [Article history](#) ▼



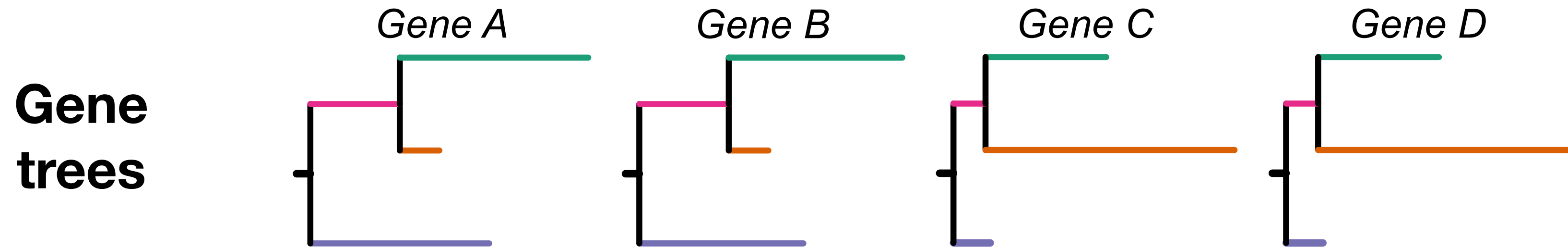
The mirror principle to detect gene coevolution



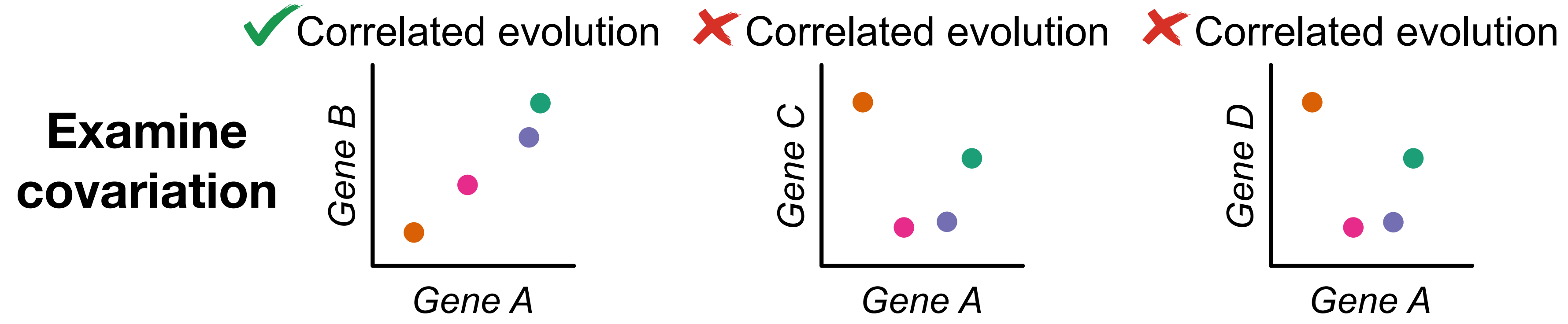
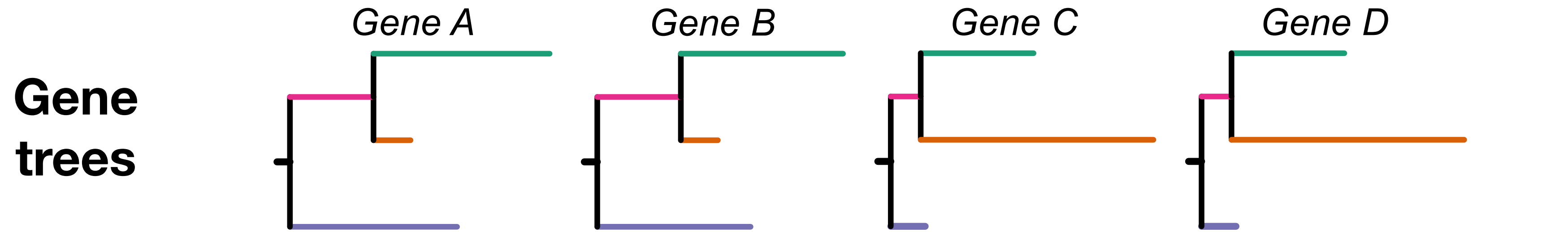
The mirror principle to detect gene coevolution



The mirror principle to detect gene coevolution

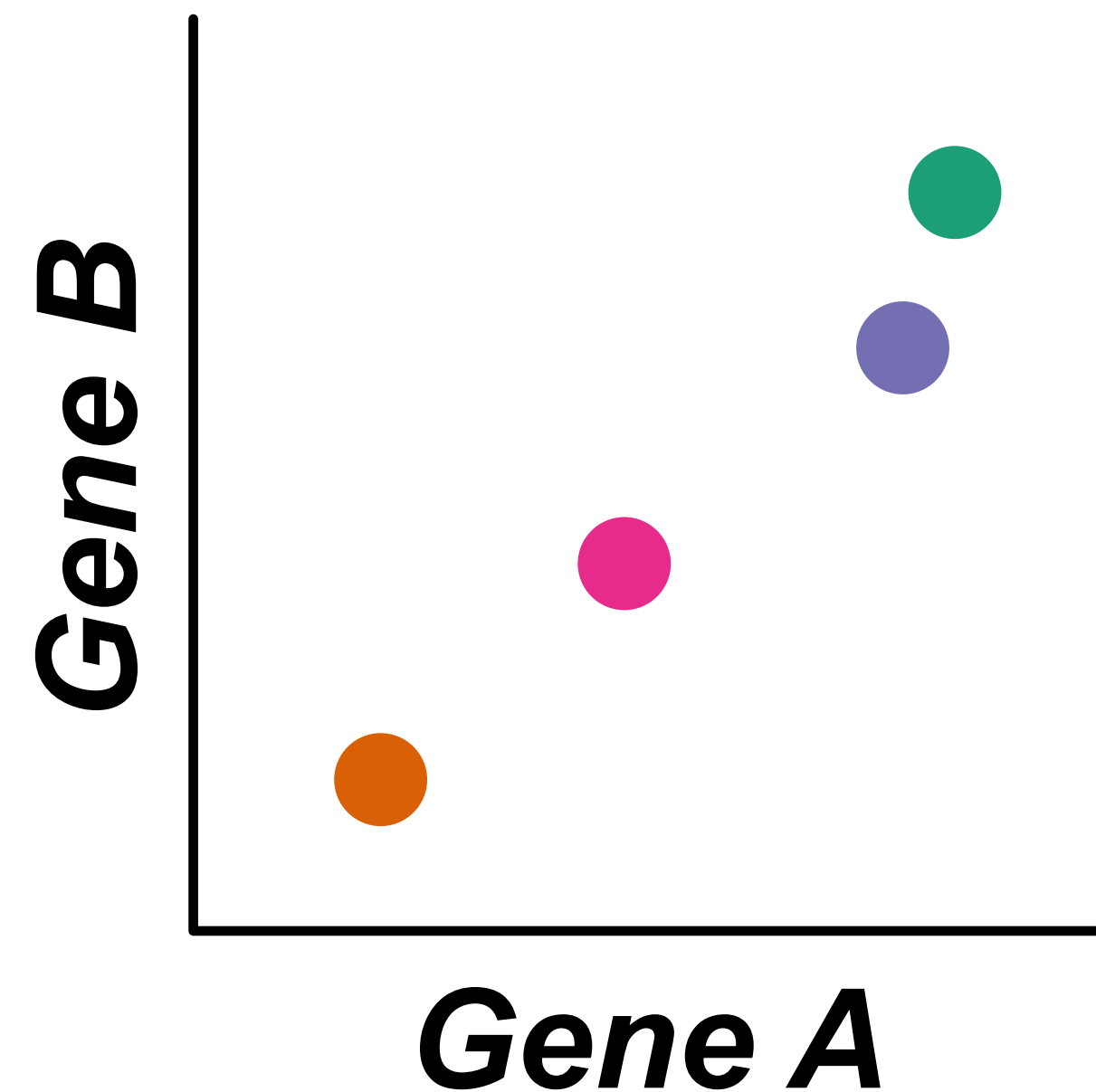


The mirror principle to detect gene coevolution



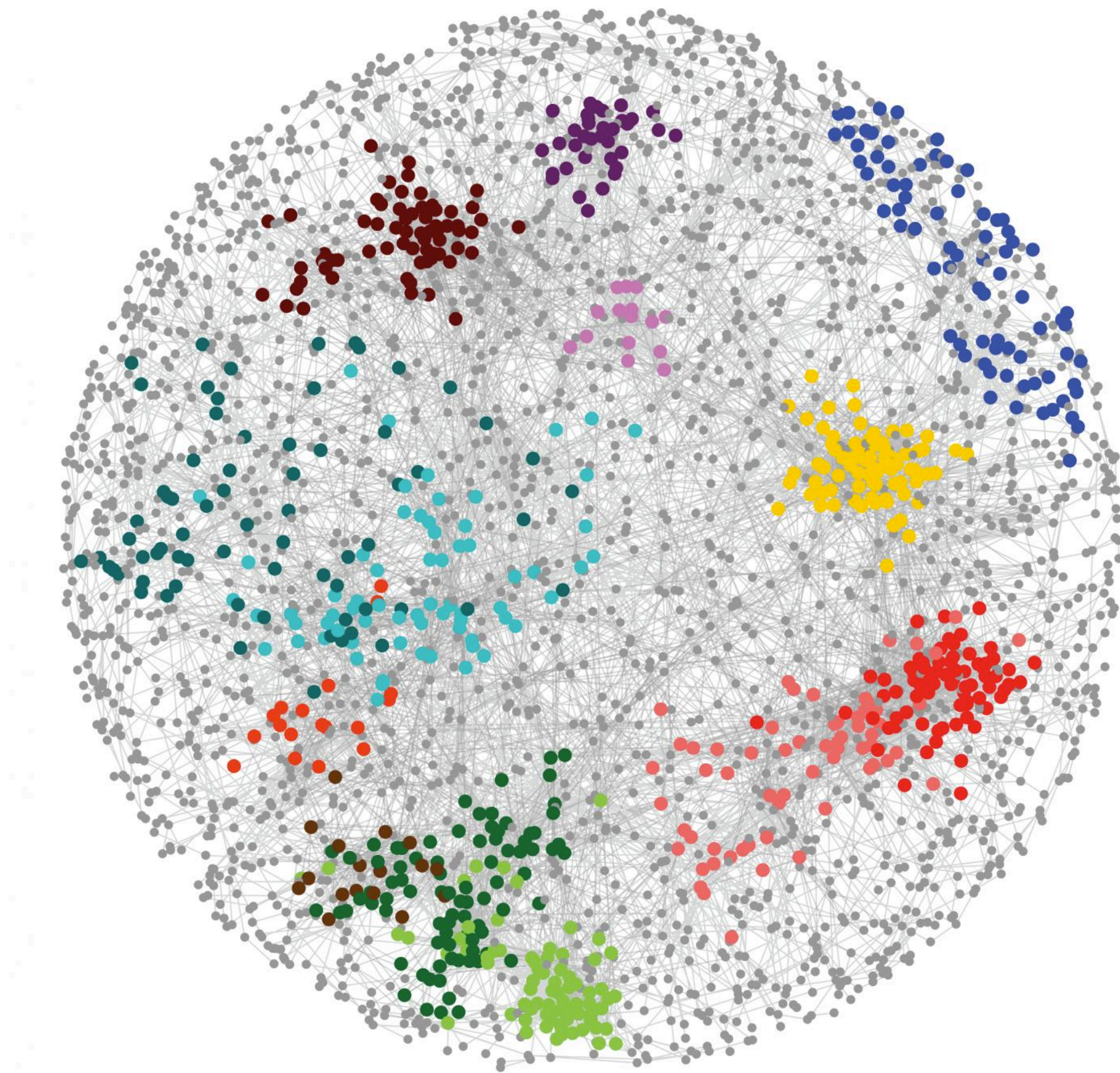
Genes of a feather evolve together

✓ Correlated evolution



- Coevolving genes tend to share function, be coexpressed, or are part of the same multimeric complexes

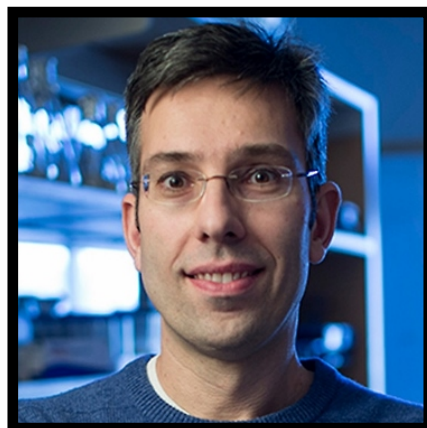
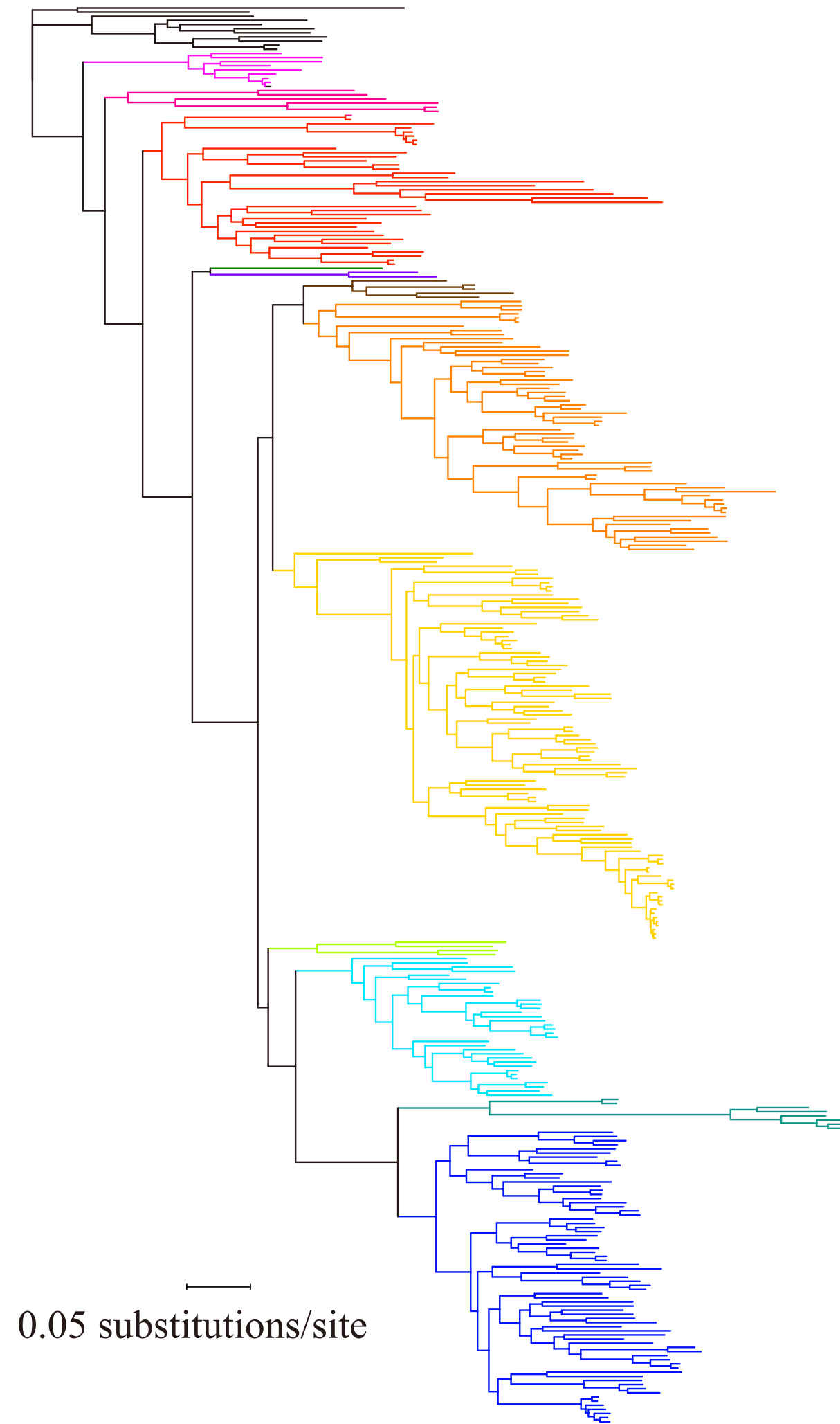
Genes of a feather evolve together



- Coevolving genes tend to share function, be coexpressed, or are part of the same multimeric complexes
- **But can we build a genetic network?**

Saccharomycotina yeast

- Saccharomycotina, a budding model subphylum



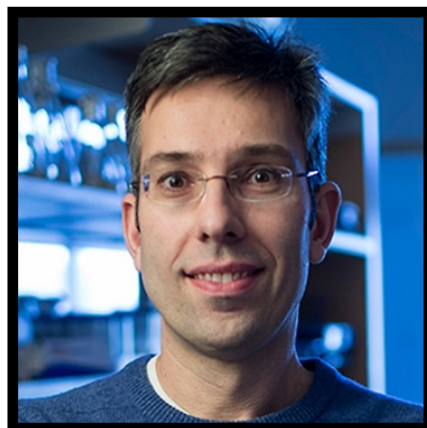
Antonis Rokas



Chris Hittinger

Saccharomycotina yeast

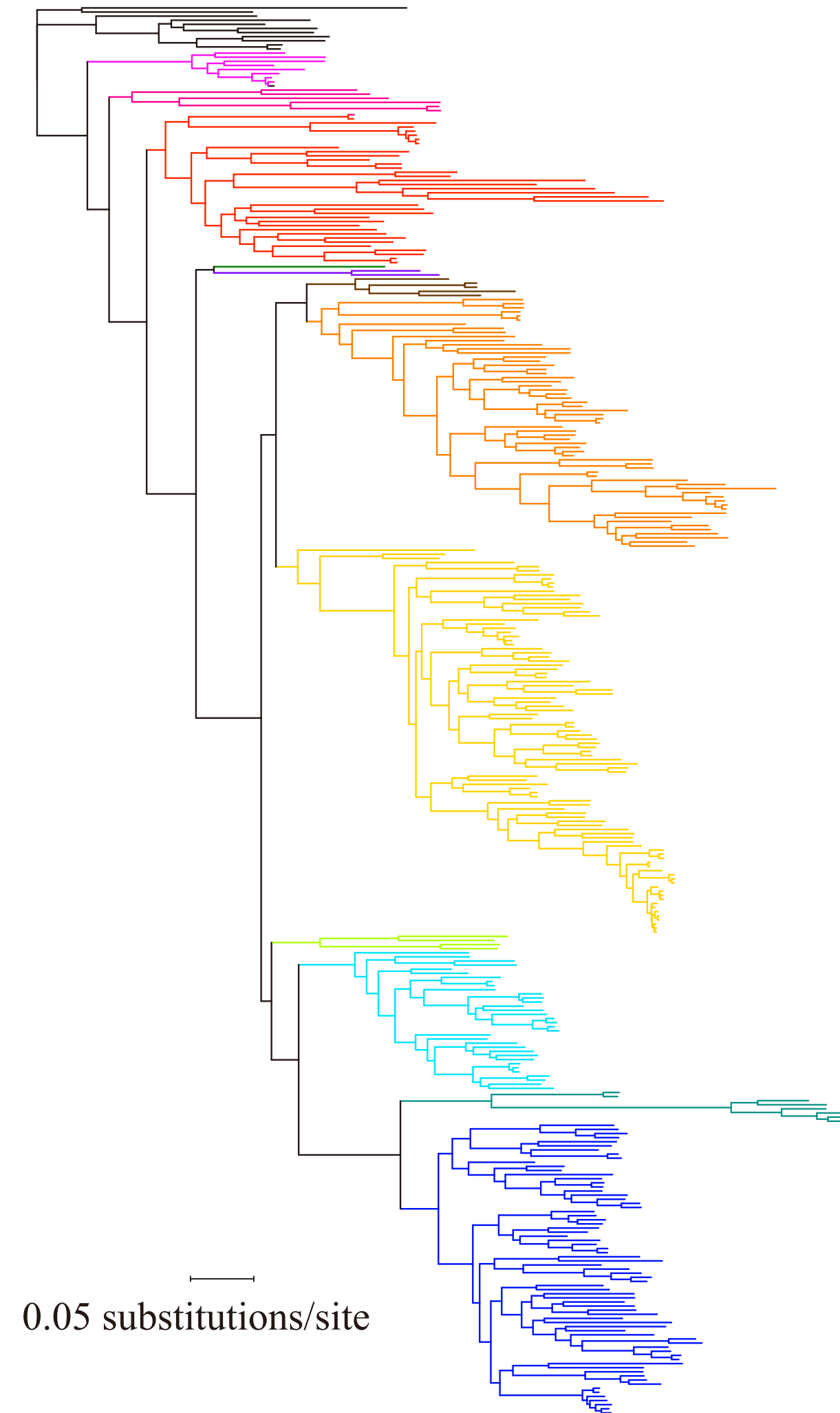
- Saccharomycotina, a budding model subphylum
- Spans 332 species of budding yeast



Antonis Rokas

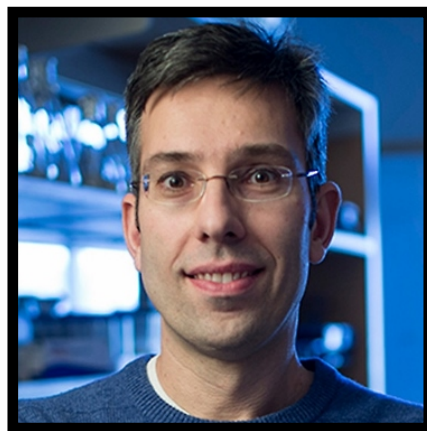


Chris Hittinger



Saccharomycotina yeast

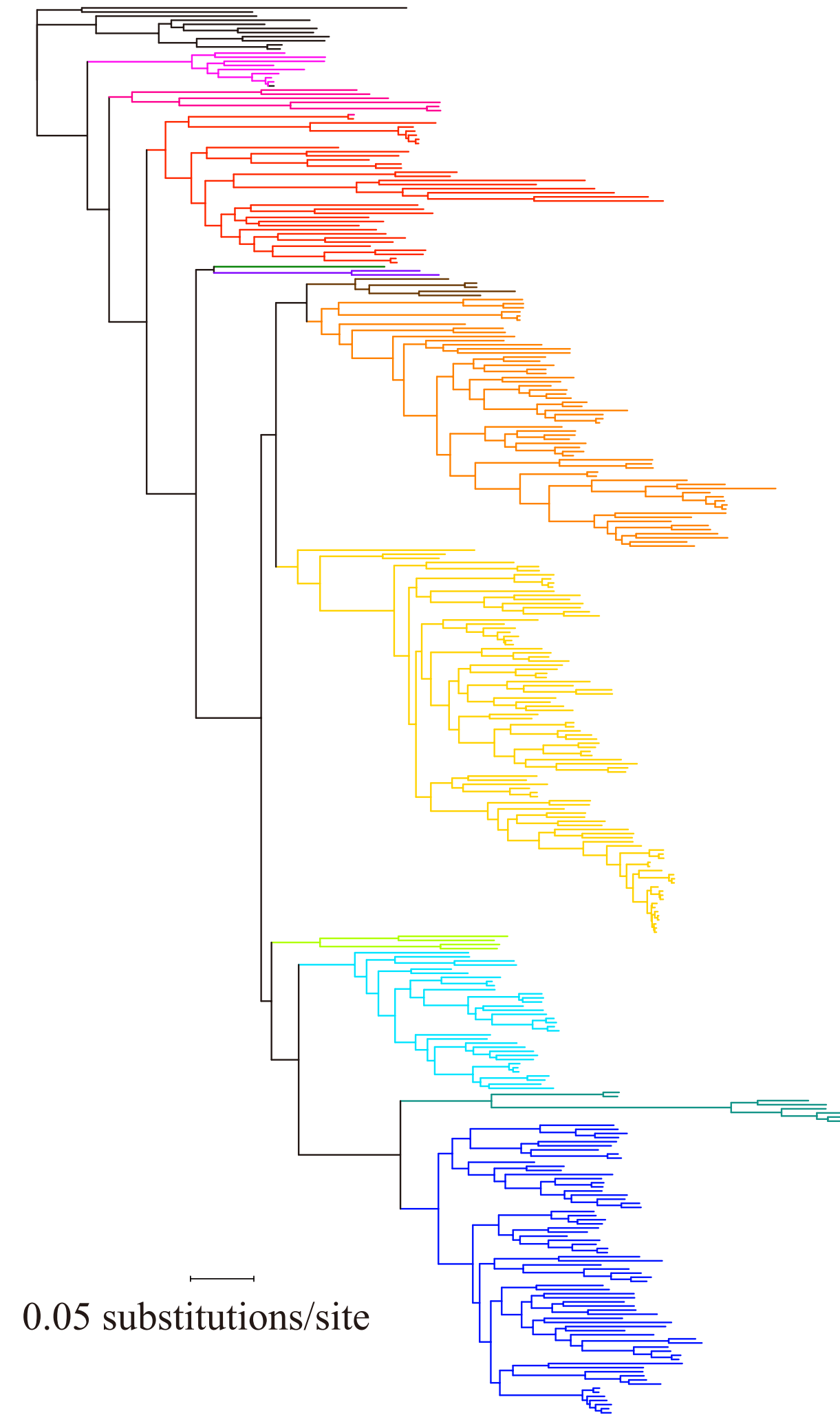
- Saccharomycotina, a budding model subphylum
- Spans 332 species of budding yeast
- 2,408 orthologous genes across all budding yeasts



Antonis Rokas

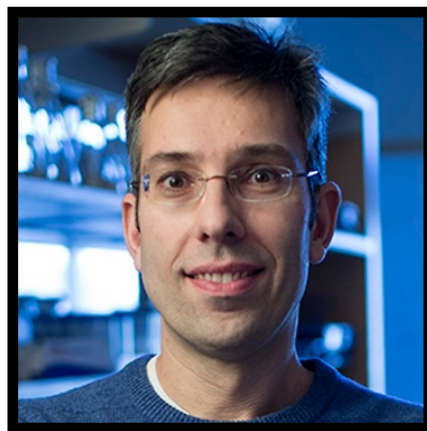


Chris Hittinger

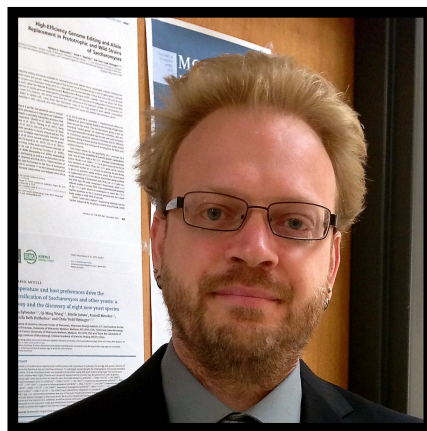


0.05 substitutions/site

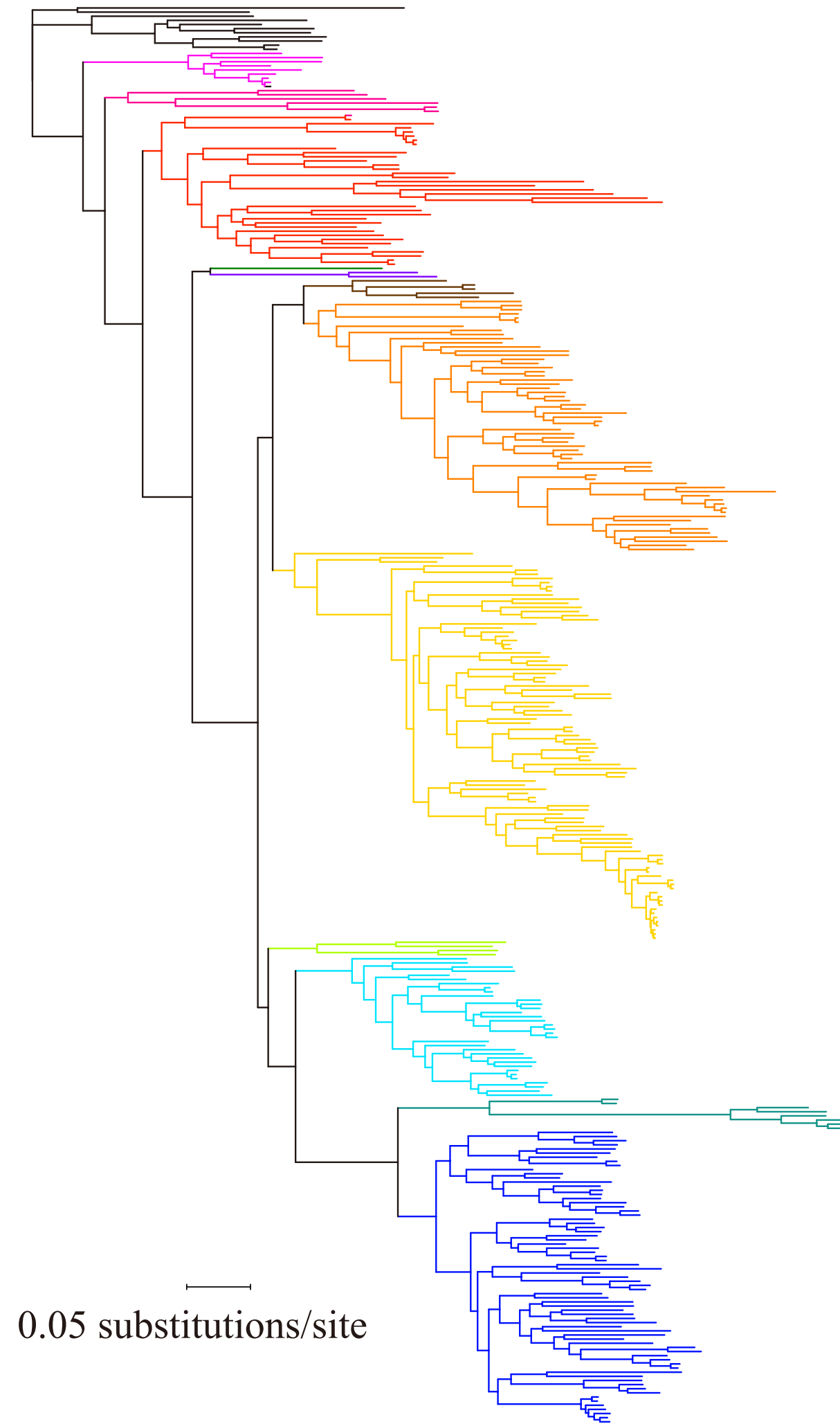
Saccharomycotina yeast



Antonis Rokas



Chris Hittinger



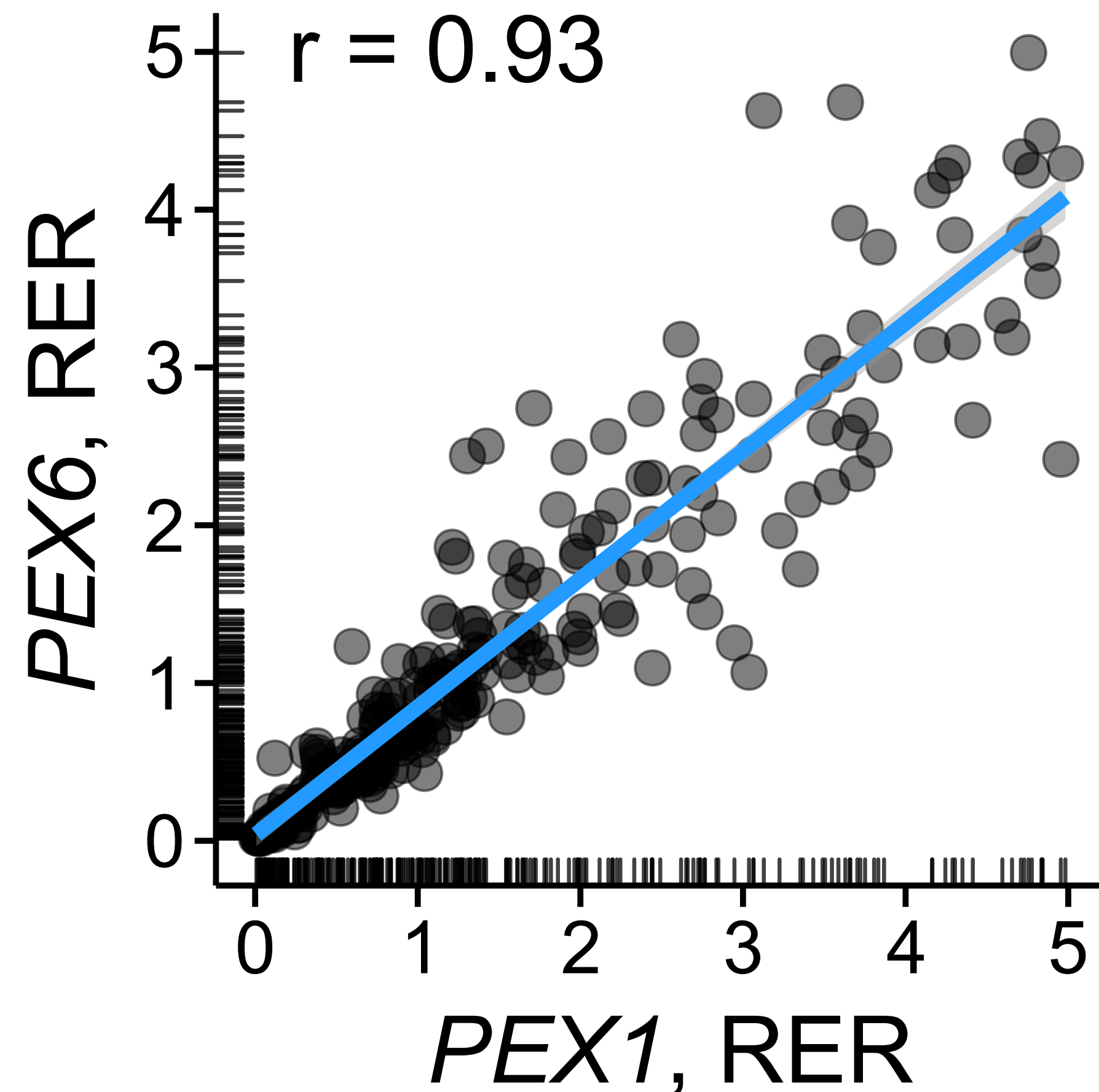
- Saccharomycotina, a budding model subphylum

- Spans 332 species of budding yeast

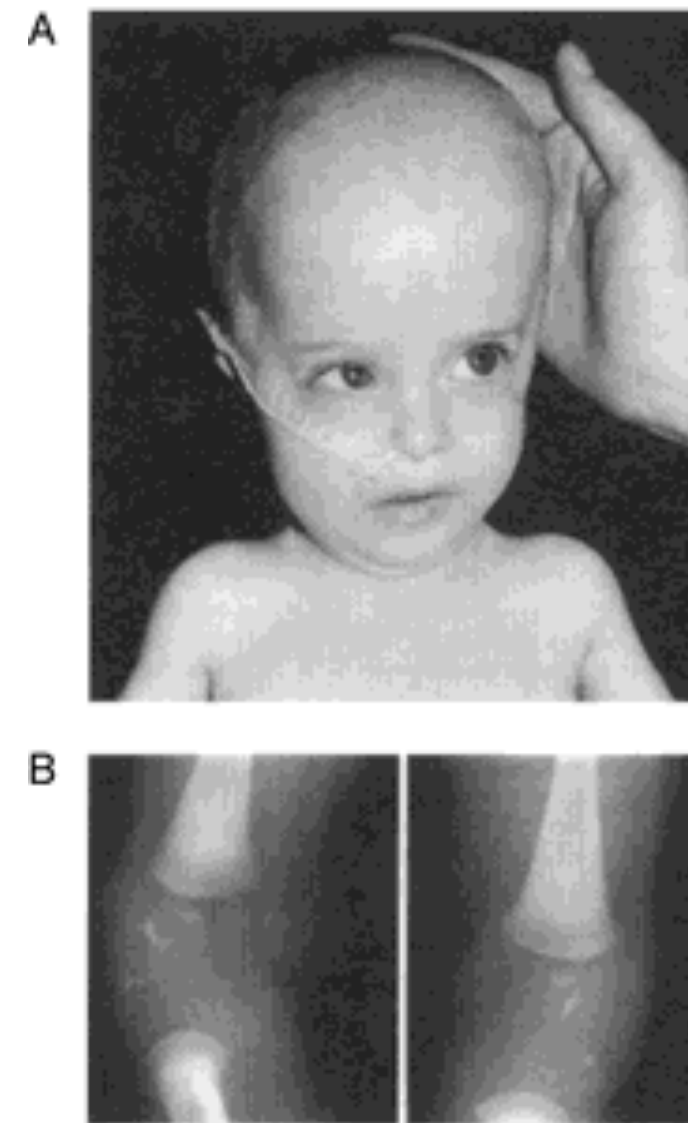
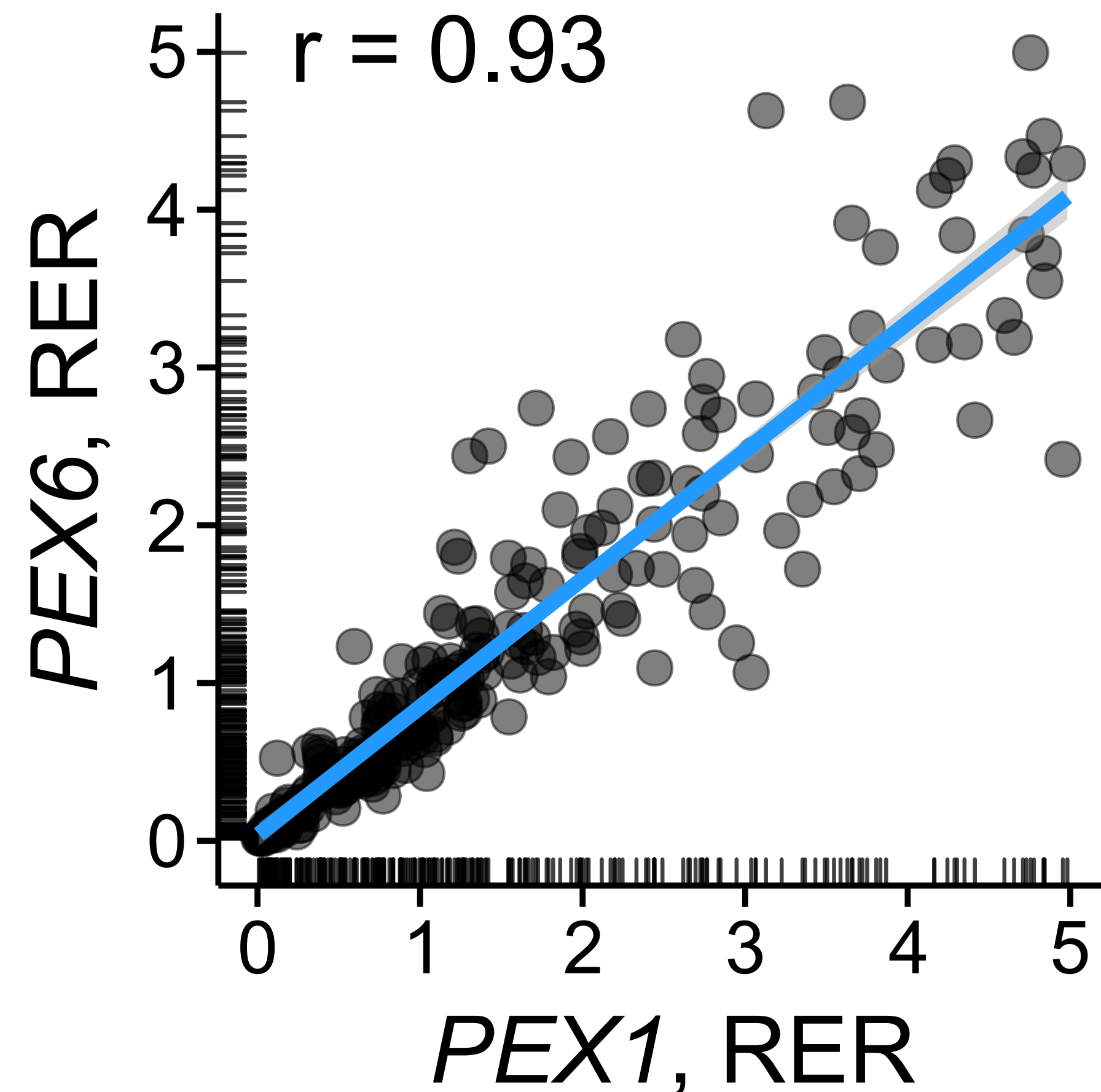
- 2,408 orthologous genes across all budding yeasts

- Calculate gene covariation across ~3 million pairwise combinations of genes

PEX1 and PEX6 are coevolving

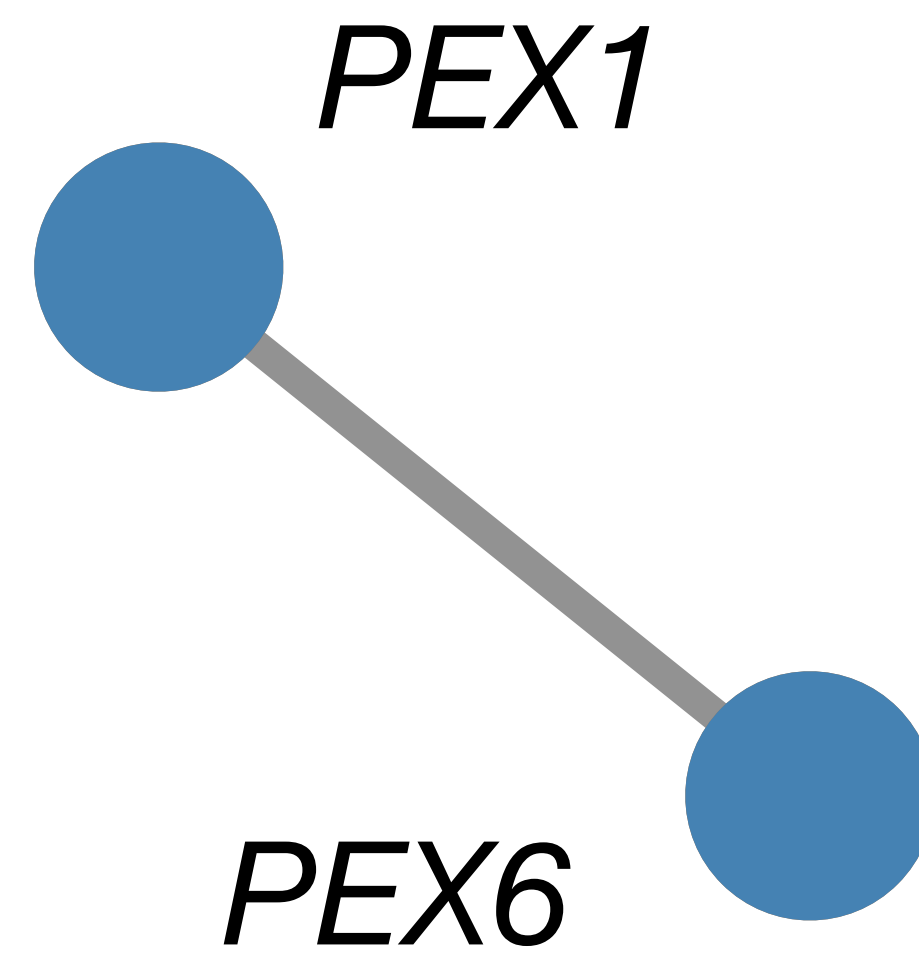


PEX1 and PEX6 are coevolving

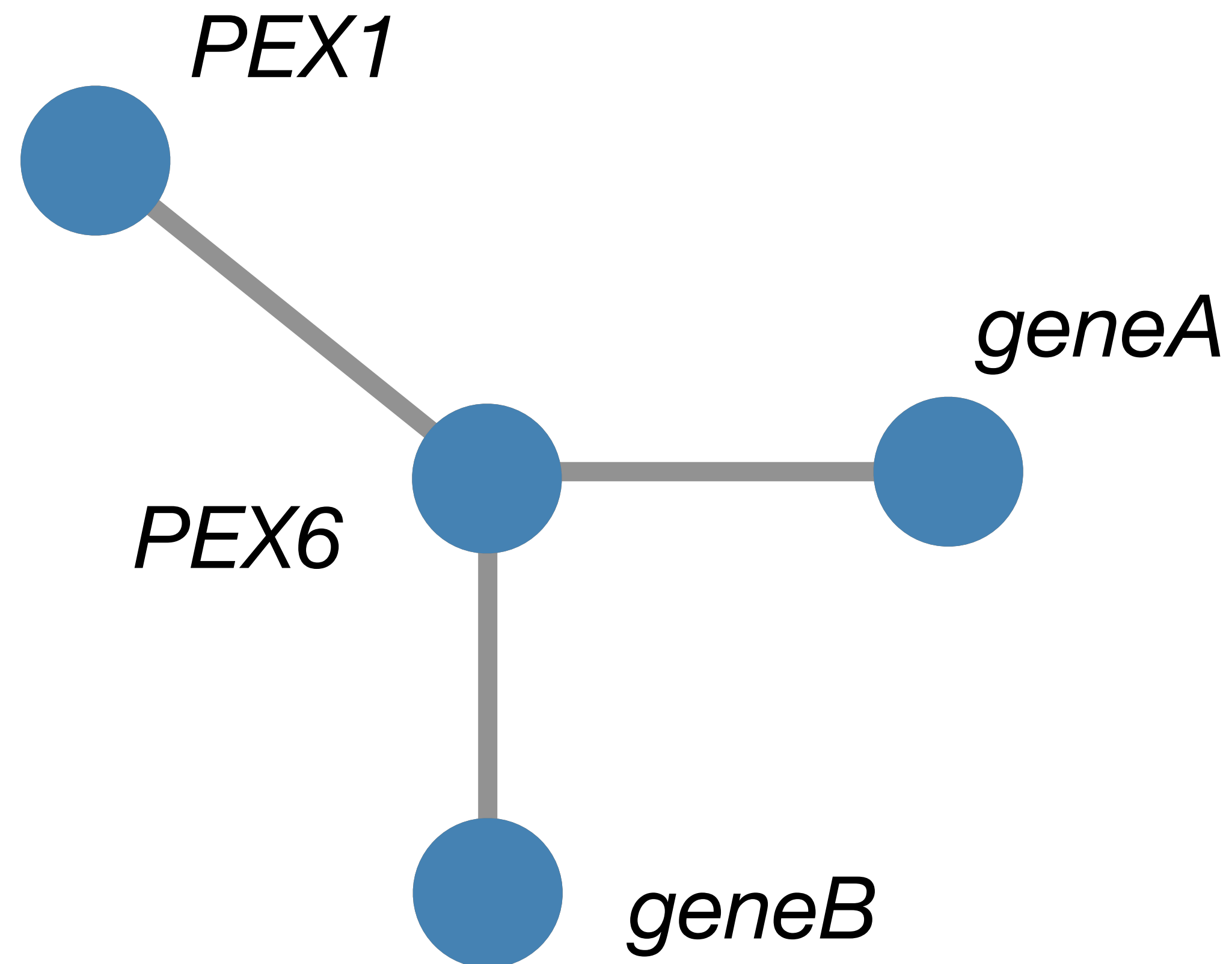


Pex1p & Pex6p: forms a heterodimer involved in recycling peroxisomal signal receptor Pex5p

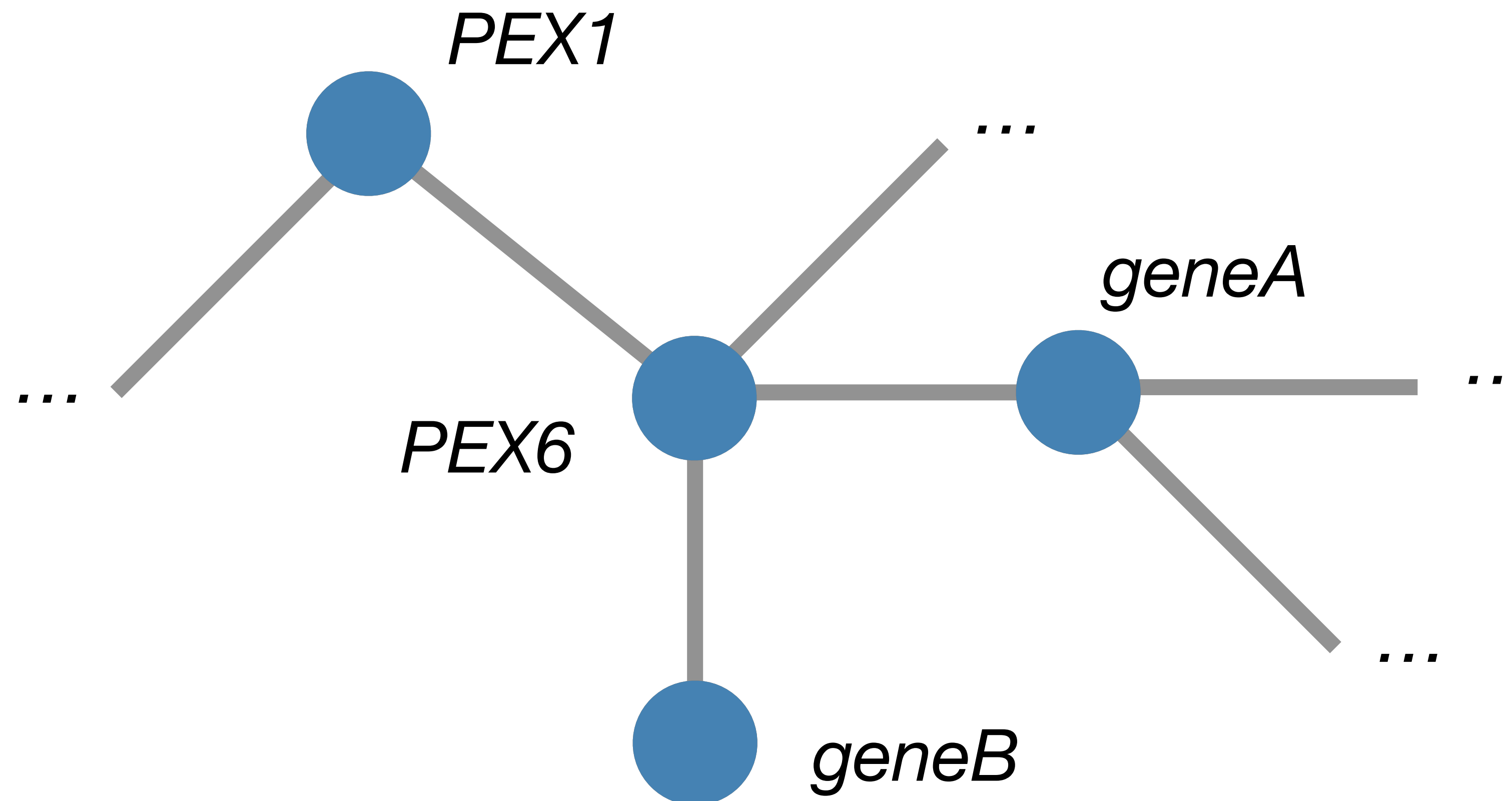
Constructing a coevolutionary genetic network



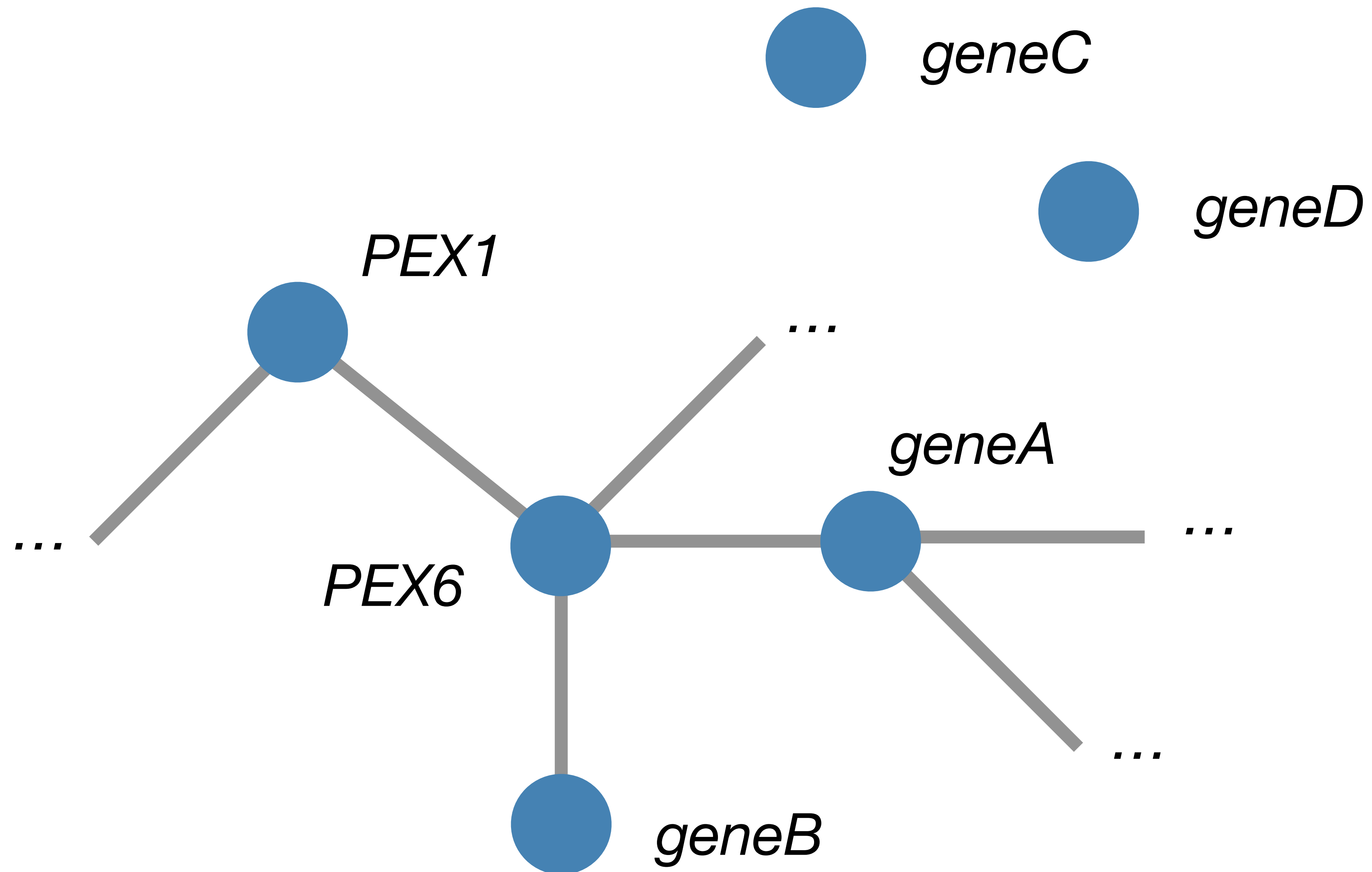
Constructing a coevolutionary genetic network



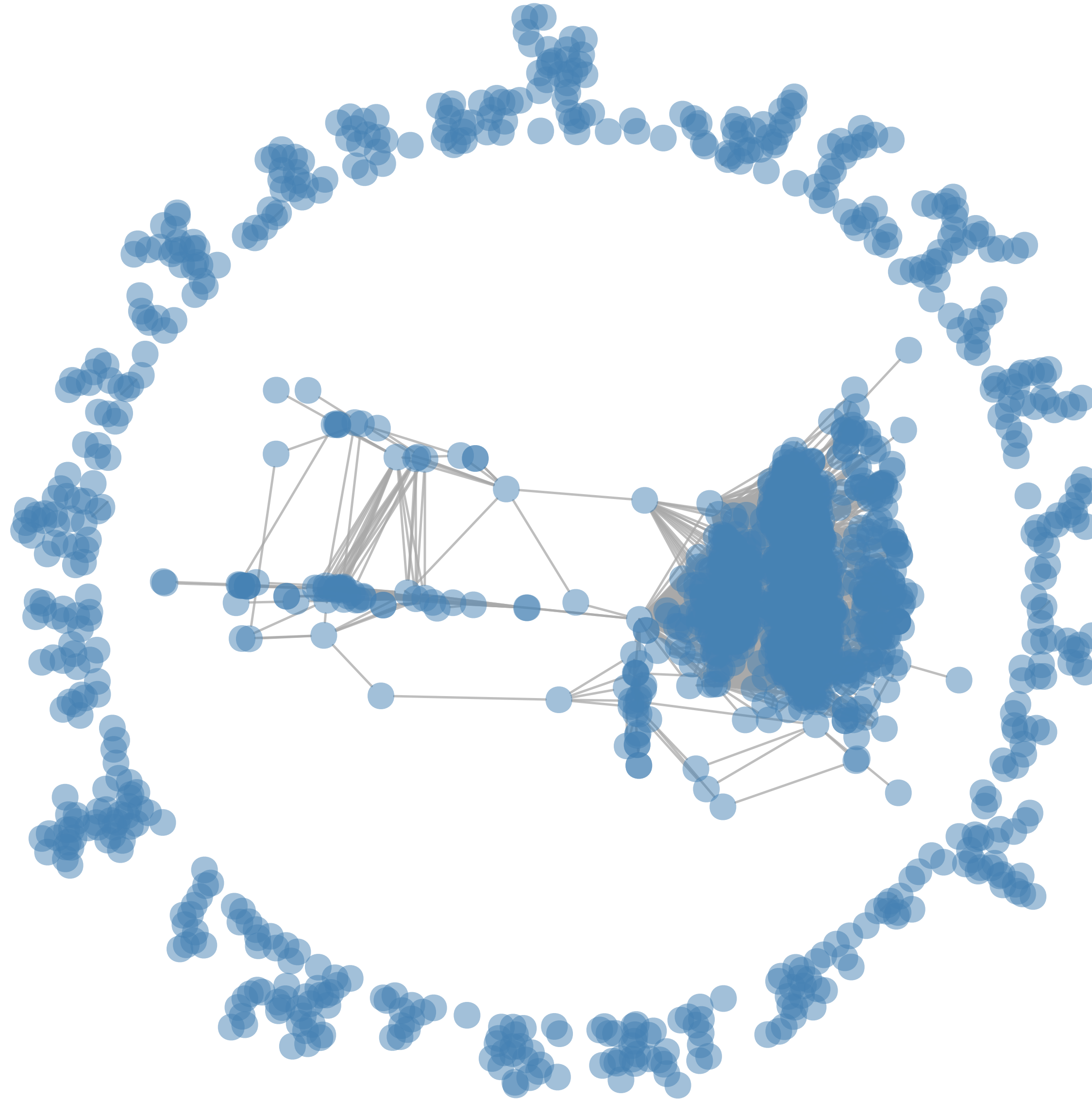
Constructing a coevolutionary genetic network



Constructing a coevolutionary genetic network



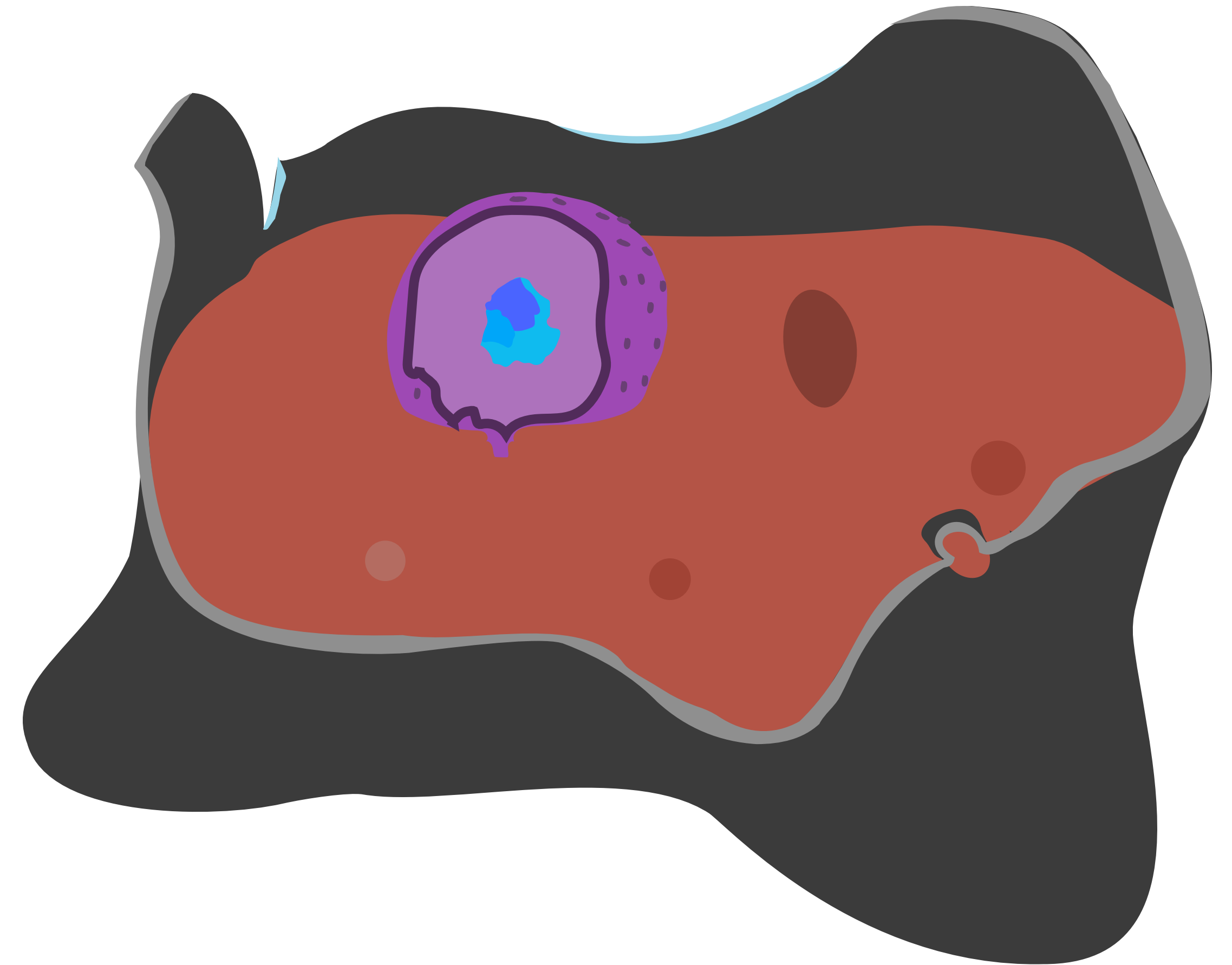
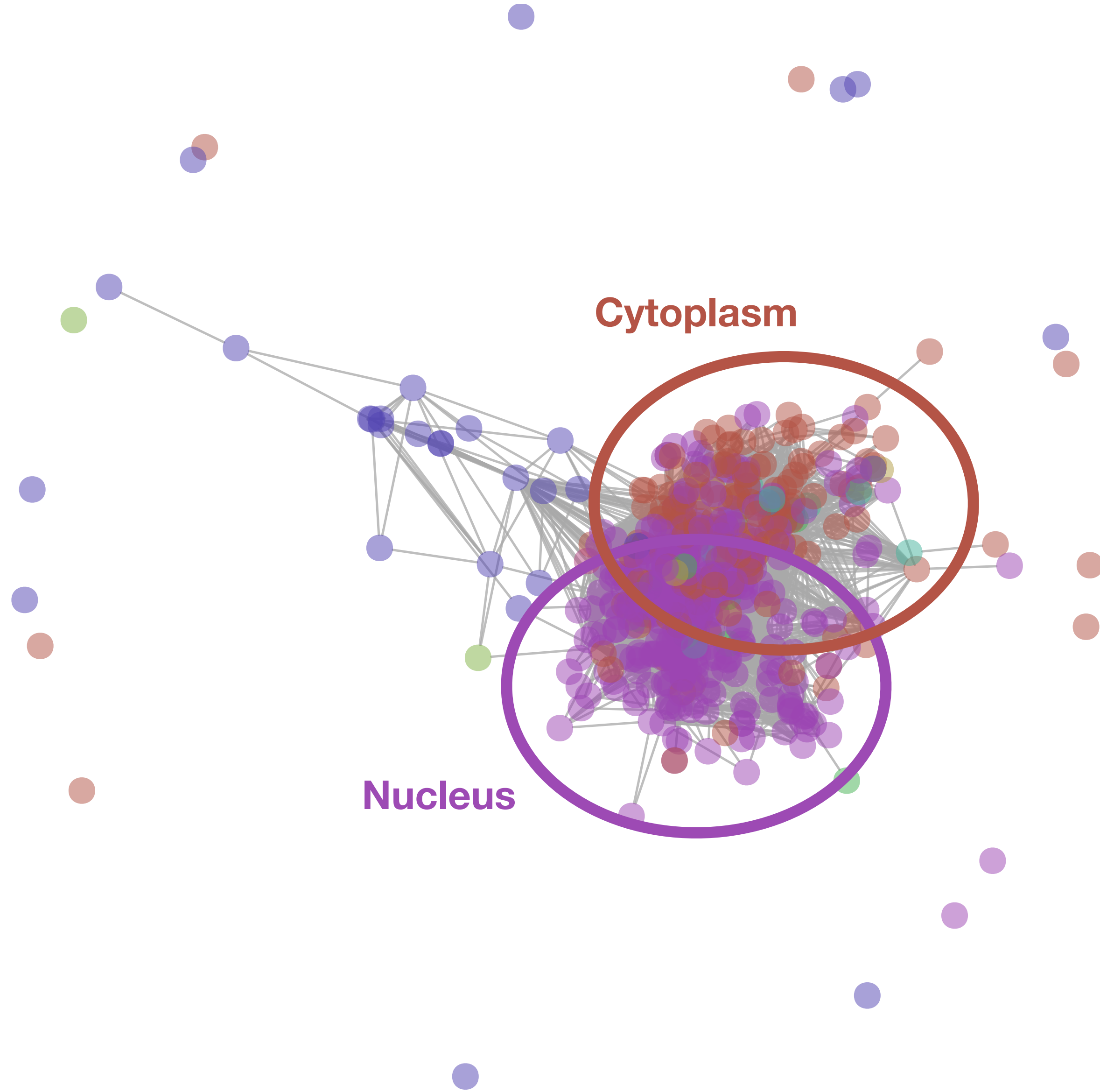
A global gene coevolutionary network



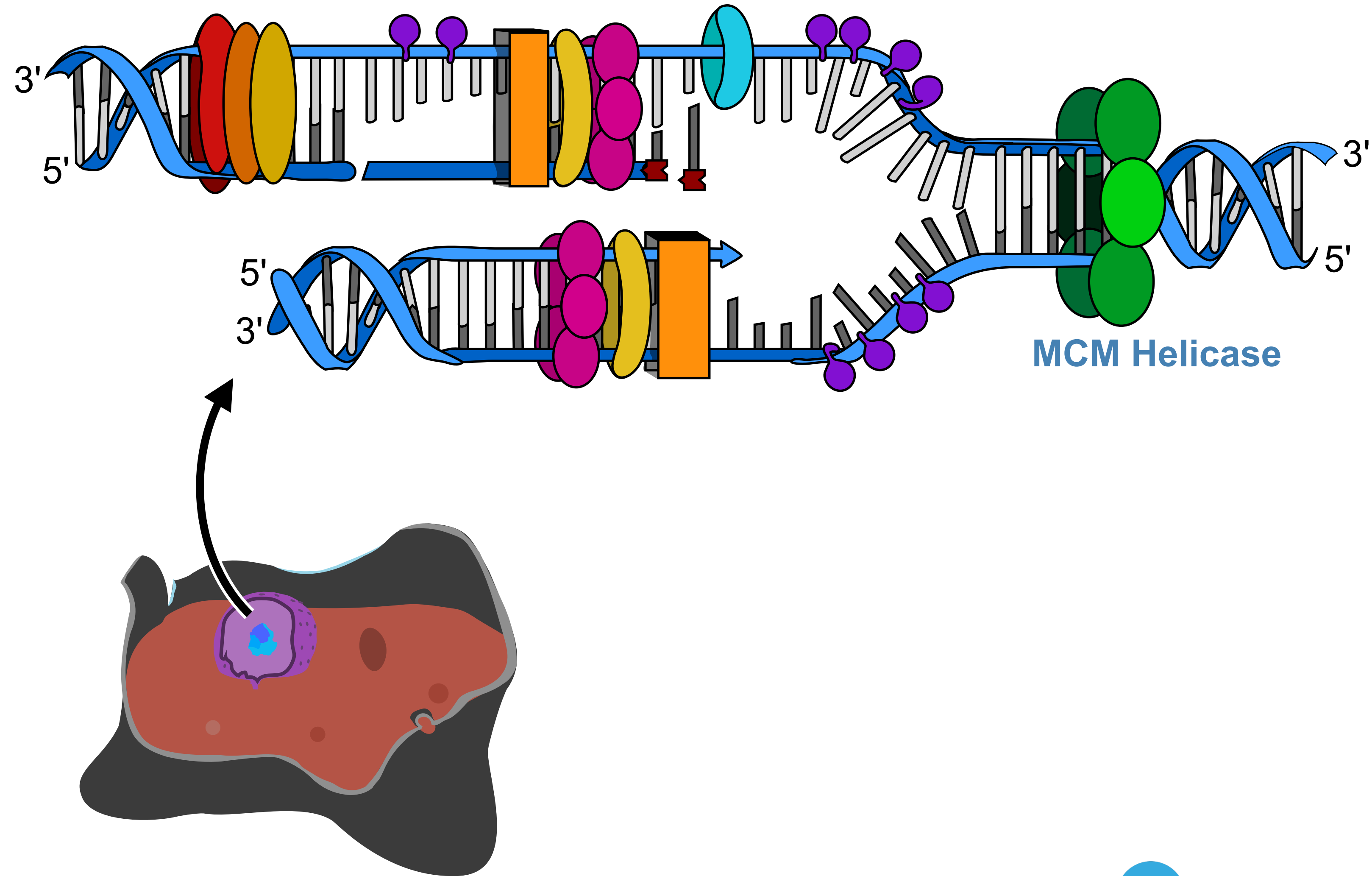
Nodes are genes

Edges connected
coevolving genes

Network reflections of cellular structure

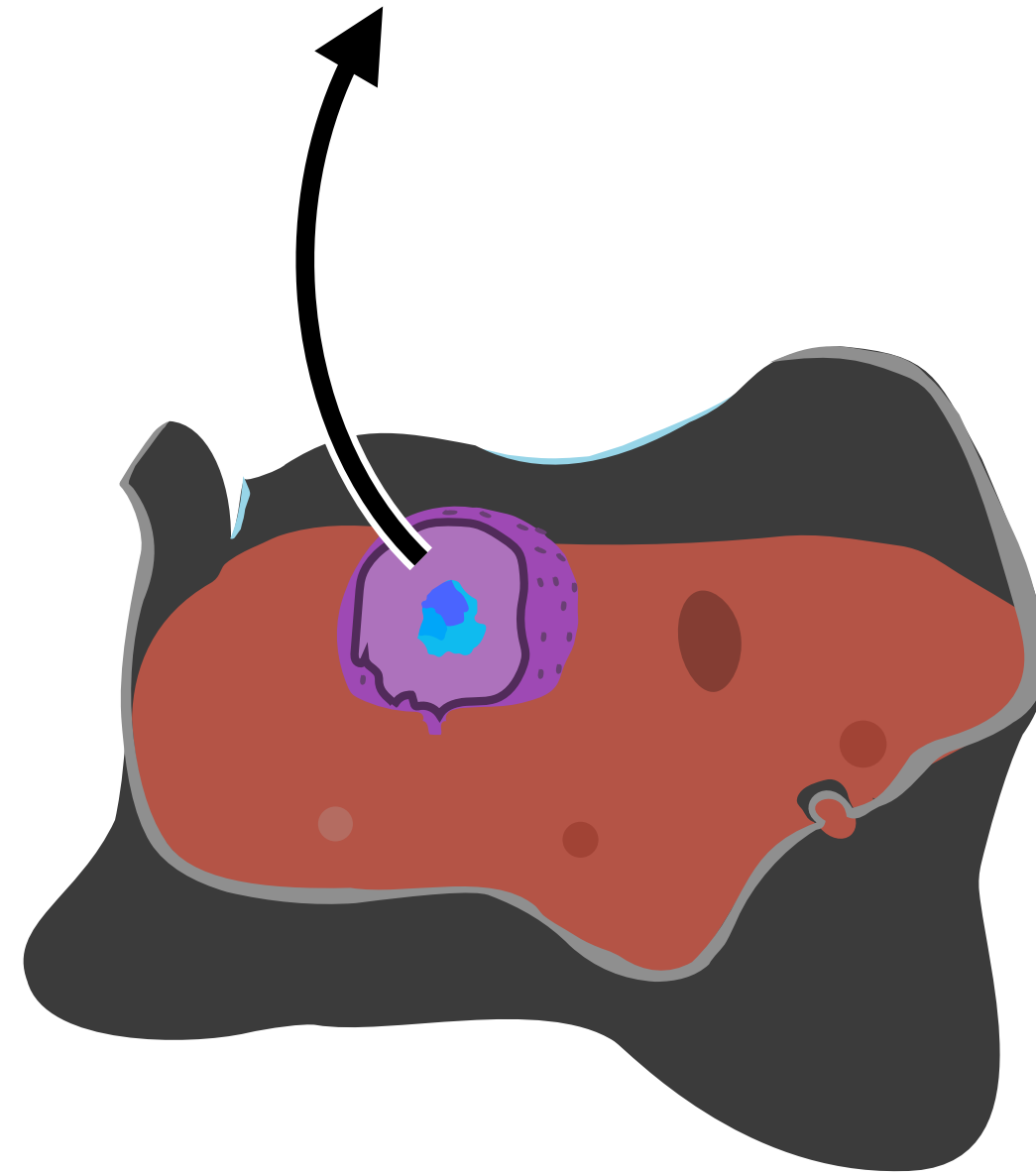
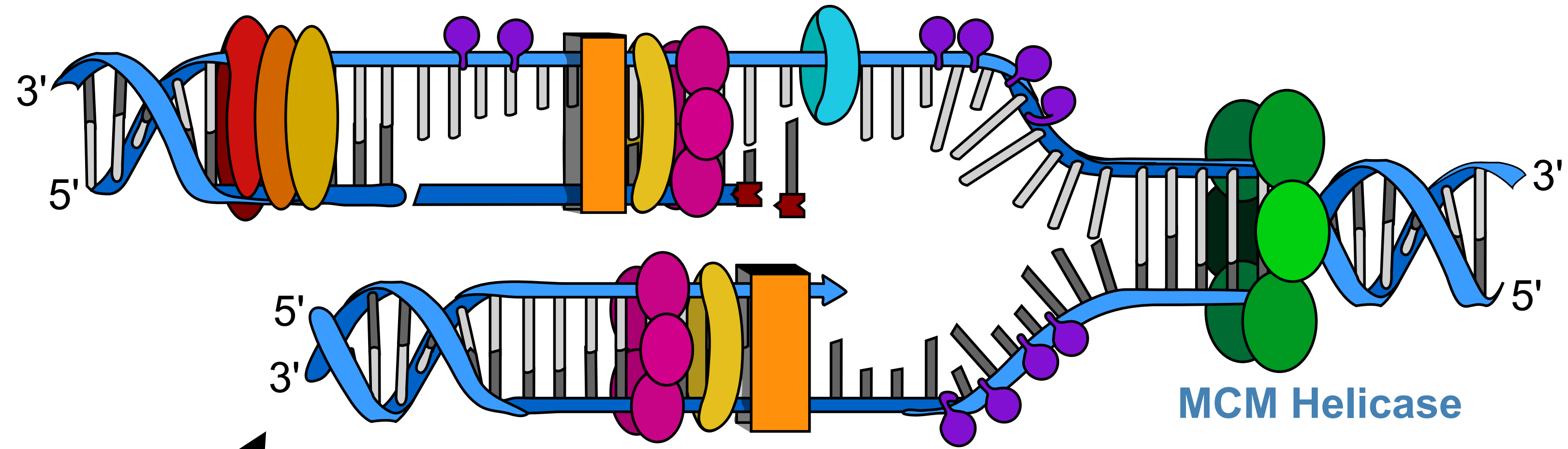
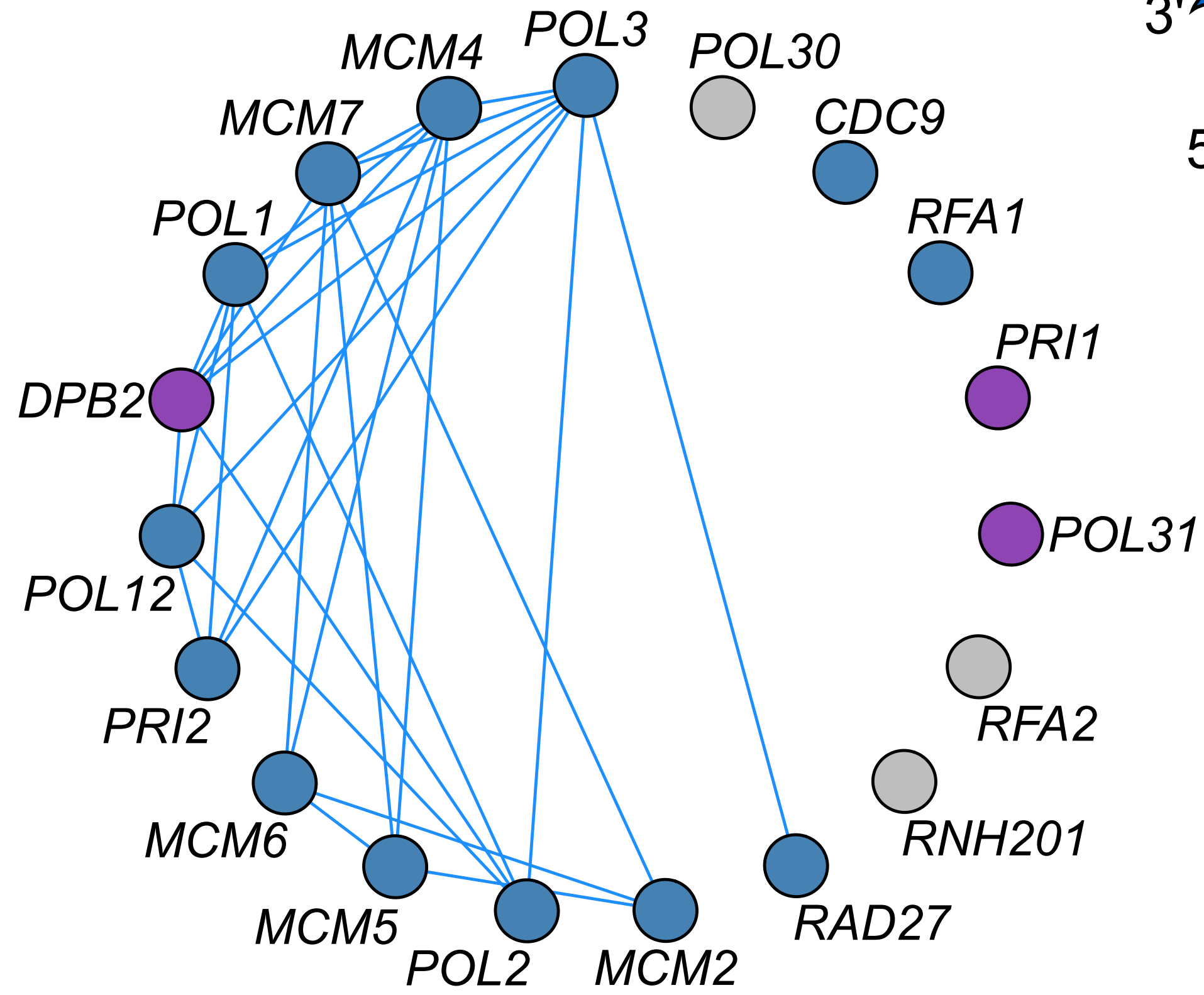


Genes from pathways are coevolving

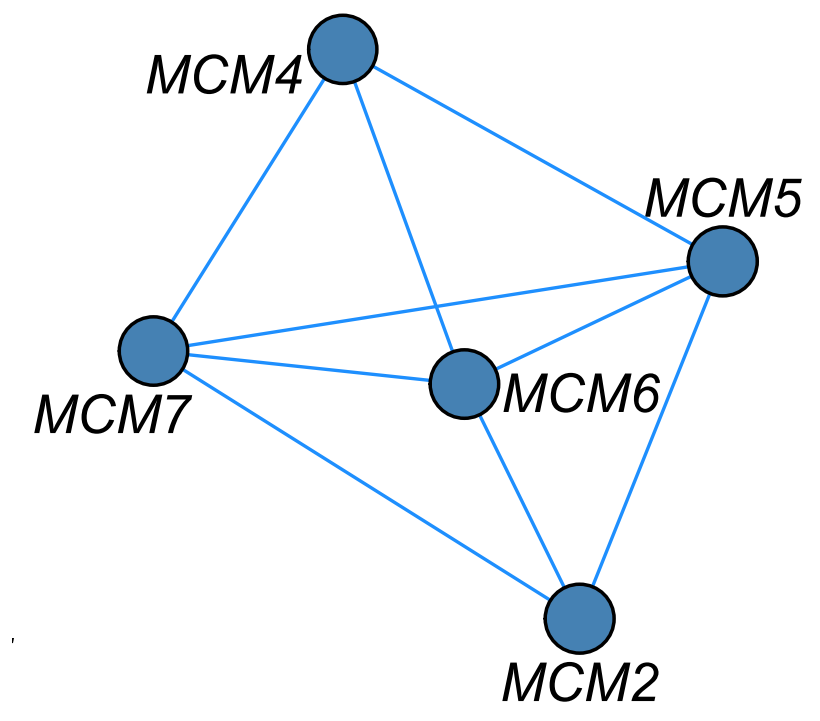
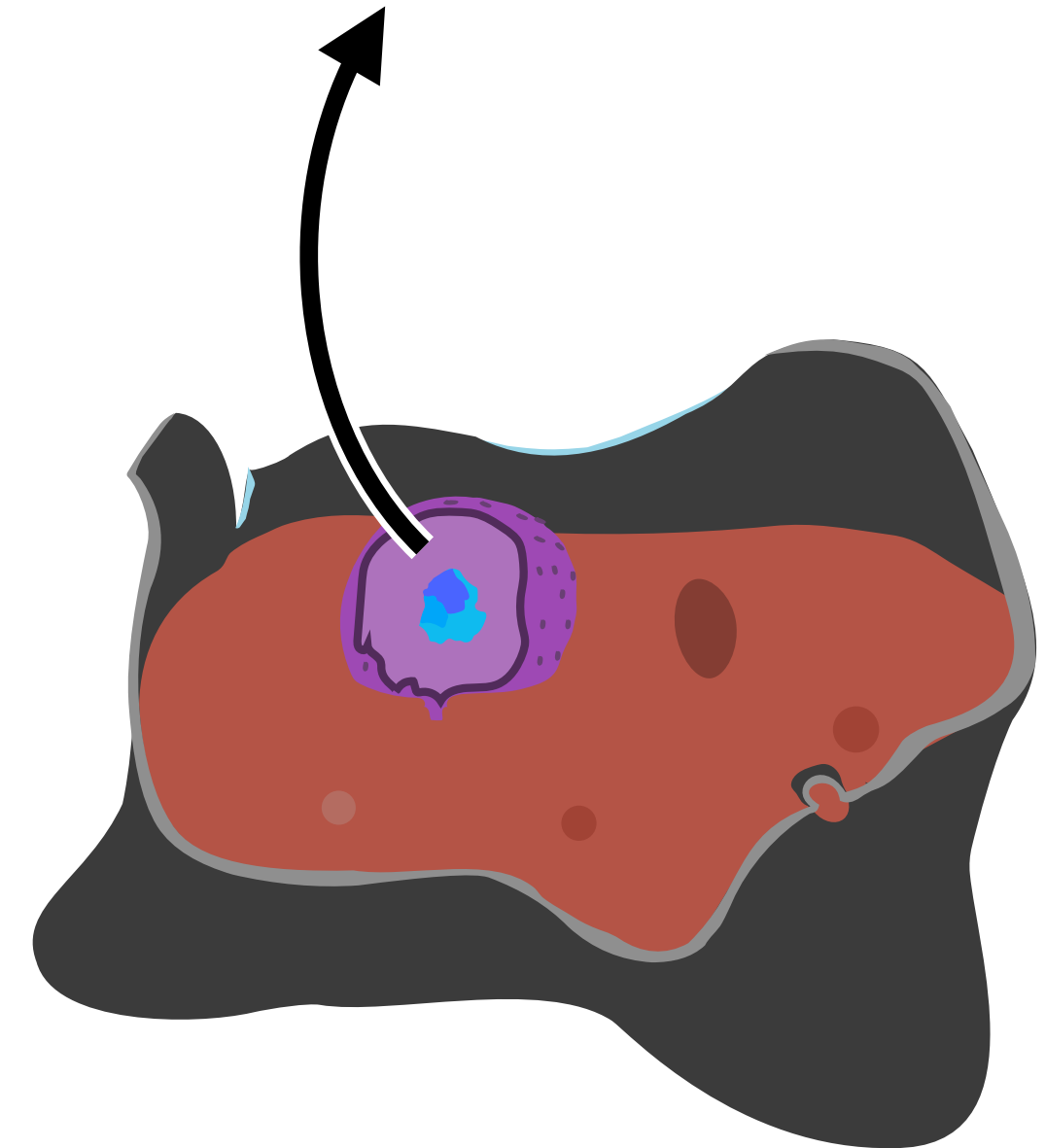
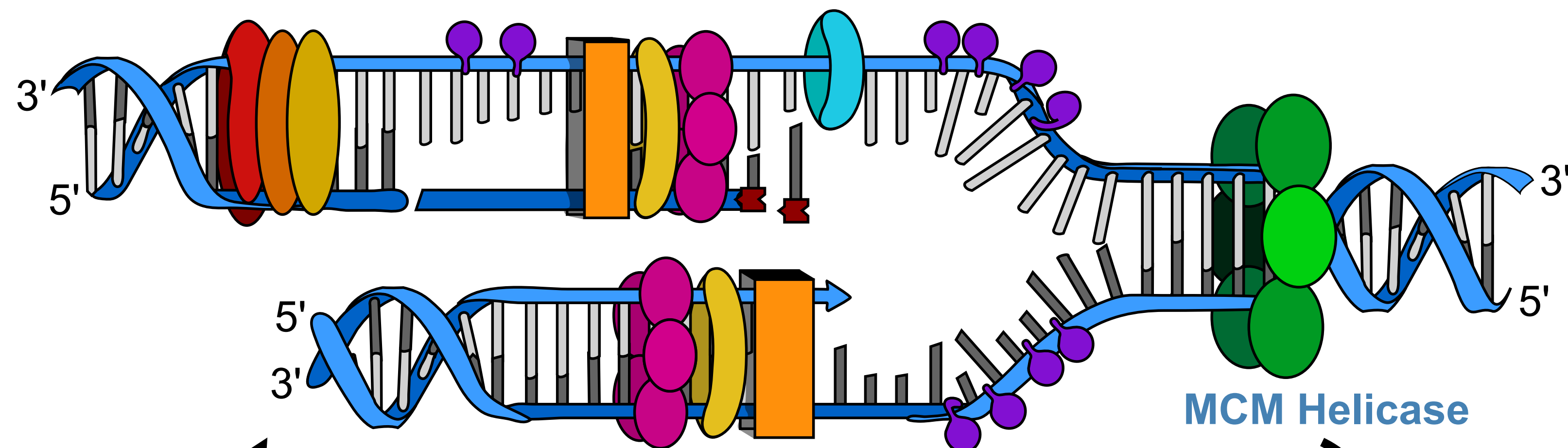
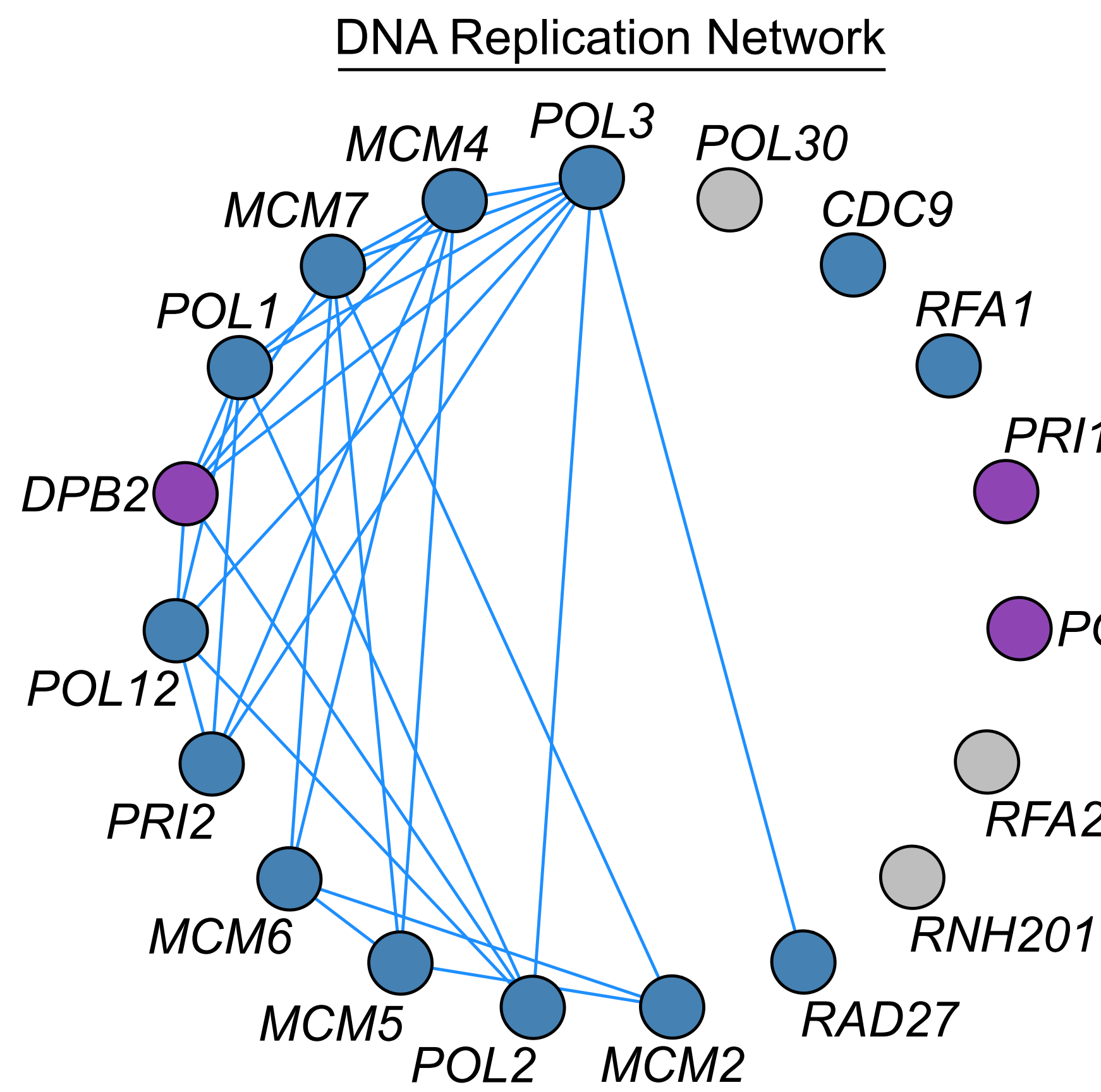


Genes from pathways are coevolving

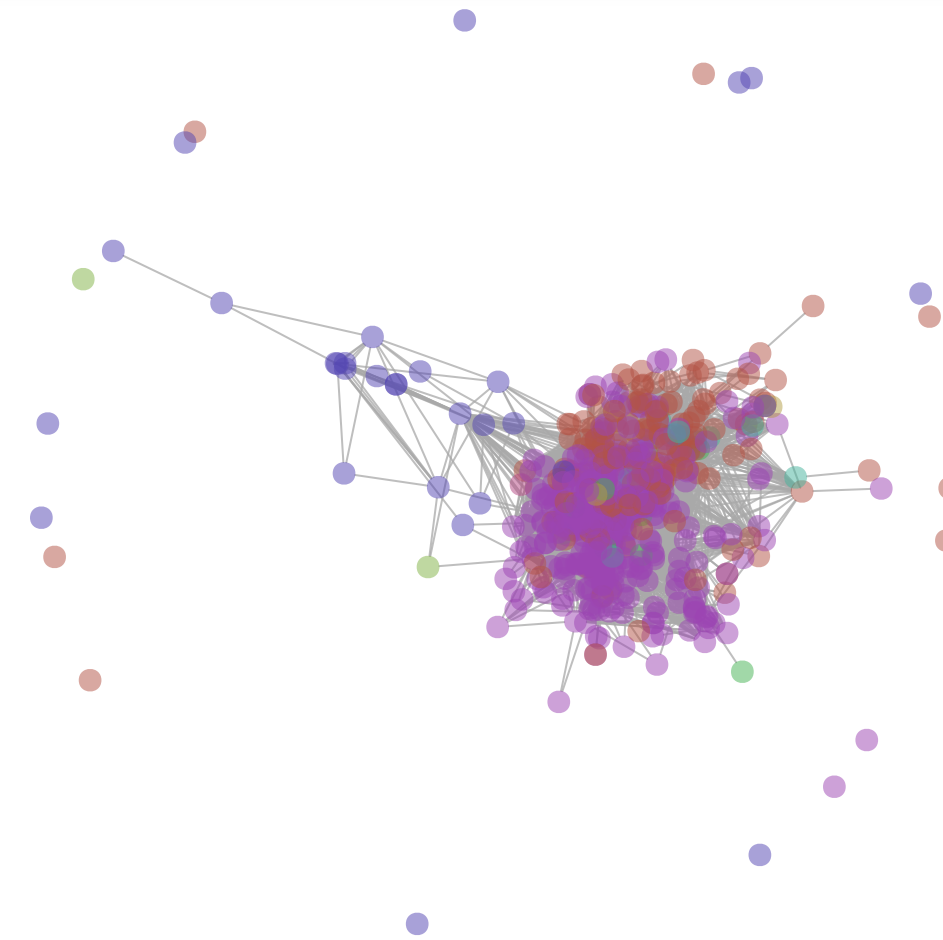
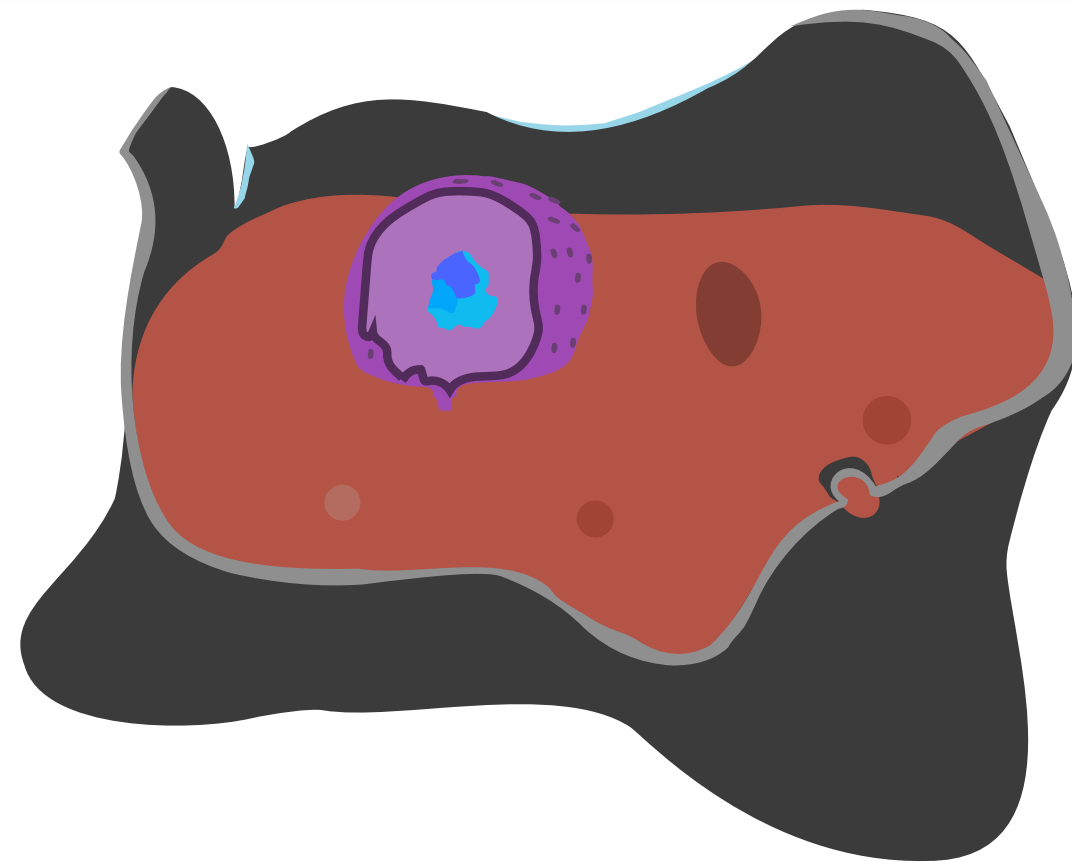
DNA Replication Network



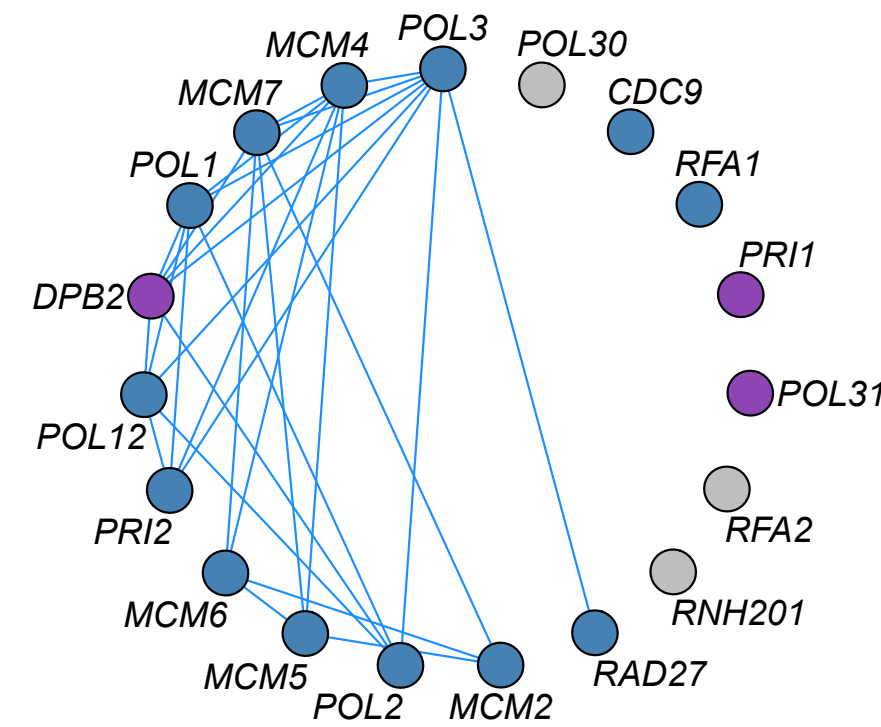
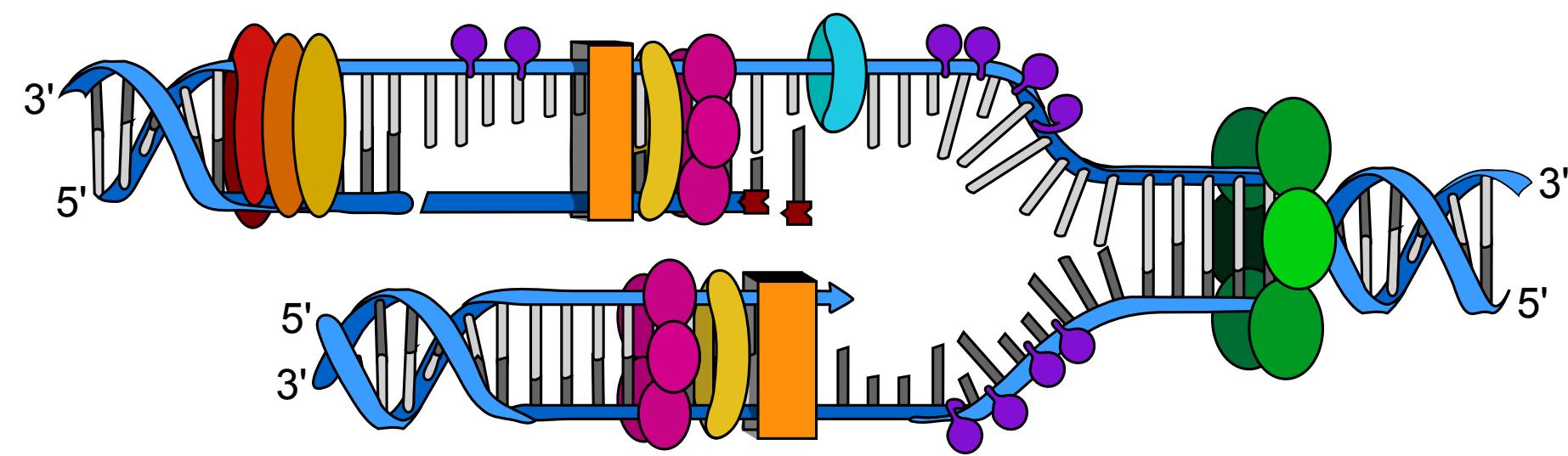
Genes from multimeric proteins are coevolving



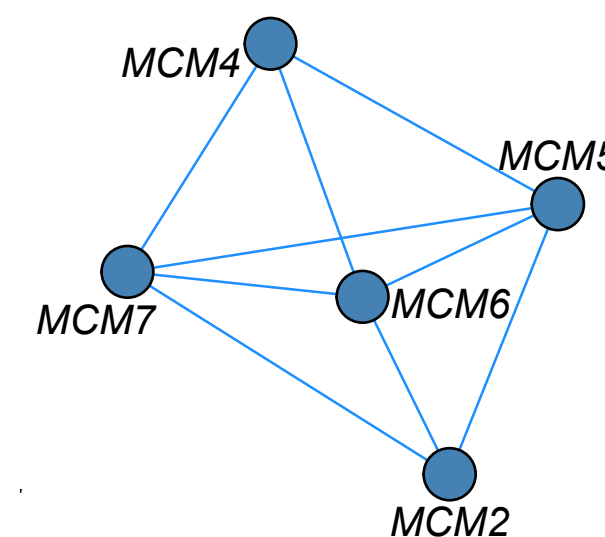
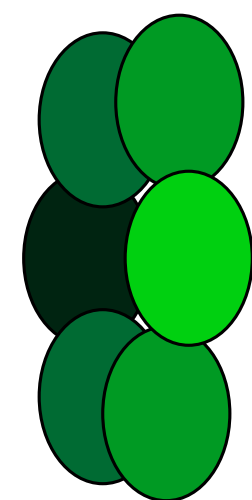
A global network provides insight to a hierarchy of function



Cellular



Bioprocess



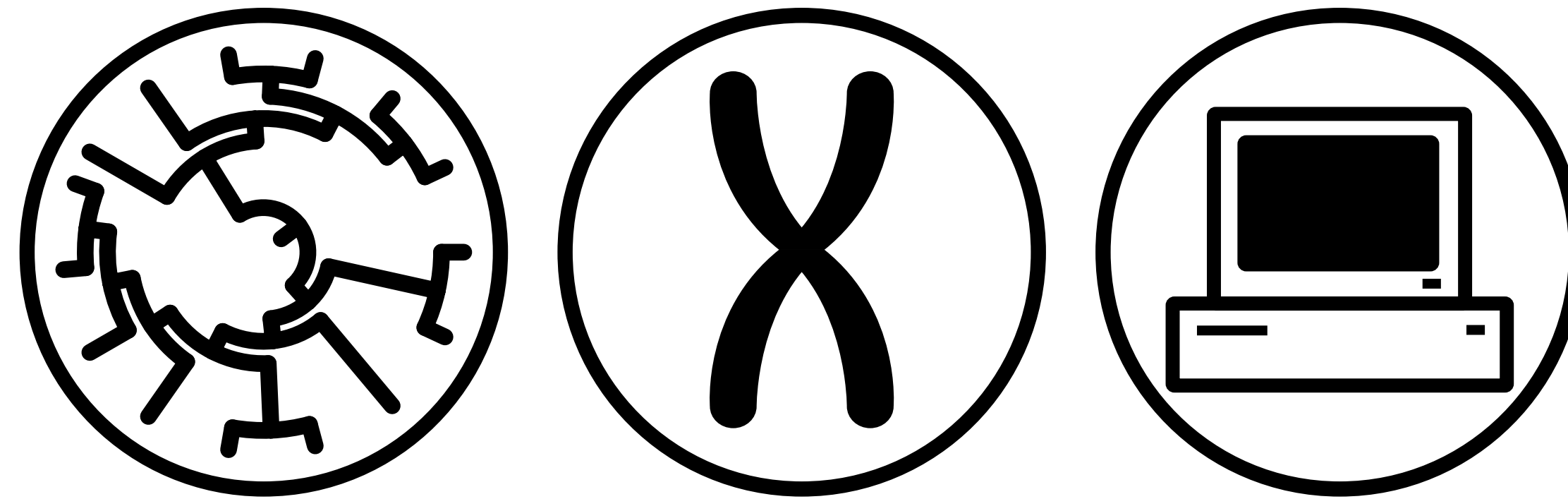
Protein complex

***Can signatures of
gene coevolution
provide insight to your
genes of interest?***

***Can signatures of
gene coevolution
provide insight to your
genes of interest?***

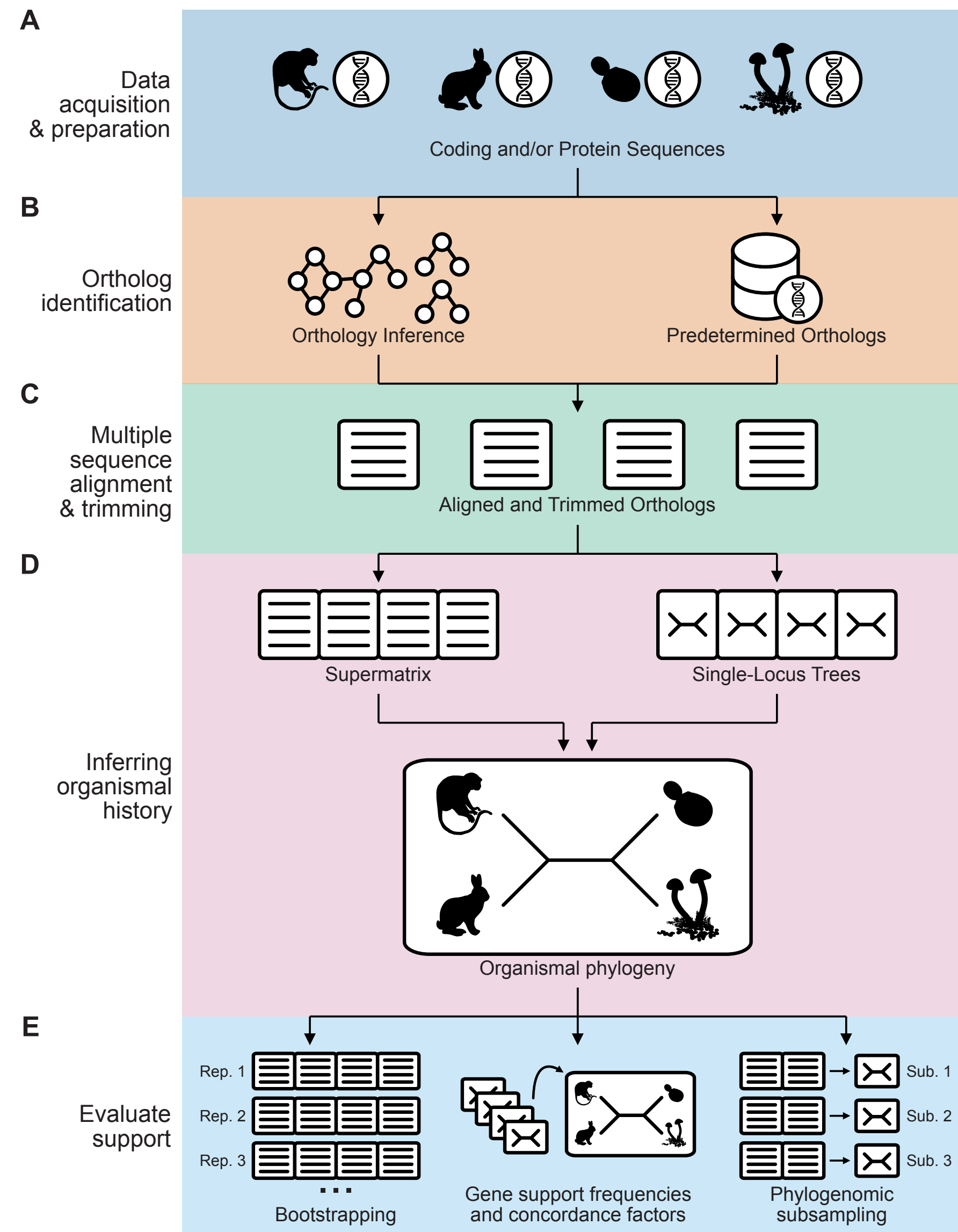
YES!

Outline

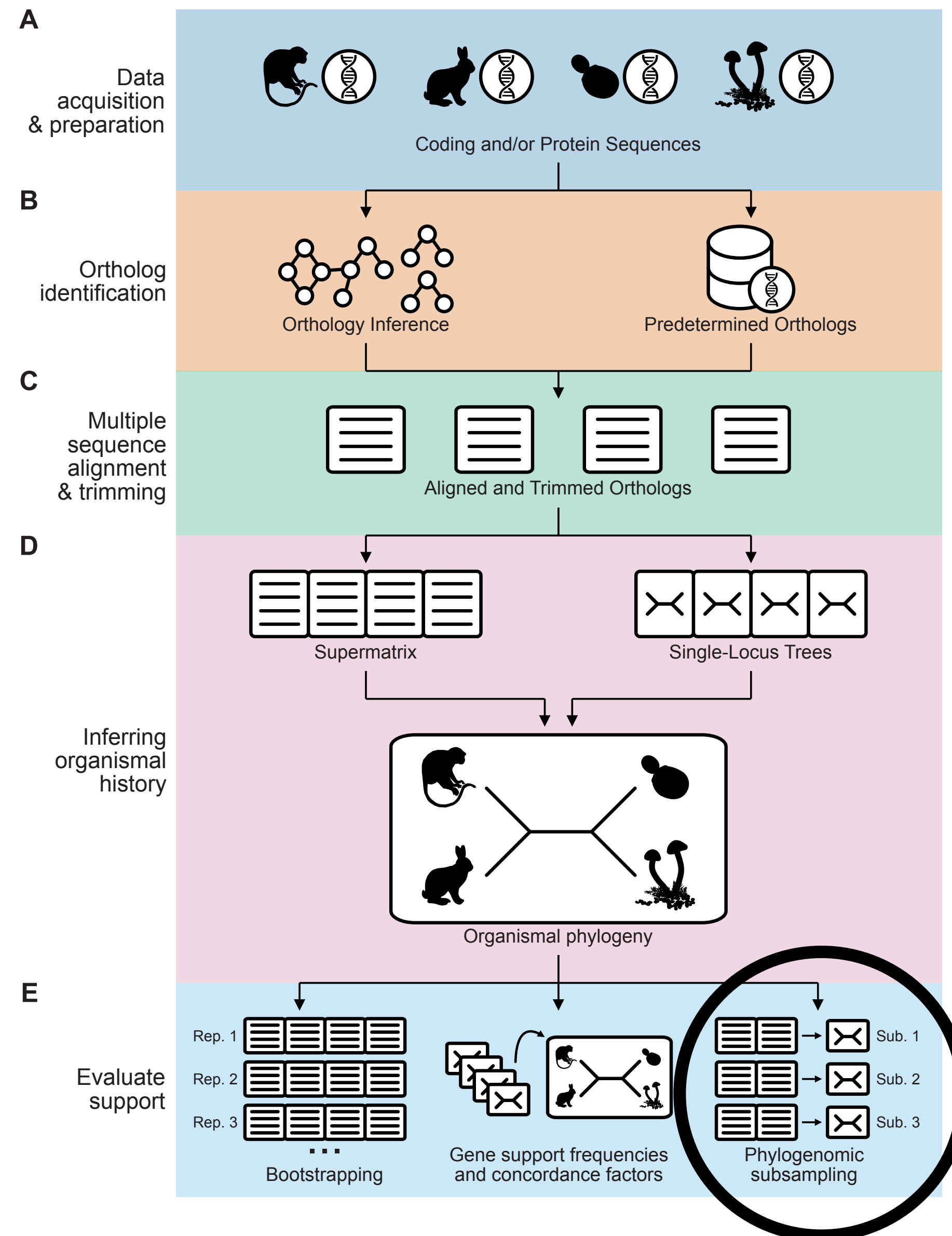


- Introduction
- Inferring genetic networks from phylogenies
- **Phylogenomic subsampling**
- Misc. notes before the tutorial

Facilitating phylogenomic workflows and beyond



Facilitating phylogenomic workflows and beyond



Phylogenomics doesn't solve everything

Review > Trends Genet. 2006 Apr;22(4):225-31. doi: 10.1016/j.tig.2006.02.003.

Epub 2006 Feb 21.

Phylogenomics: the beginning of incongruence?


Olivier Jeffroy ¹, Henner Brinkmann, Frédéric Delsuc, Hervé Philippe

Affiliations + expand

PMID: 16490279 DOI: 10.1016/j.tig.2006.02.003

Free article

Resolving Difficult Phylogenetic Questions: Why More Sequences Are Not Enough

Hervé Philippe , Henner Brinkmann, Dennis V. Lavrov, D. Timothy J. Littlewood, Michael Manuel, Gert Wörheide, Denis Baurain

Incongruence is to be celebrated!

nature reviews genetics

[Explore content](#) ▾ [About the journal](#) ▾ [Publish with us](#) ▾

[nature](#) > [nature reviews genetics](#) > [review articles](#) > [article](#)

Review Article | Published: 27 June 2023

Incongruence in the phylogenomics era

[Jacob L. Steenwyk](#), [Yuanning Li](#), [Xiaofan Zhou](#), [Xing-Xing Shen](#) & [Antonis Rokas](#) 

[Nature Reviews Genetics](#) **24**, 834–850 (2023) | [Cite this article](#)

8371 Accesses | **69** Altmetric | [Metrics](#)

Phylogenomic subsampling, in brief

Phylogenomic subsampling, in brief

1. Unstable bipartitions will be sensitive to gene/taxon/site selection

Phylogenomic subsampling, in brief

1. Unstable bipartitions will be sensitive to gene/taxon/site selection
2. Subsample the full data matrix and reinfer the species tree using fewer (but typically still several dozen to hundreds of genes)

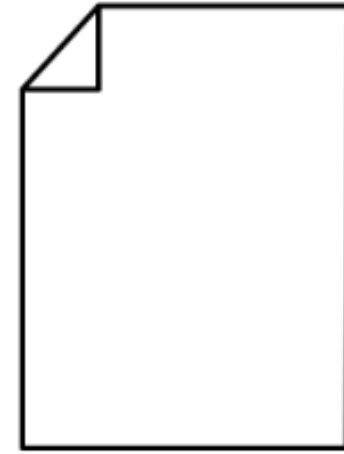
Phylogenomic subsampling, in brief

1. Unstable bipartitions will be sensitive to gene/taxon/site selection
2. Subsample the full data matrix and reinfer the species tree using fewer (but typically still several dozen to hundreds of genes)
3. Compare resulting phylogenies and determine which bipartition are unstable

Phylogenomic subsampling, in brief

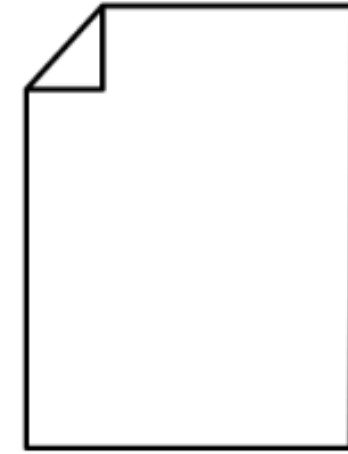
1. Unstable bipartitions will be sensitive to gene/taxon/site selection
2. Subsample the full data matrix and reinfer the species tree using fewer (but typically still several dozen to hundreds of genes)
3. Compare resulting phylogenies and determine which bipartition are unstable
4. Examine potential drivers of incongruence thereafter. Incongruence will be examined in a later lab

Phylogenetic subsampling



Complete
phylogenomic
data matrix

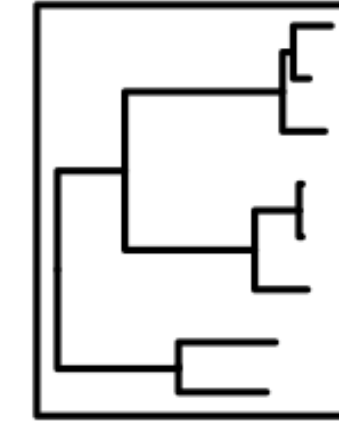
Phylogenetic subsampling



Complete
phylogenomic
data matrix

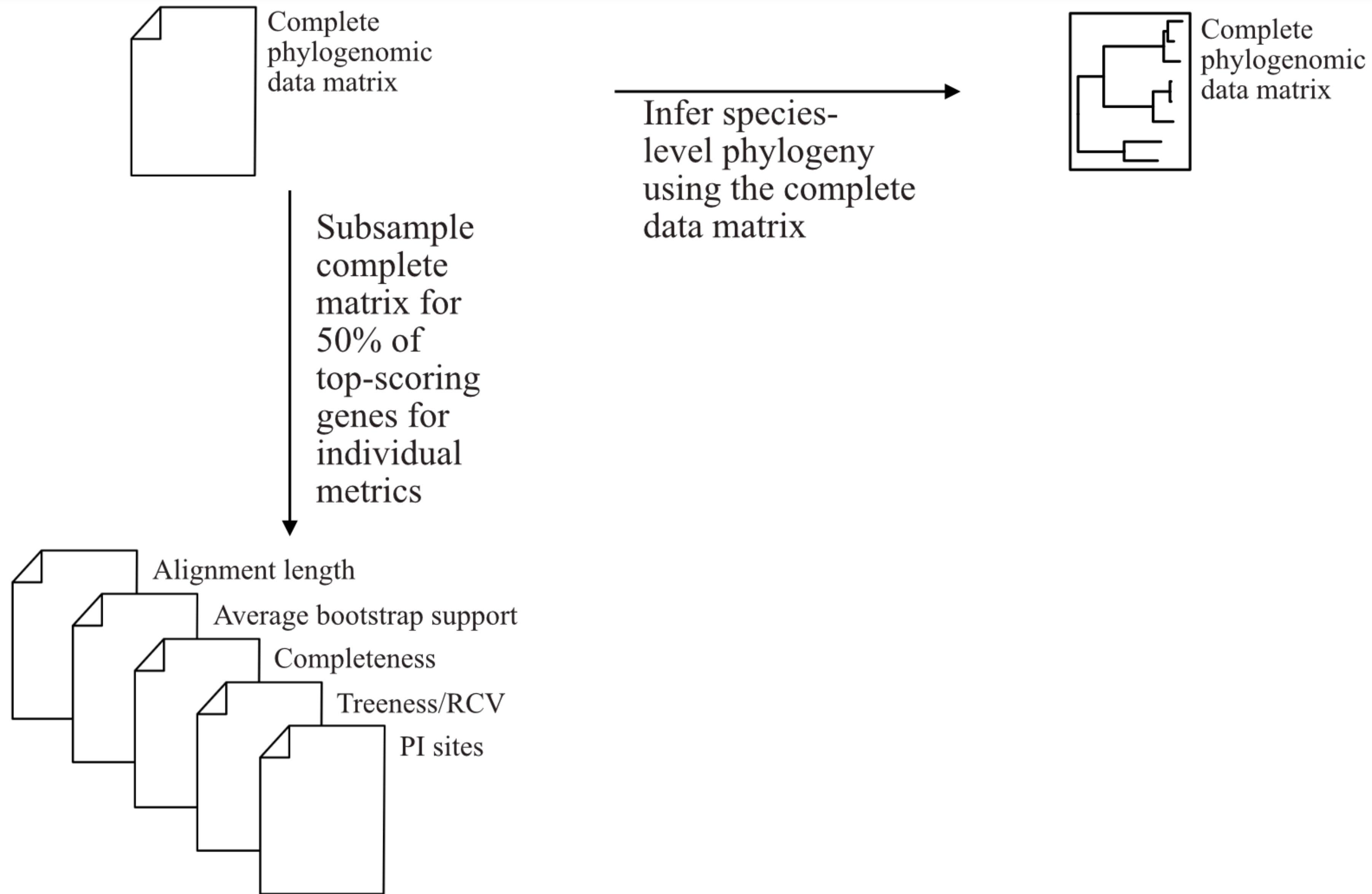


Infer species-
level phylogeny
using the complete
data matrix

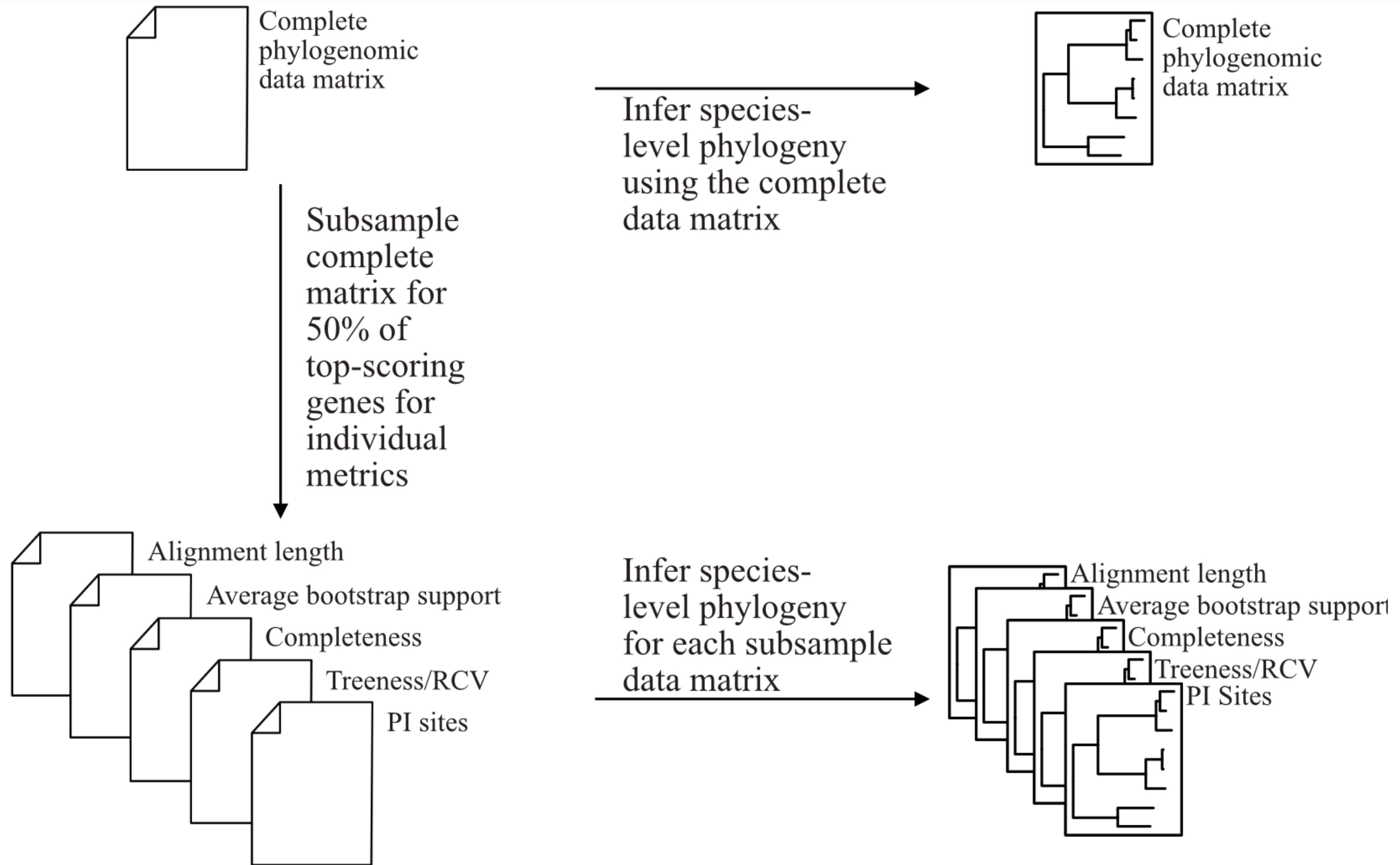


Complete
phylogenomic
data matrix

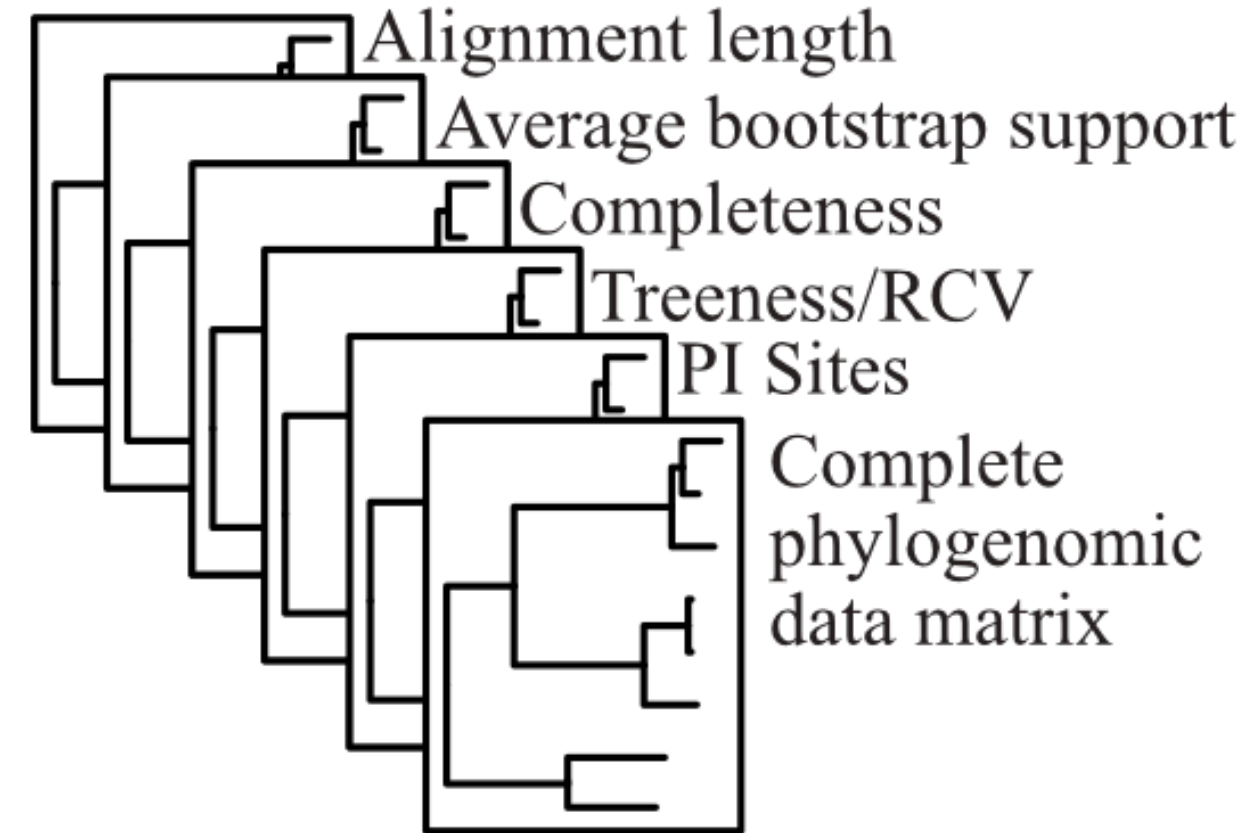
Phylogenetic subsampling



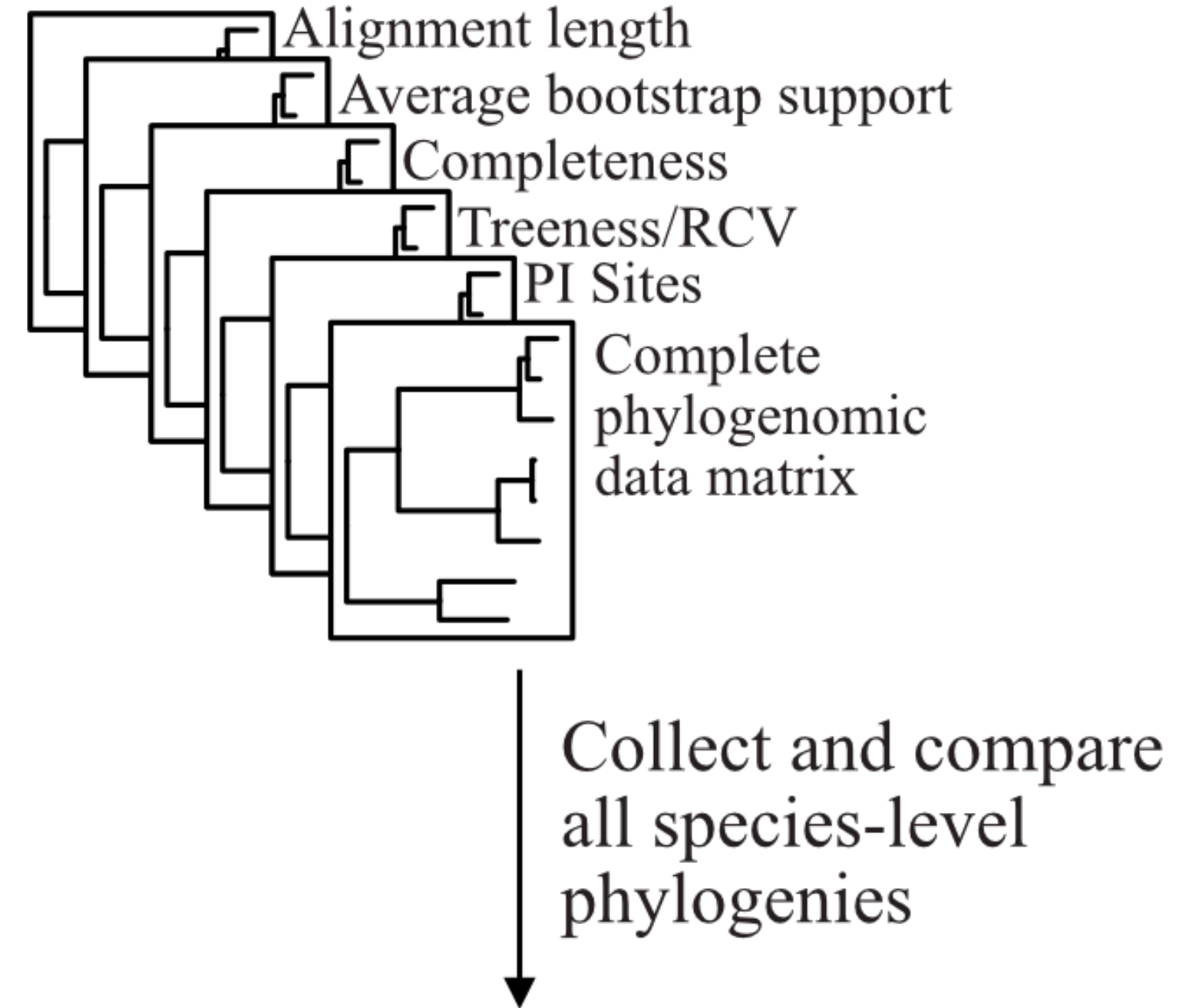
Phylogenetic subsampling



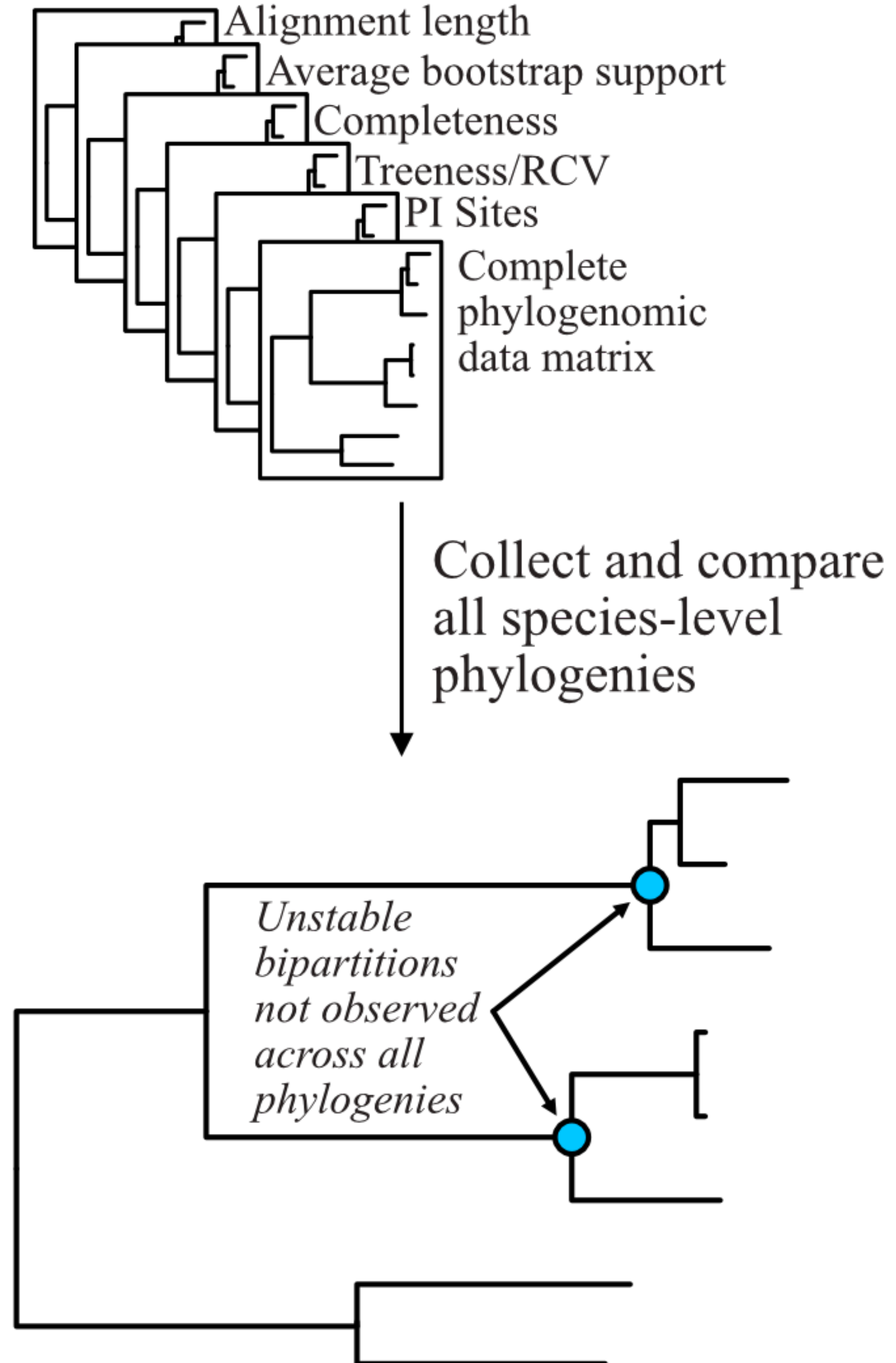
Phylogenetic subsampling



Phylogenetic subsampling



Phylogenetic subsampling



Metrics that capture phylogenetic signal

1. Alignment length
2. Alignment length with no gaps
3. GC content (for NTs)
4. Pairwise identity
5. # of parsimony informative sites
6. # of variable sites
7. Relative composition variability
8. Average bootstrap support value
9. Degree of violation of a molecular clock
10. Evolutionary rate
11. Long branch score
12. Treeness
13. Saturation
14. Treeness / RCV
15. RCVT
16. Compositional bias per site
17. Evolutionary rate per site

Metrics that capture phylogenetic signal

- 1. Alignment length**
2. Alignment length with no gaps
3. GC content (for NTs)
4. Pairwise identity
5. # of parsimony informative sites
6. # of variable sites
- 7. Relative composition variability**
8. Average bootstrap support value
9. Degree of violation of a molecular clock
10. Evolutionary rate
- 11. Long branch score**
12. Treeness
- 13. Saturation**
- 14. Treeness / RCV**
- 15. RCVT**
- 16. Compositional bias per site**
- 17. Evolutionary rate per site**

Phylogenetic signal across genes

- 1. Alignment length**
2. Alignment length with no gaps
3. GC content (for NTs)
4. Pairwise identity
5. # of parsimony informative sites
6. # of variable sites
- 7. Relative composition variability**
8. Average bootstrap support value
9. Degree of violation of a molecular clock
10. Evolutionary rate
- 11. Long branch score**
12. Treeness
- 13. Saturation**
- 14. Treeness / RCV**

Alignment length

```
>sp1
A C G T A G C G - T C G A T C
>sp2
A C G T - G C G A T C G A T C
>sp3
A C G T A G C - A T C G A T C
>sp4
A C G T A G C G A T C G A T G
>sp5
A C - - A G C G A T C G A T C
>sp6
A C G T A G C G A - - - A T C
```

The length of this alignment is 15 sites

Alignment length

```
>sp1
A C G T A G C G - T C G A T C
>sp2
A C G T - G C G A T C G A T C
>sp3
A C G T A G C - A T C G A T C
>sp4
A C G T A G C G A T C G A T G
>sp5
A C - - A G C G A T C G A T C
>sp6
A C G T A G C G A - - - A T C
```

Higher values are better!

The length of this alignment is 15 sites

Relative composition variability

- Average variability in the sequence composition among taxa in an MSA

Relative composition variability

- Average variability in the sequence composition among taxa in an MSA
- Evaluates potential composition biases
 - violate assumptions of site composition homogeneity in standard substitution models

Relative composition variability

- Average variability in the sequence composition among taxa in an MSA
- Evaluates potential composition biases
 - violate assumptions of site composition homogeneity in standard substitution models

$$\sum_{i=1}^c \sum_{j=1}^n \frac{|c_{ij} - \bar{c}_i|}{s \times n}$$

Relative composition variability

- Average variability in the sequence composition among taxa in an MSA
- Evaluates potential composition biases
 - violate assumptions of site composition homogeneity in standard substitution models

$$\sum_{i=1}^c \sum_{j=1}^n \frac{|c_{ij} - \bar{c}_i|}{s \times n}$$

- c is the number of different character states per sequence type
- n is the number of taxa in an MSA
- s is the number of sites in an MSA

Relative composition variability

>Seq 1
MKGATTLAK

>Seq 2
MK-AITLAK

>Seq 3
MKGATT--K

>Seq 4
MK-AITLA-

RCV = 0.375

Relative composition variability

>Seq_1
MKGATTLAK

>Seq_2
MK-AITLAK

>Seq_3
MKGATT--K

>Seq_4
MK-AITLA-

RCV = 0.375



>Seq_1
MKTTTTTTT

>Seq_2
MKQQQQQQQ

>Seq_3
MKKKKKKKK

>Seq_4
MKLLLLLLL

RCV = 1.1667

Relative composition variability

Lower compositional bias

```
>Seq_1  
MKGATTLAK
```

```
>Seq_2  
MK-AITLAK
```

```
>Seq_3  
MKGATT--K
```

```
>Seq_4  
MK-AITLA-
```

RCV = 0.375



Higher compositional bias

```
>Seq_1  
MKTTTTTTT
```

```
>Seq_2  
MKQQQQQQQ
```

```
>Seq_3  
MKKKKKKKK
```

```
>Seq_4  
MKLLLLLLL
```

RCV = 1.1667

Relative composition variability

Lower compositional bias

Higher compositional bias

***Lower RCV
values are better***

>Seq_1
MKGATTLAK

>Seq_2
MK-AITLAK

>Seq_3
MKGATT--K

>Seq_4
MK-AITLA-

RCV = 0.375



>Seq_1
MKTTTTTTT

>Seq_2
MKQQQQQQQ

>Seq_3
MKKKKKKKK

>Seq_4
MKLLLLLLL

RCV = 1.1667

Relative composition variability

Lower compositional bias

Higher compositional bias

***Lower RCV
values are better***

>Seq_1
MKGATTLAK

>Seq_2
MK-AITLAK

>Seq_3
MKGATT--K

>Seq_4
MK-AITLA-

RCV = 0.375



>Seq_1
MKTTTTTTT

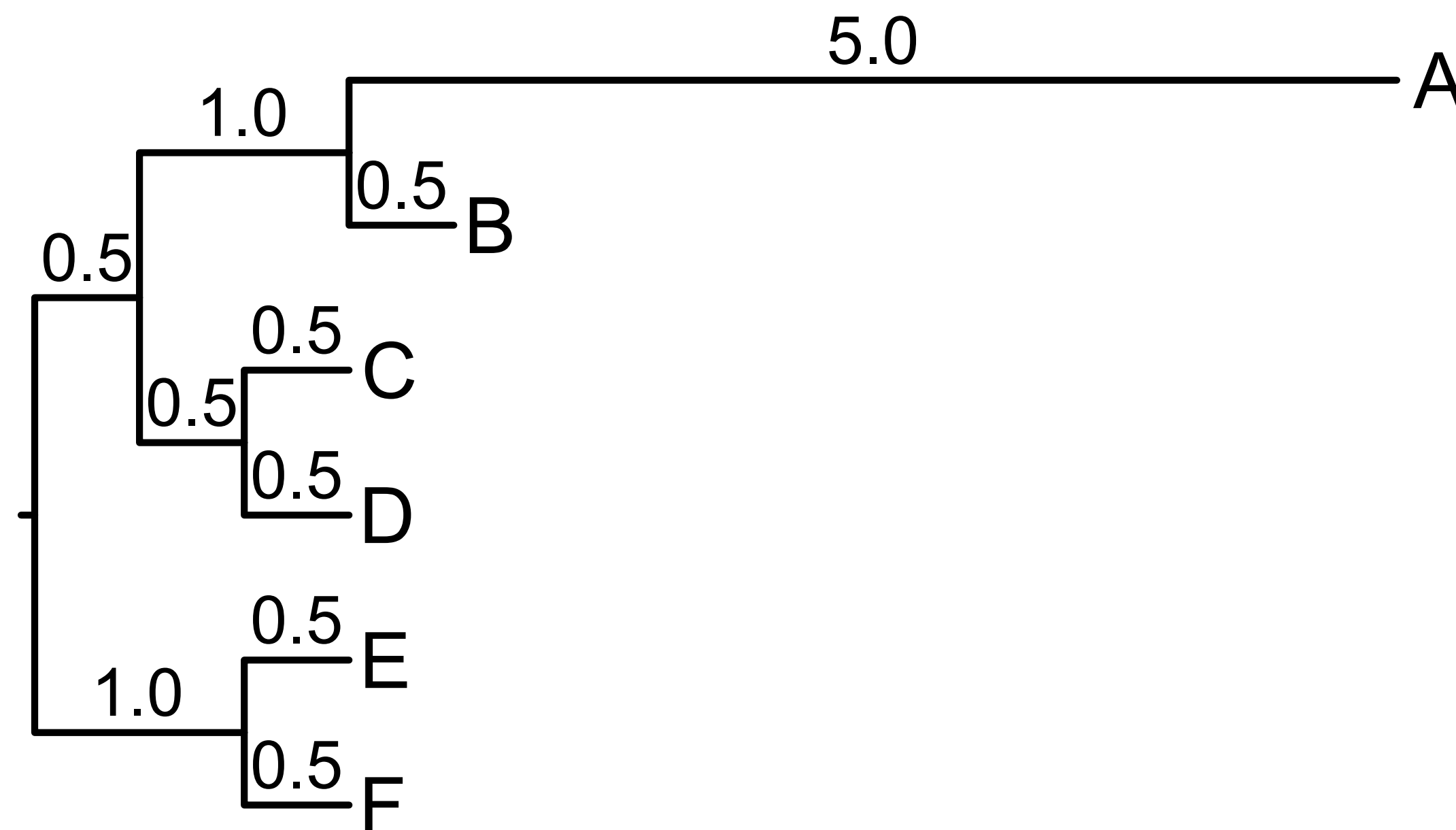
>Seq_2
MKQQQQQQQ

>Seq_3
MKKKKKKKK

>Seq_4
MKLLLLLLL

RCV = 1.1667

Long branch score



Long branch scores

A: 73.17

B: -14.63

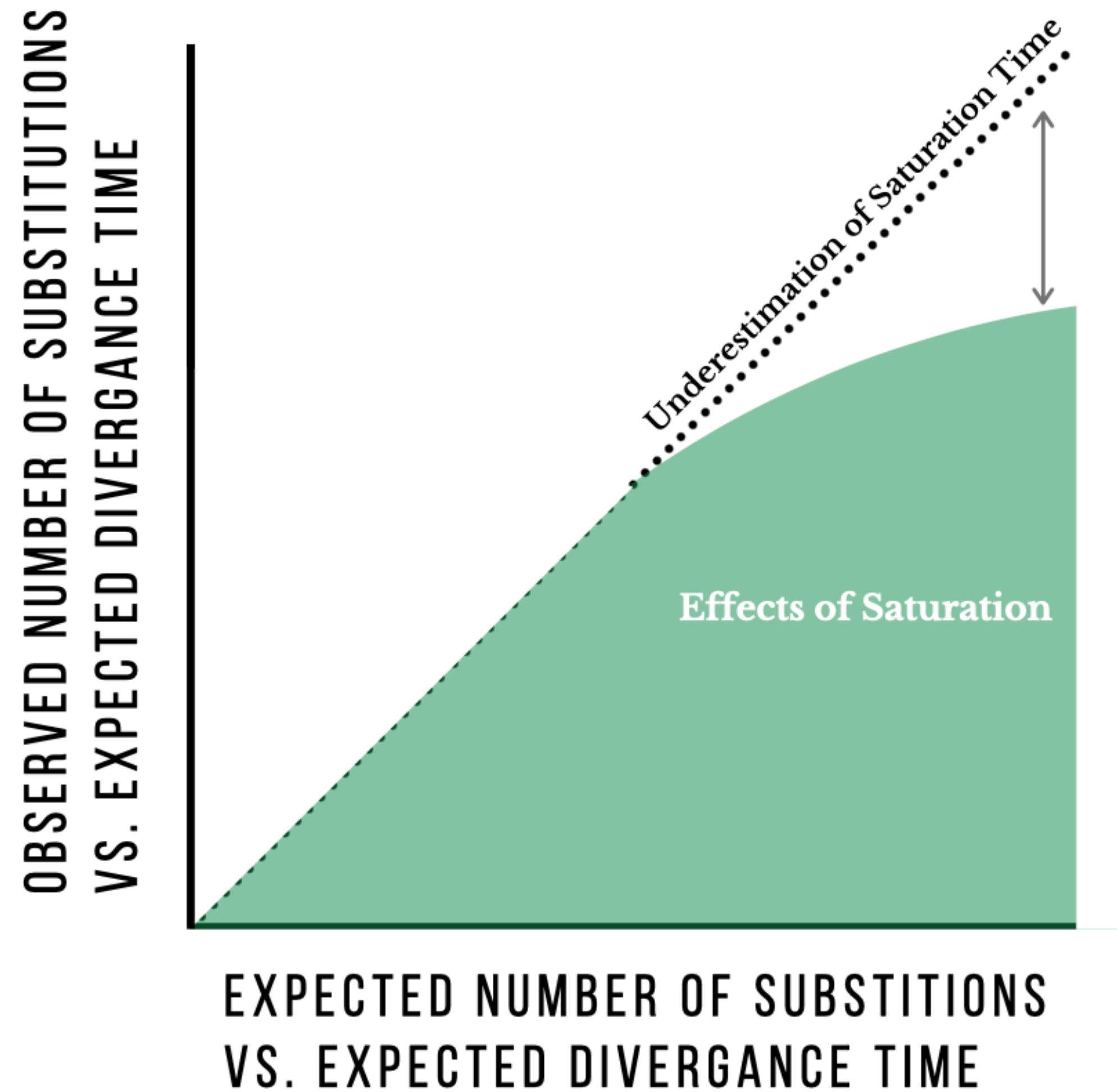
C: -19.51

D: -19.51

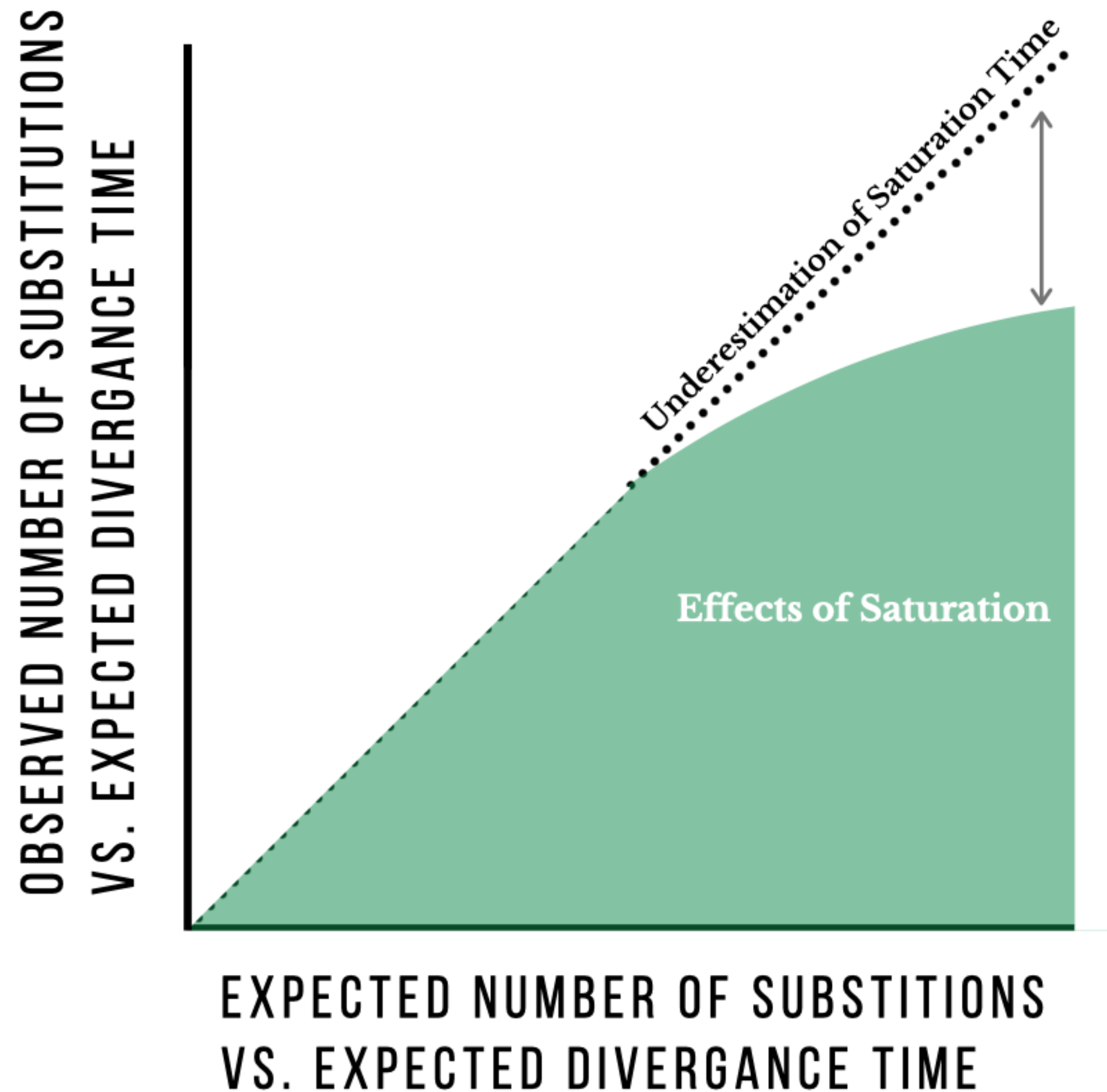
E: -9.76

F: -9.76

Saturation by multiple substitutions

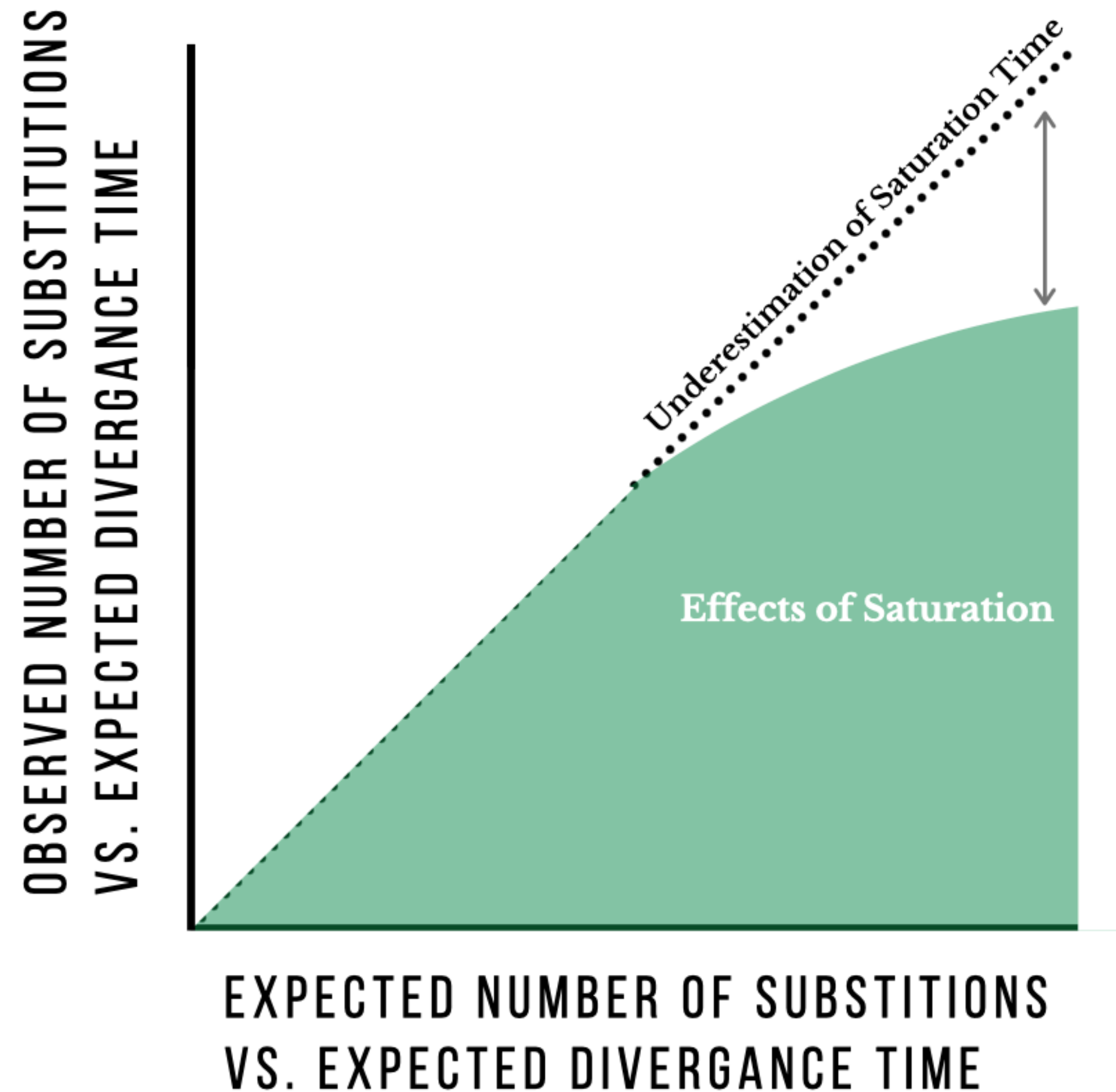


Saturation by multiple substitutions



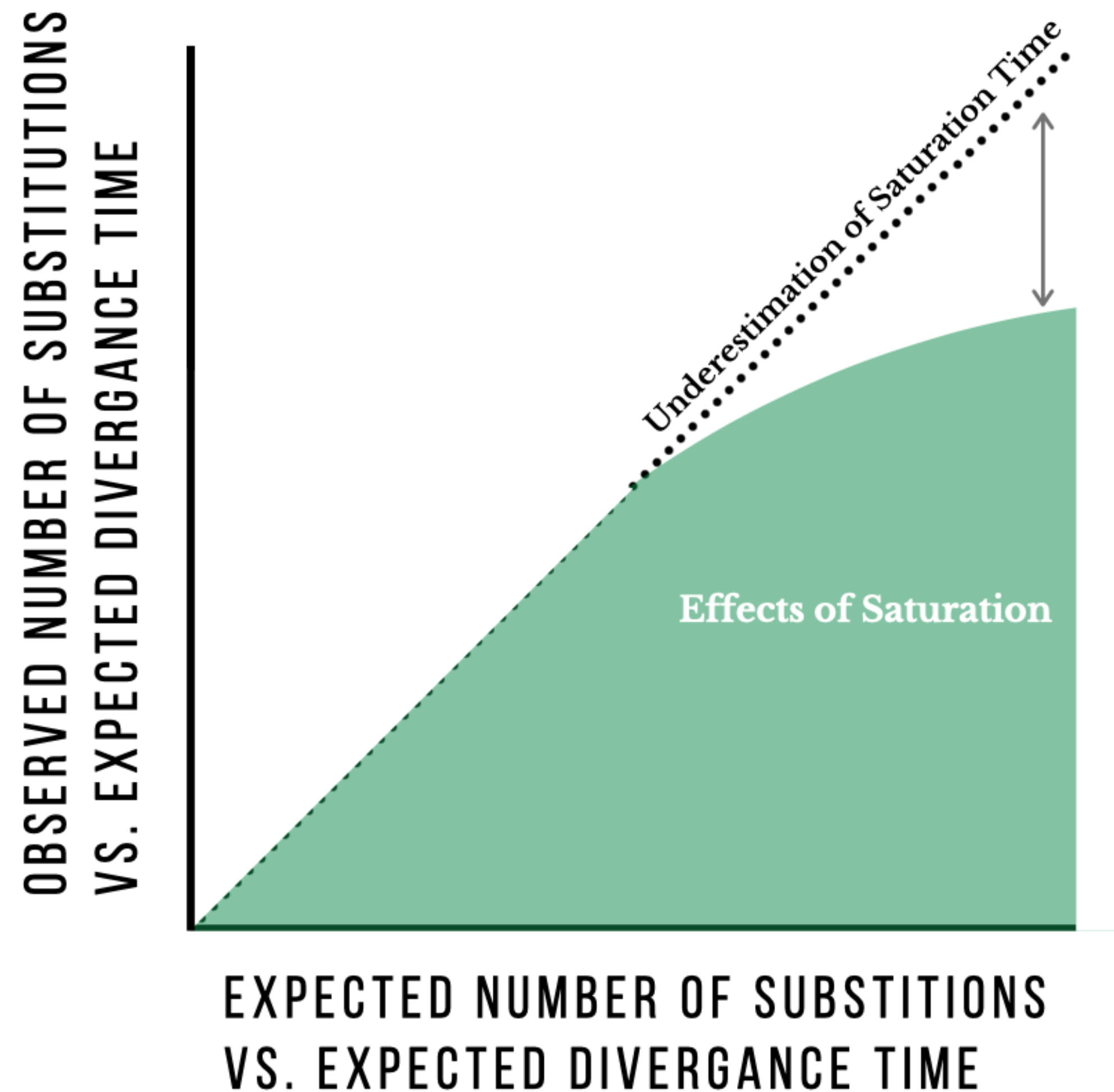
- X-axis can be approximated using phylogenetic distances
 - Tip-to-tip distances in a tree

Saturation by multiple substitutions



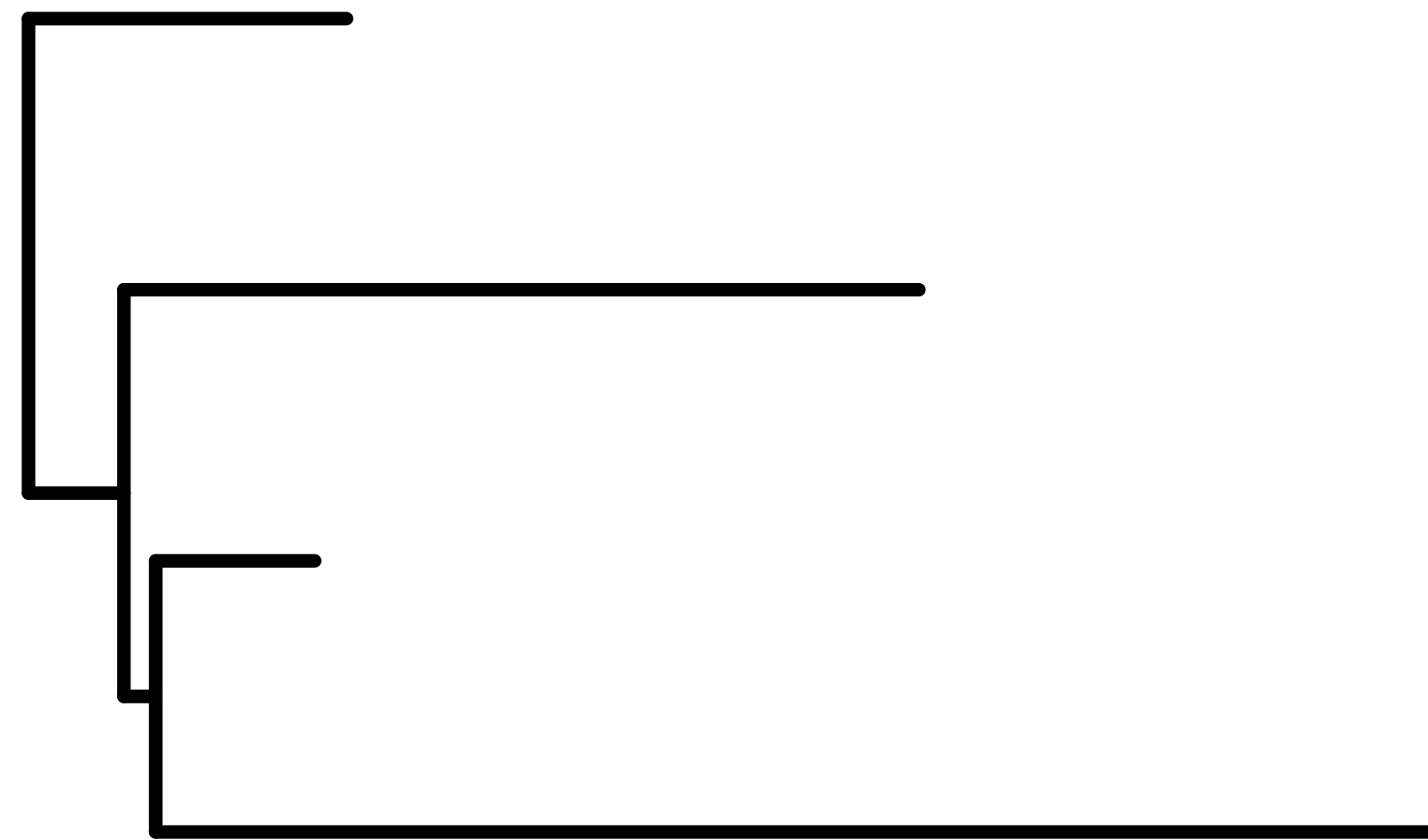
- X-axis can be approximated using phylogenetic distances
 - Tip-to-tip distances in a tree
- Y-axis can be approximated using pairwise identity
 - Distance in an MSA

Saturation by multiple substitutions



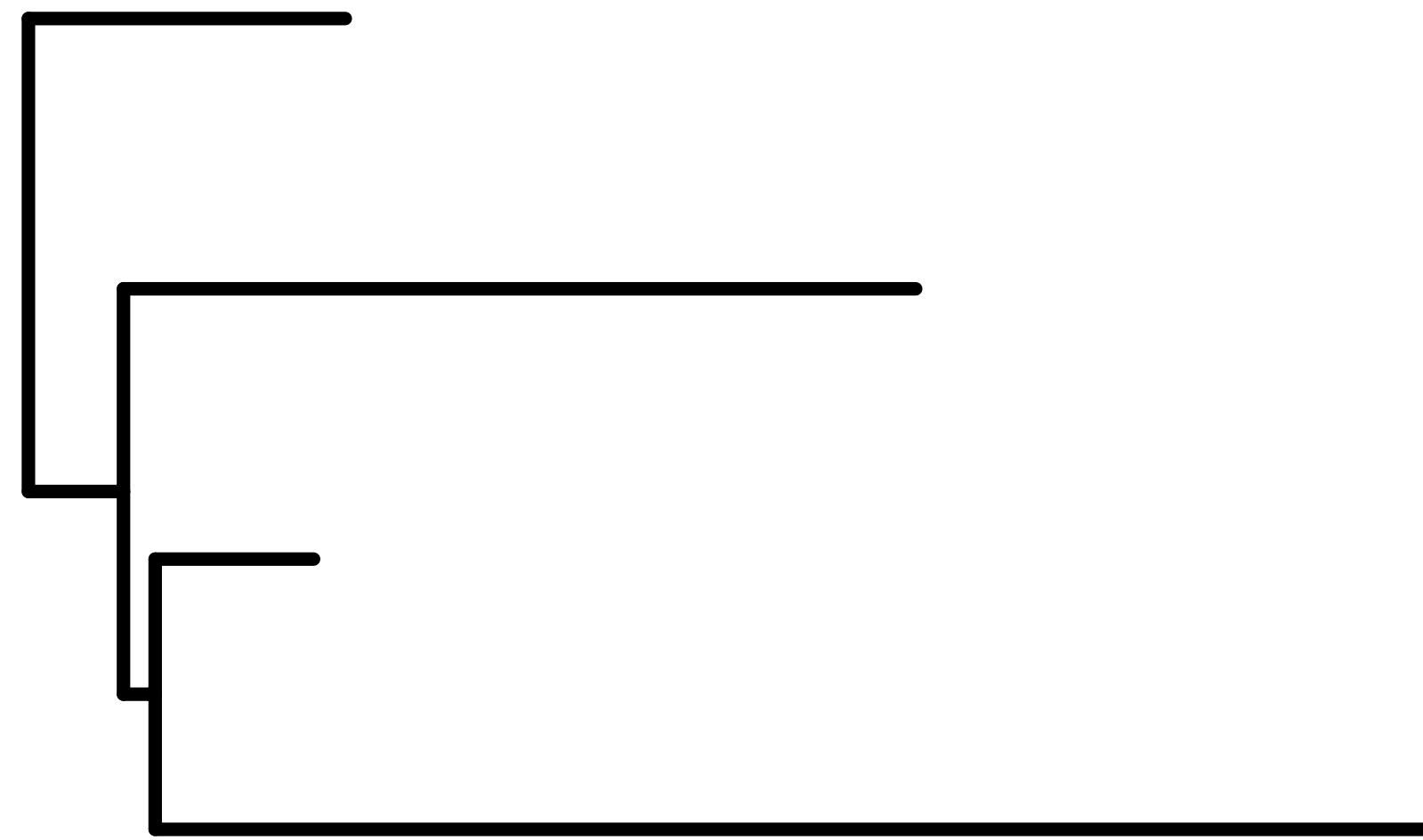
- The closer the slope is to 1, the better.
- PhyKIT reports the slope
- PhyKIT also reports the absolute difference between the slope and 1
 - Thus, the lower the value the better

Treeness

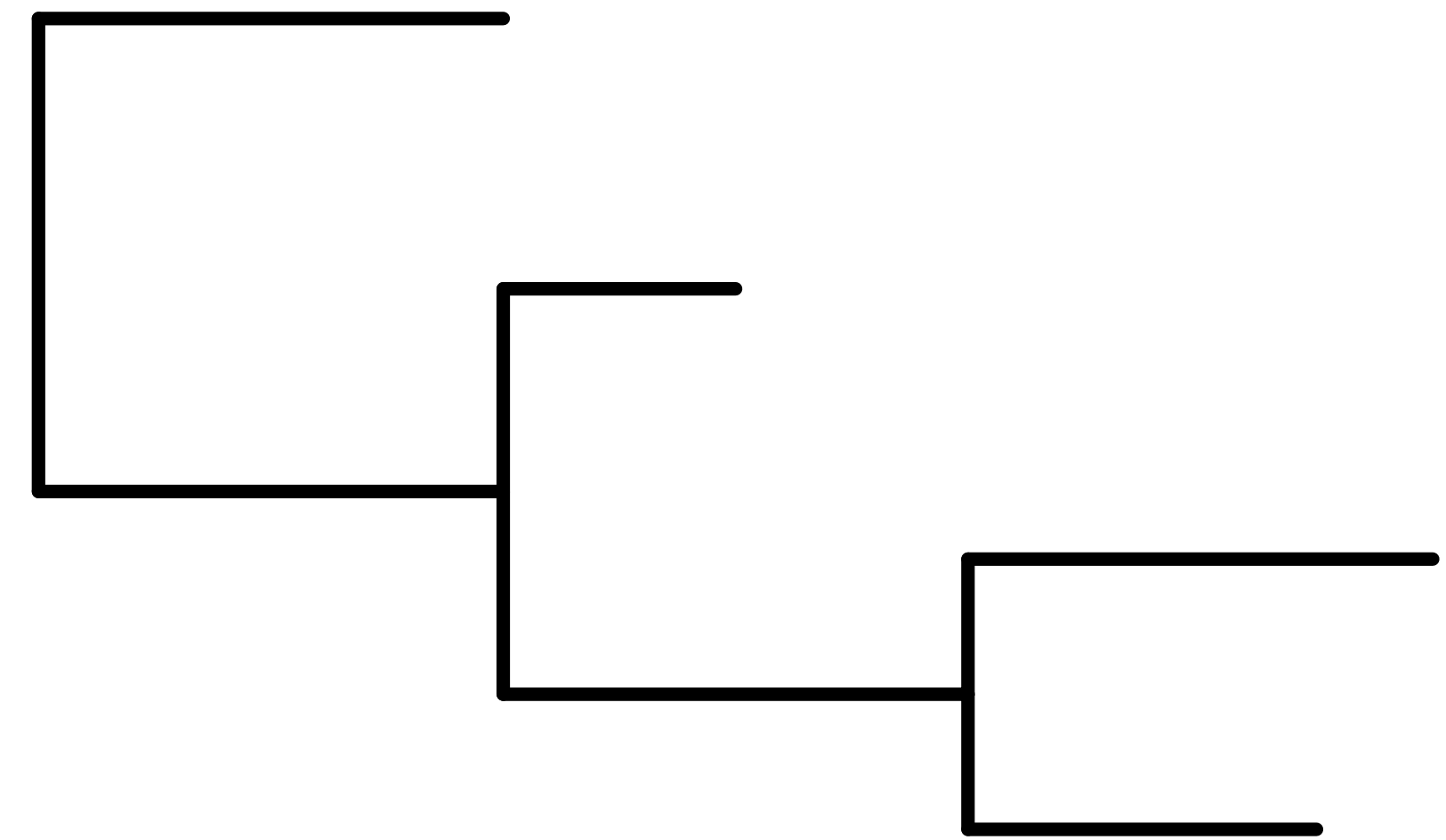


treeness = 0.0476

Treeness



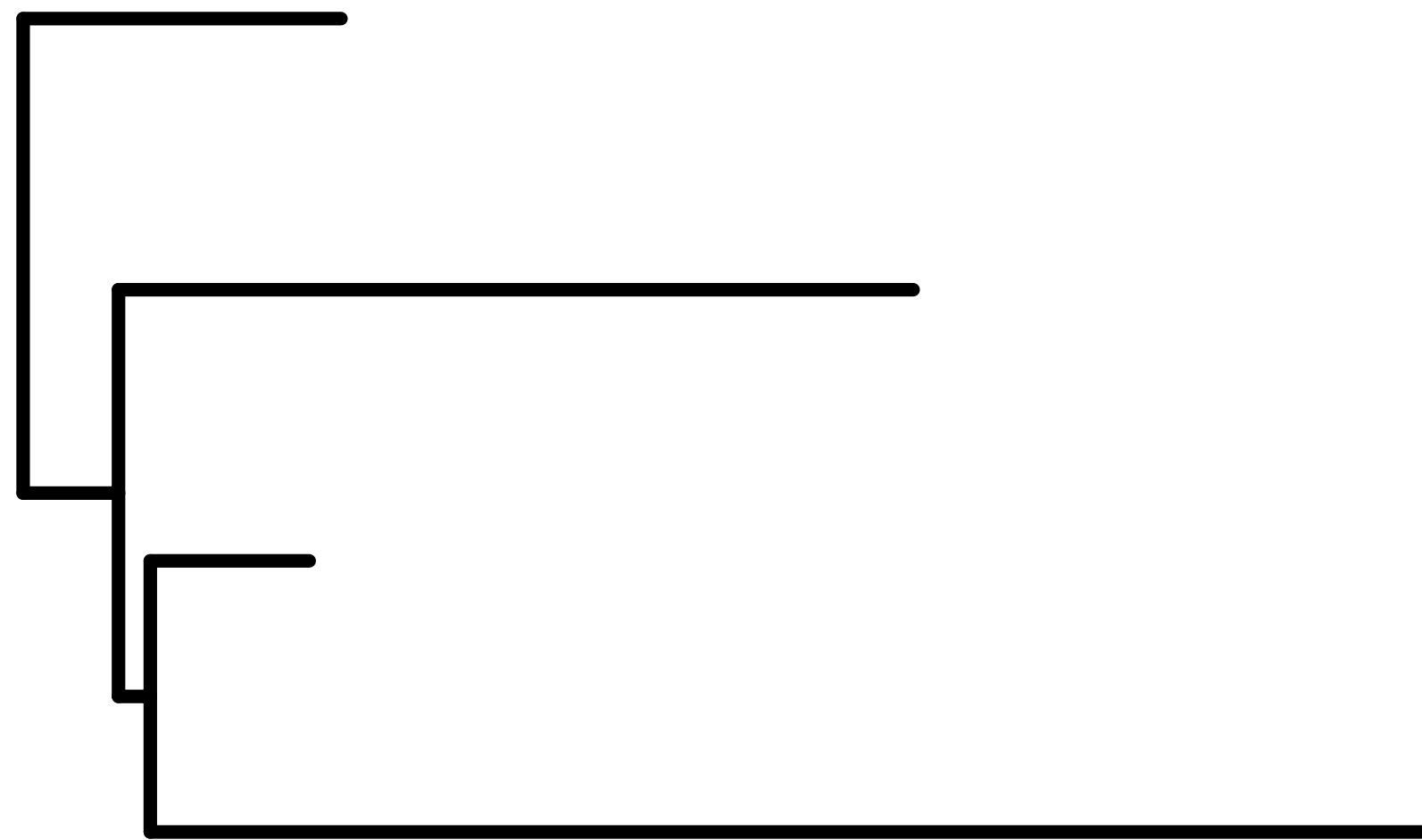
treeness = 0.0476



treeness = 0.381

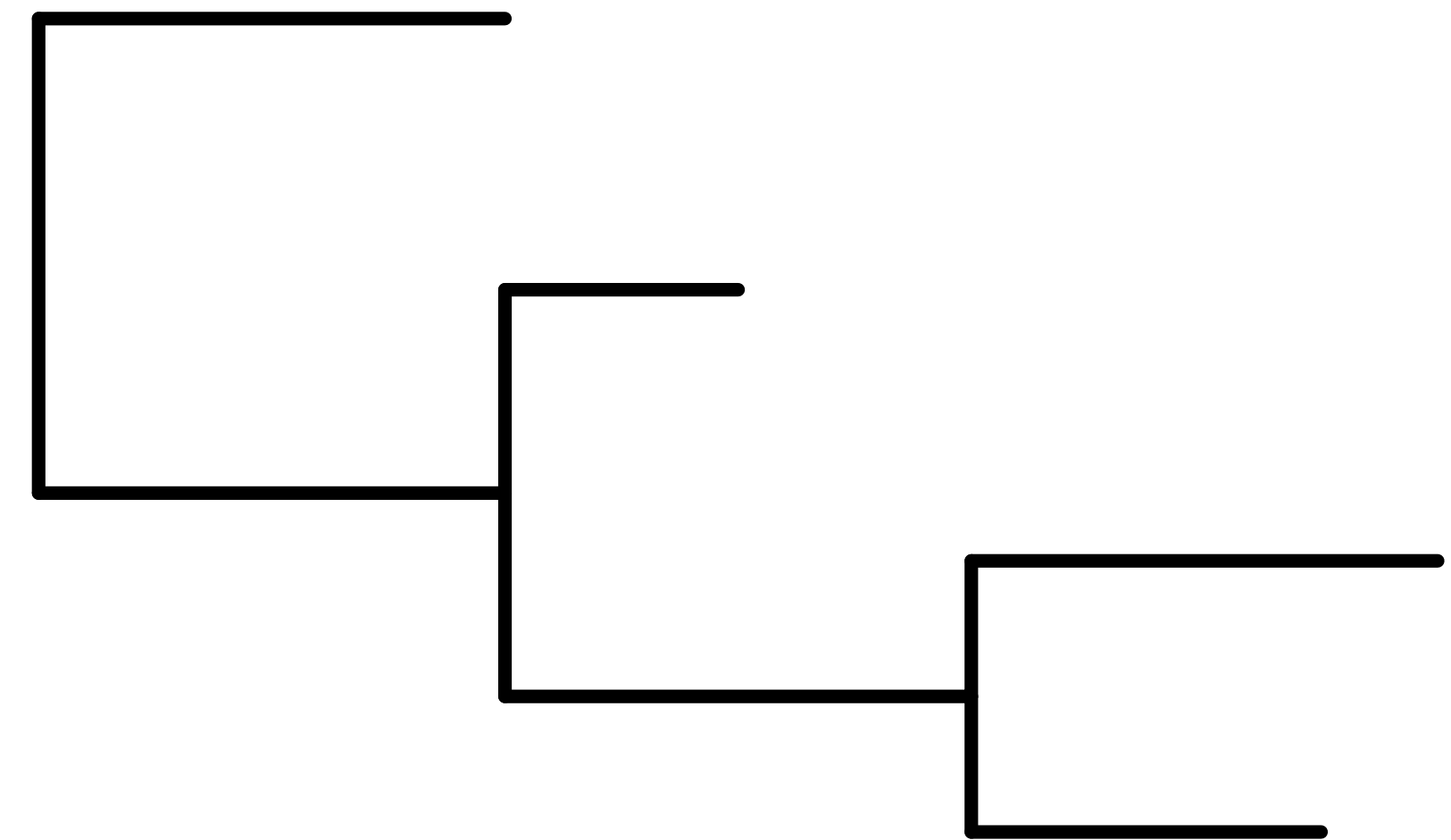
Treeness

Low treeness



treeness = 0.0476

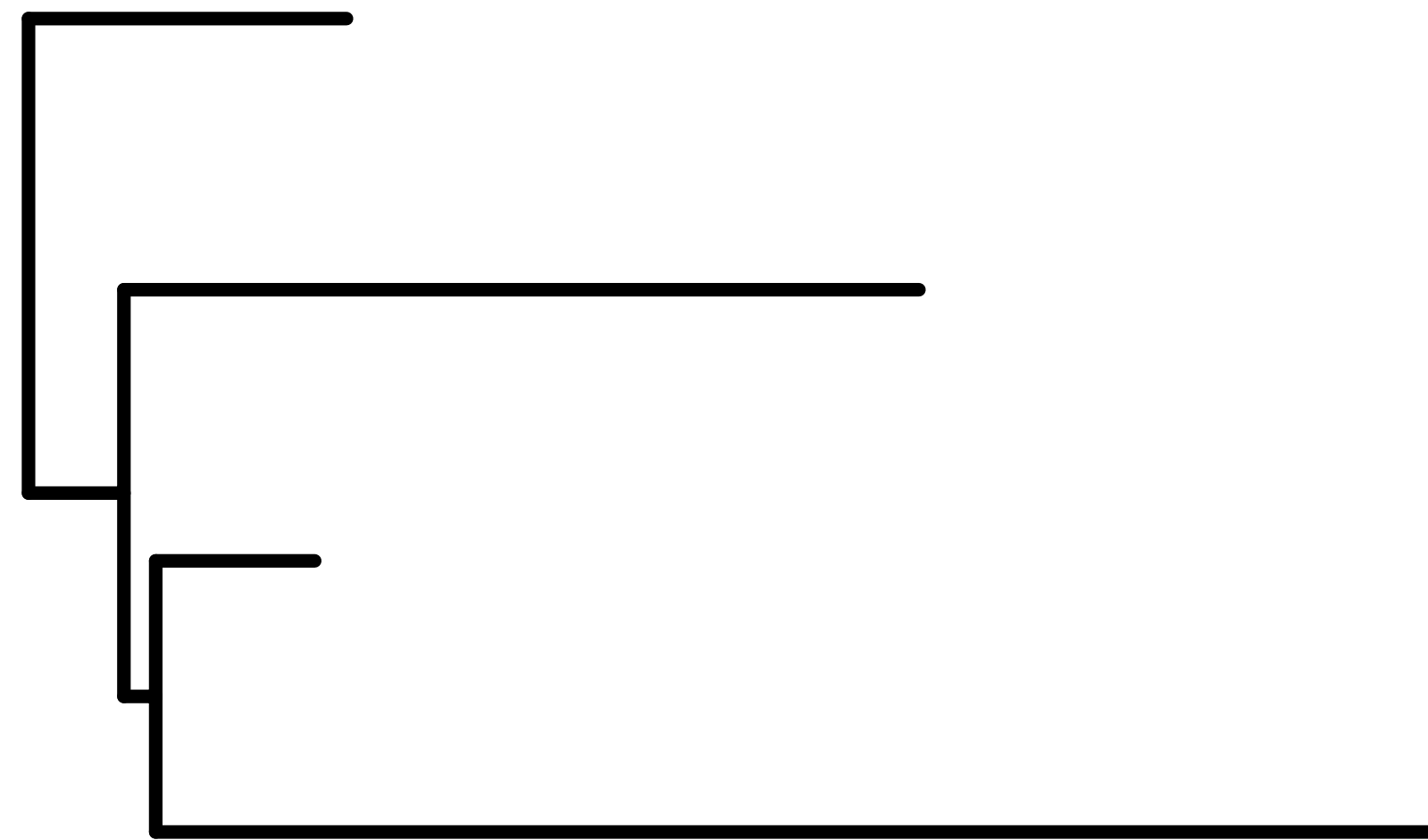
High treeness



treeness = 0.381

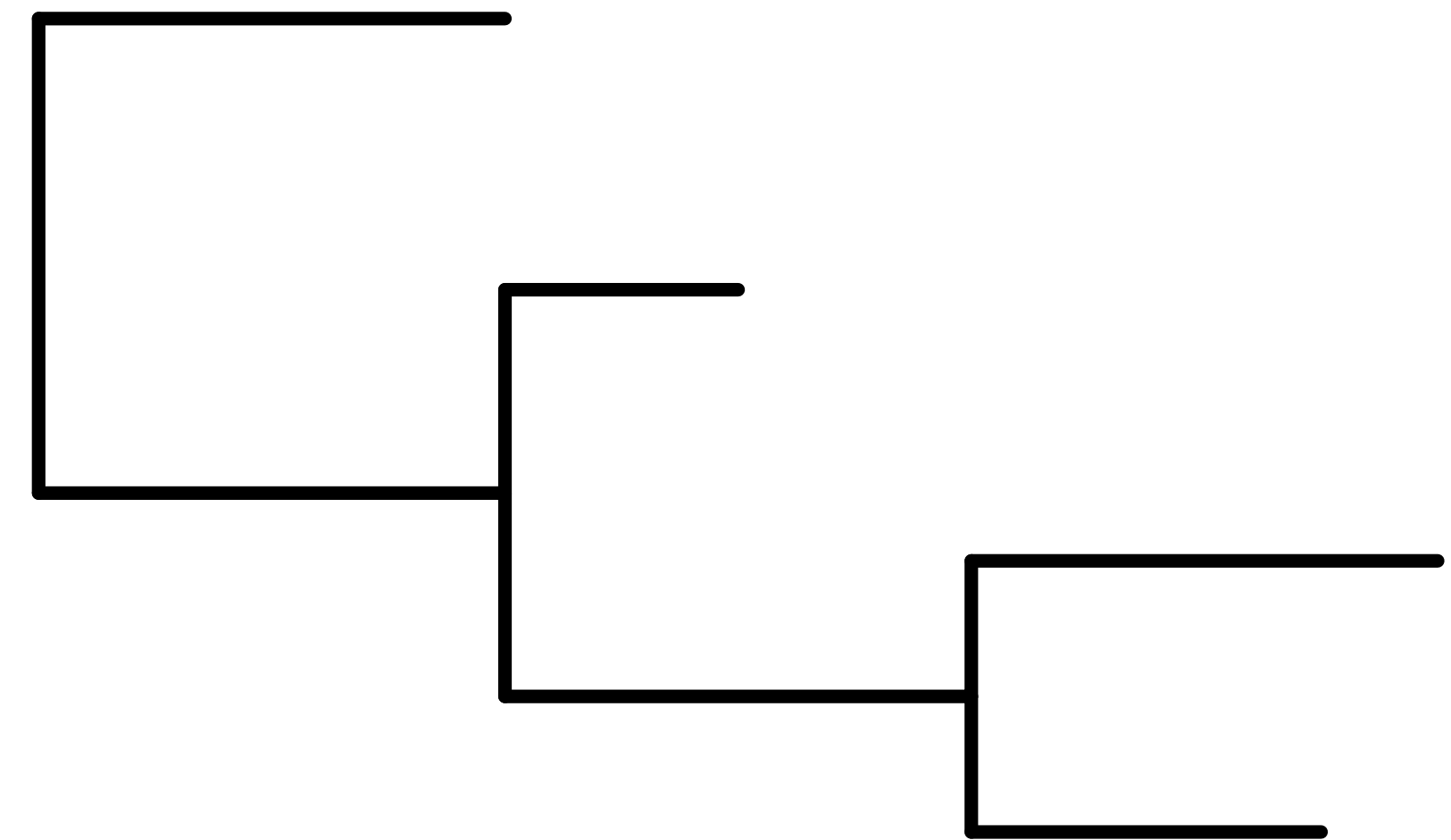
Treeness

Low treeness



treeness = 0.0476

High treeness

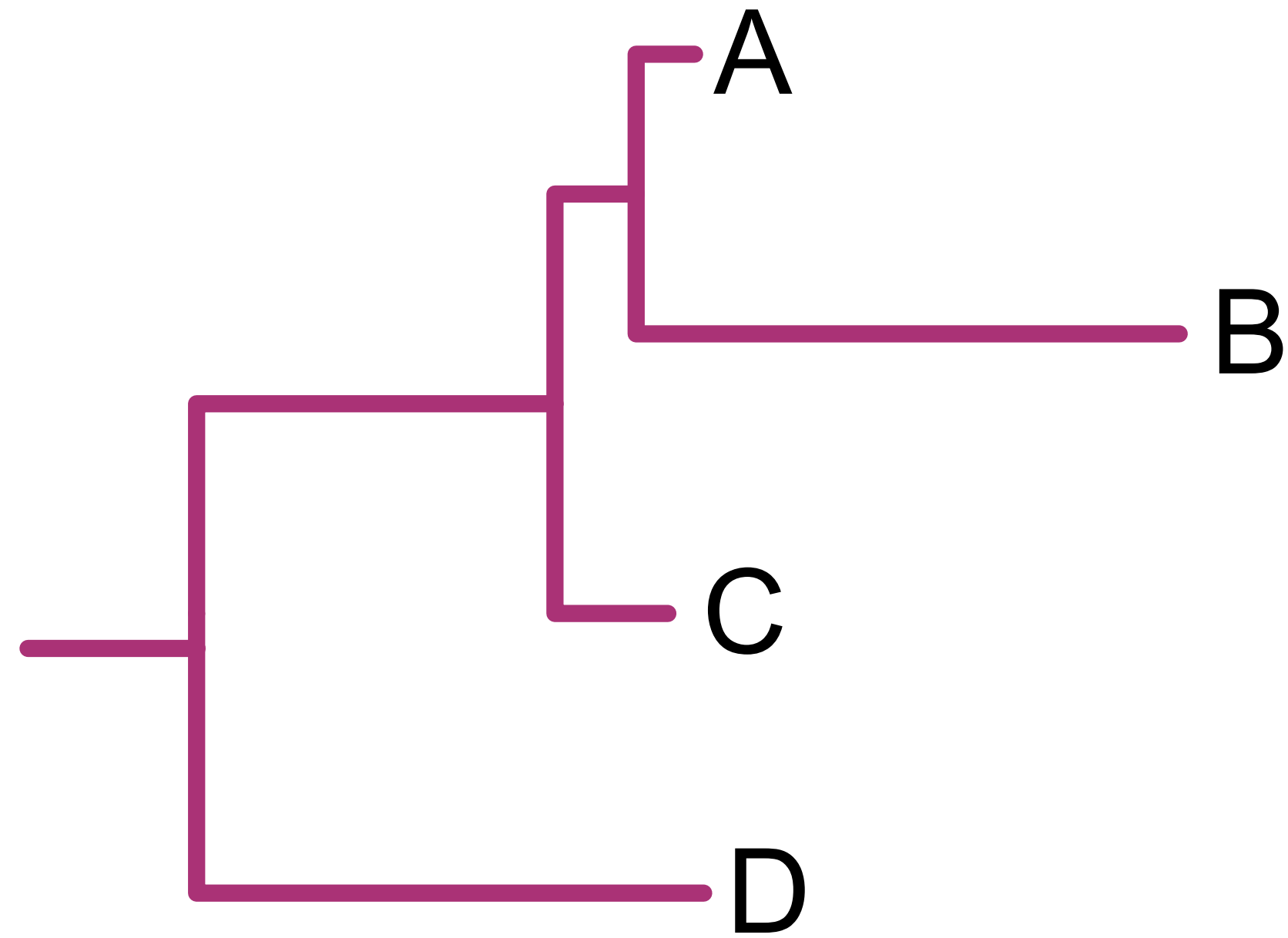


treeness = 0.381

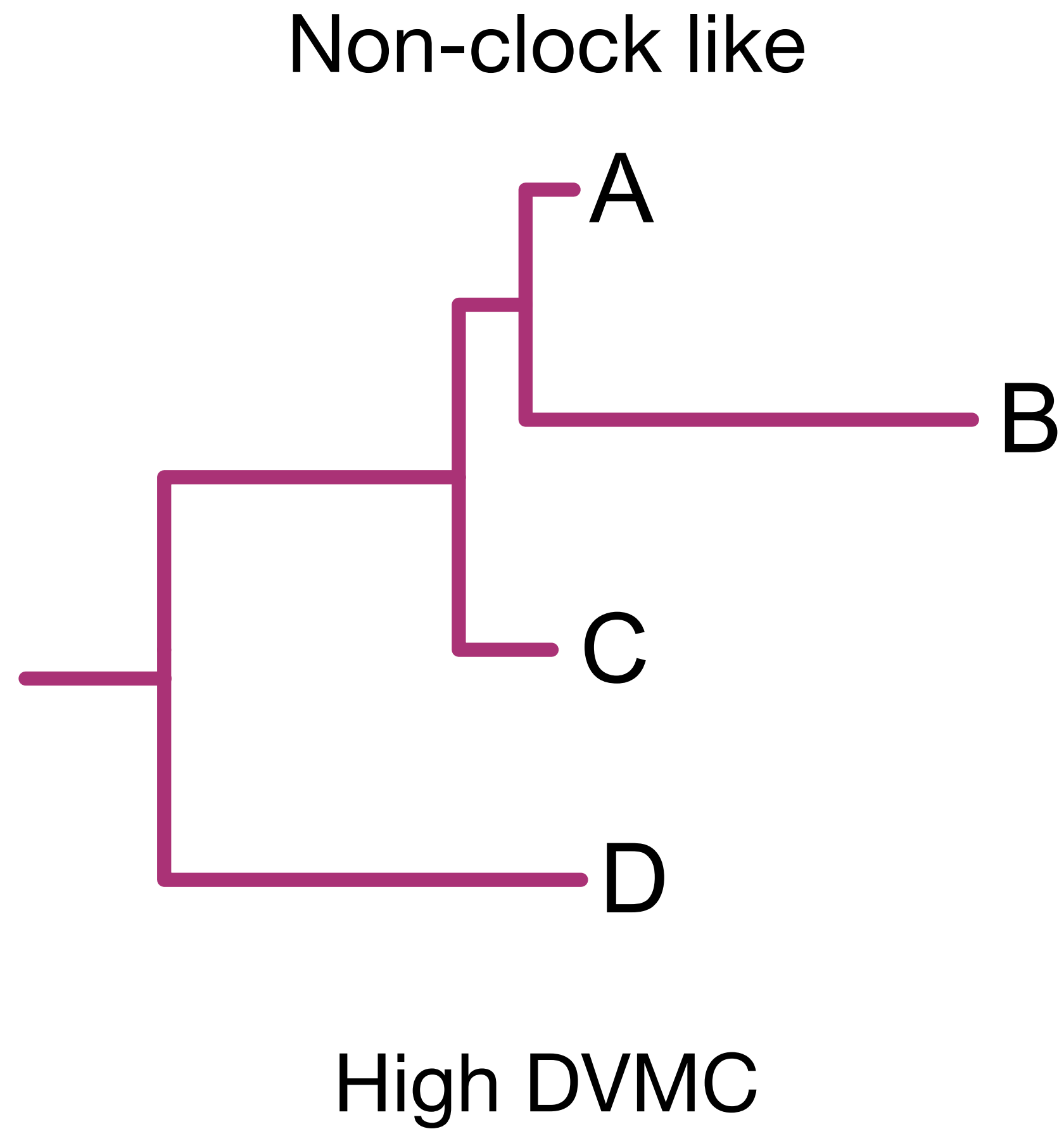
Higher treeness values are better

Degree of violation of a molecular clock

Non-clock like

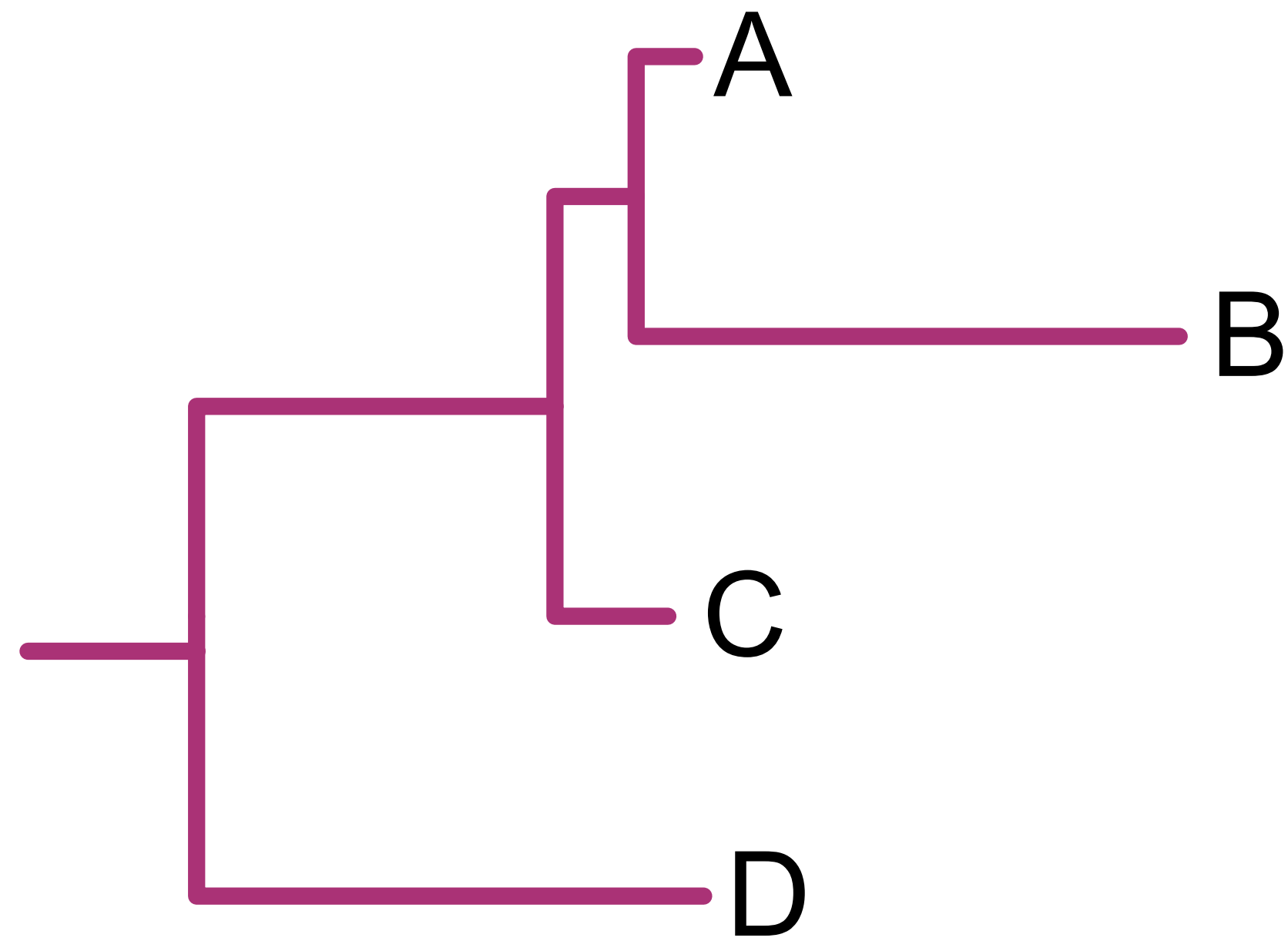


Degree of violation of a molecular clock



Degree of violation of a molecular clock

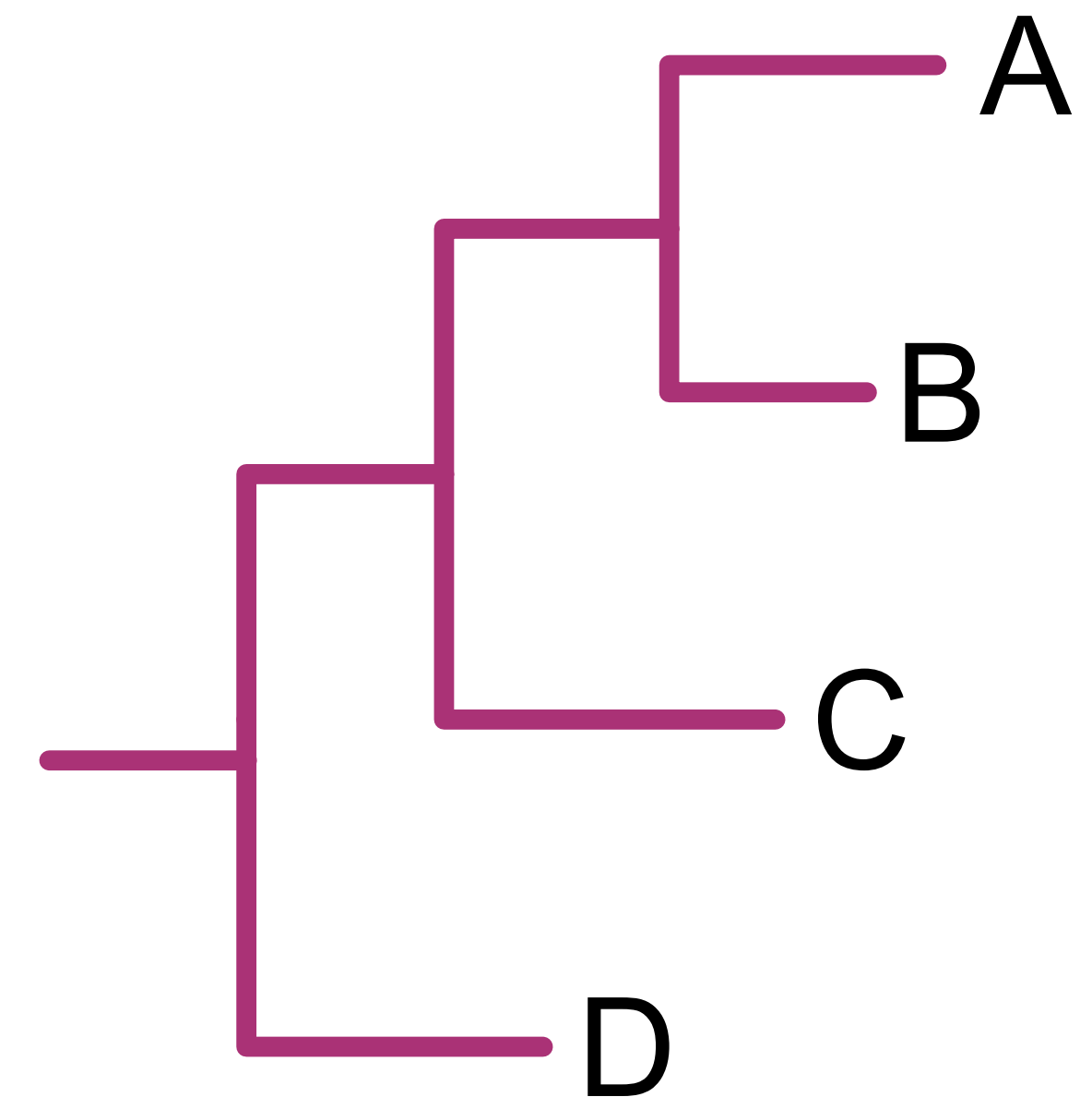
Non-clock like



High DVMC

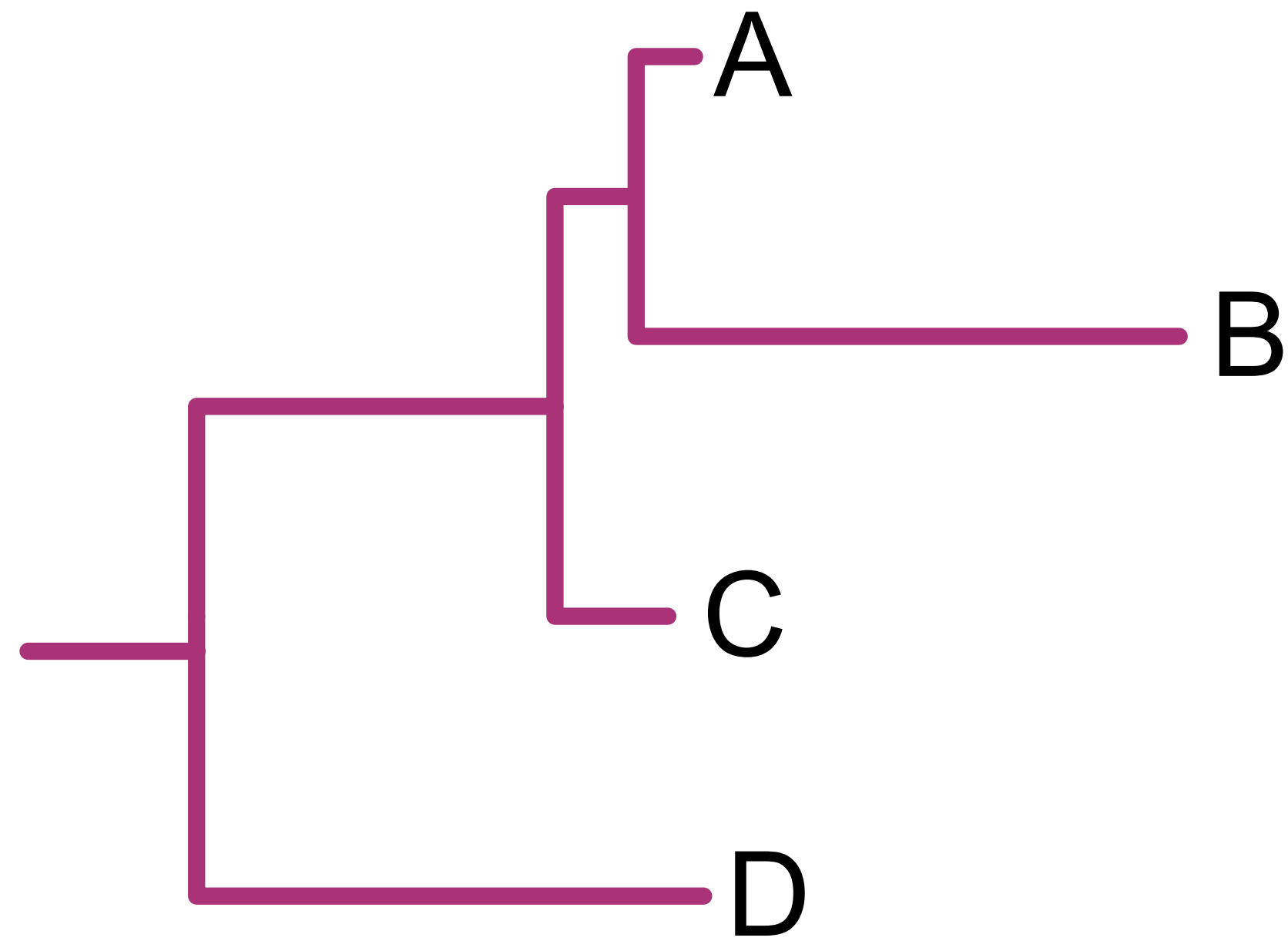


Clock-like



Degree of violation of a molecular clock

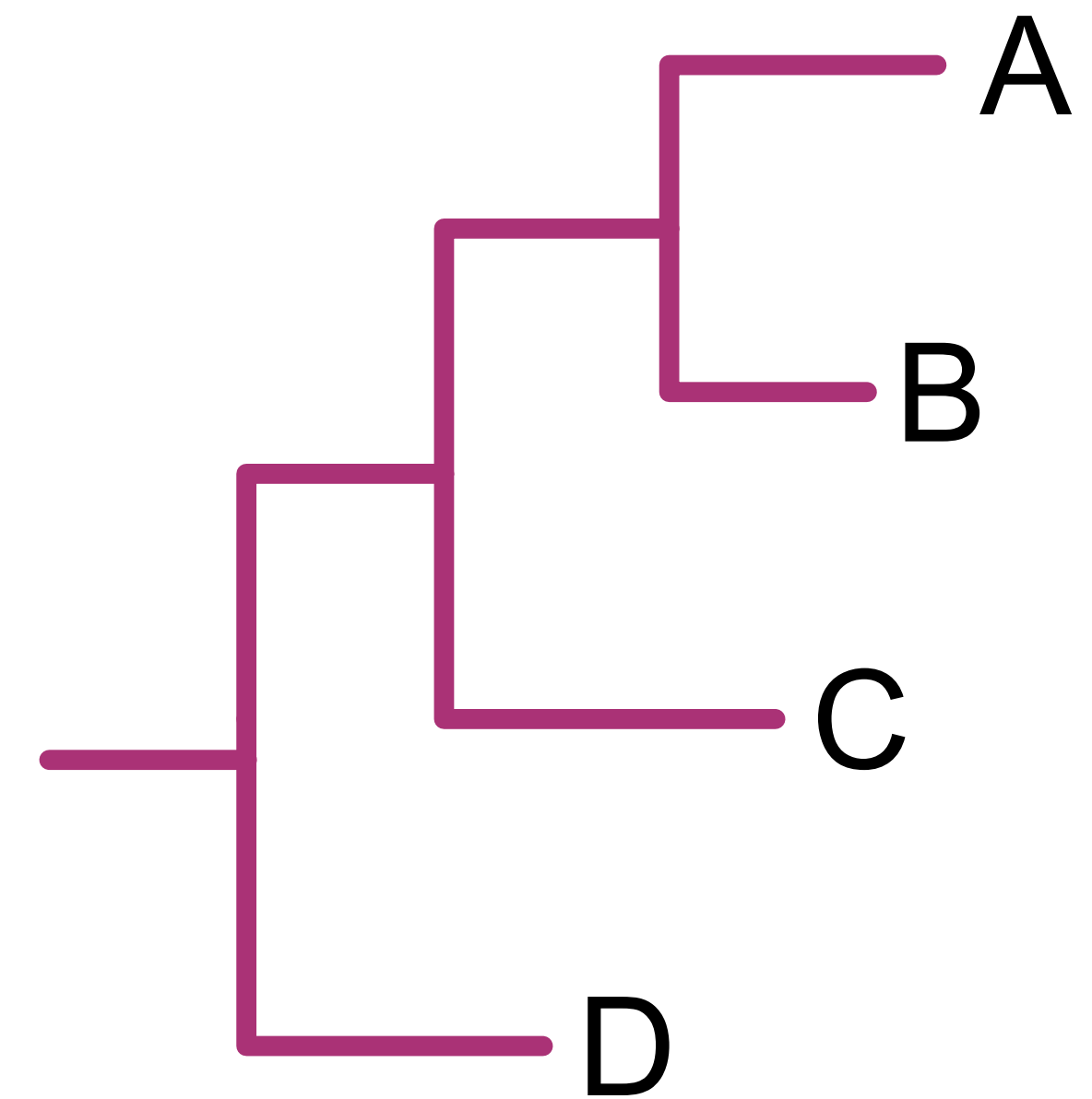
Non-clock like



High DVMC



Clock-like



Low DVMC

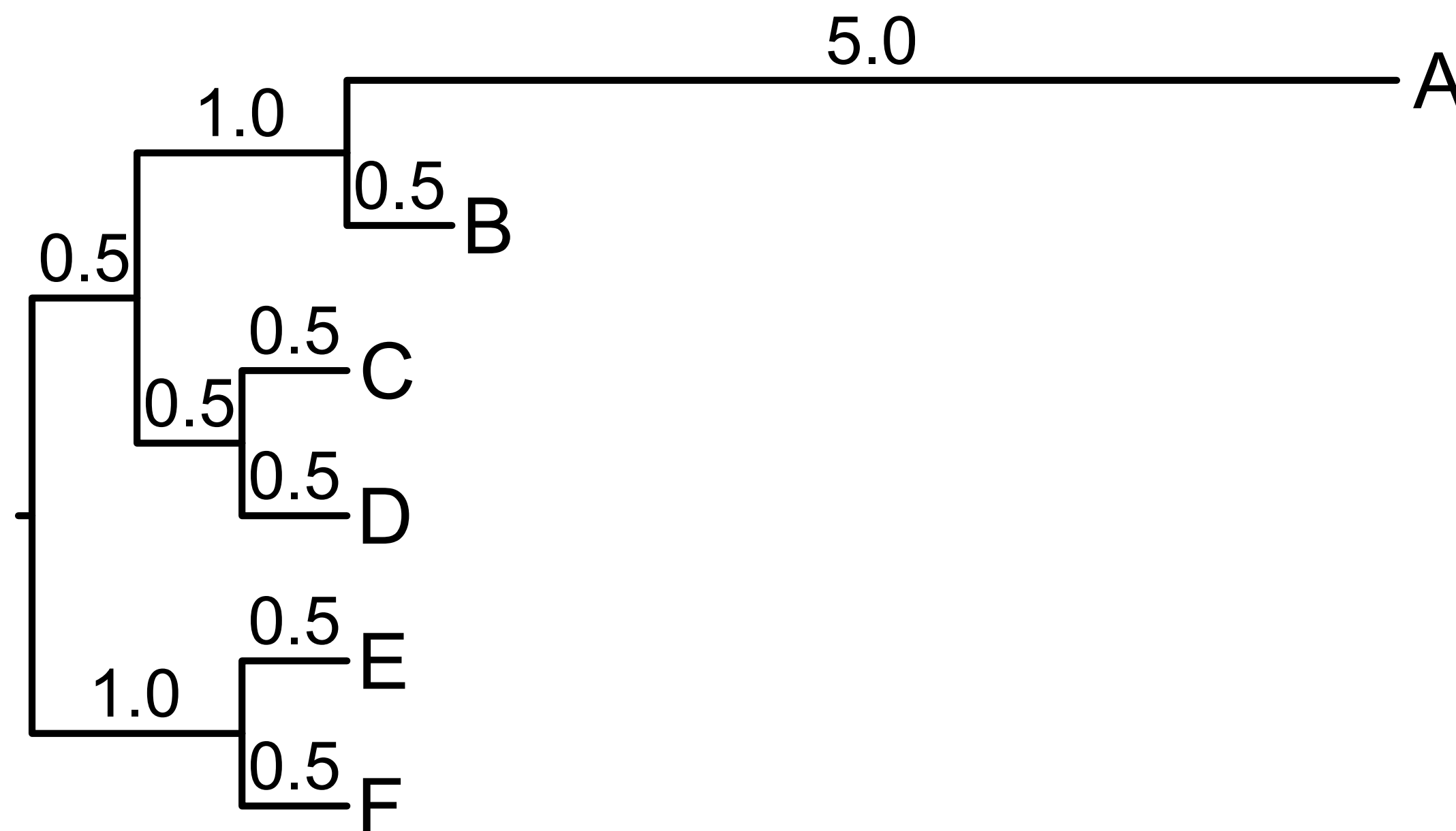
Degree of violation of a molecular clock

Genes with low DVMC may be more useful for divergence time analysis

Phylogenetic signal across taxa

1. Long branch score
2. RCVT

Long branch score



Long branch scores

A: 73.17

B: -14.63

C: -19.51

D: -19.51

E: -9.76

F: -9.76

Relative composition variability per taxon

$$RCVT_j = \sum_{i=1}^c \frac{c_{ij} - \bar{c}_i}{s \times n}$$

- $RCVT_j$: relative composition variability of j th taxon
- c : the number of different character states per sequence type in an alignment
- c_{ij} : number of occurrences of the i th character state for the j th taxon
- \bar{c}_i : the average number of the i th c character state across n taxa
- s : total number of sites
- n : number of taxa

Relative composition variability per taxon

>1
GGGGGCC

>2
ATGCATGC

>3
ATGCATGC

>4
ATGCATGC

>5
GGGGGGGG

Relative composition variability per taxon

>1
GGGGGCCCC

>2
ATGCATGC

Sequences
1 and 5 are
GC rich

>3
ATGCATGC

>4
ATGCATGC

>5
GGGGGGGG

Relative composition variability per taxon

>1
GGGGGCCCC

>2
ATGCATGC

Sequences
1 and 5 are
GC rich

>3
ATGCATGC

→
Calc RCVT


>4
ATGCATGC

>5
GGGGGGGG

Relative composition variability per taxon

	>1 GGGGGCCCC		1 0.12
	>2 ATGCATGC		2 0.09
Sequences 1 and 5 are GC rich	>3 ATGCATGC	→ Calc RCVT	3 0.09
	>4 ATGCATGC		4 0.09
	>5 GGGGGGGGG		5 0.21

Relative composition variability per taxon

Sequences 1 and 5 are GC rich	>1 GGGGGCCCC	 Calc RCVT	1 0.12	Lower values indicate lower biases
	>2 ATGCATGC		2 0.09	
	>3 ATGCATGC		3 0.09	
	>4 ATGCATGC		4 0.09	
	>5 GGGGGGGGG		5 0.21	

Phylogenetic signal across sites

1. Compositional bias
2. Evolutionary rate

Compositional bias per site

>1
GGGGGCC

>2
ATGCATGC

>3
ATGCATGC

>4
ATGCATGC

>5
GGGGGGGG

Compositional bias per site

>1 GGGGGCCCC		1	0.2	0.6547	0.6547
>2 ATGCATGC		2	0.2	0.6547	0.6547
>3 ATGCATGC		3	0.0	nan	nan
>4 ATGCATGC	→ Calc comp bias per site	4	0.2	0.6547	0.6547
>5 GGGGGGGG		5	0.2	0.6547	0.6547
		6	1.6	0.6547	0.4493
		7	1.8	0.6290	0.1797
		8	1.8	0.6290	0.1797

Compositional bias per site

		chi-square	p-val	Multi-test corrected p-val
>1		1	0.2	0.6547
GGGGGCCCC		2	0.2	0.6547
>2		3	0.0	nan
ATGCATGC		4	0.2	0.6547
>3		5	0.2	0.6547
ATGCATGC	→ Calc comp bias per site	6	1.6	0.6547
>4		7	1.8	0.6290
ATGCATGC		8	1.8	0.6290
>5				
GGGGGGGGG				

Compositional bias per site

		chi-square	p-val	Multi-test corrected p-val
>1	GGGGGCCCC	1	0.2	0.6547
>2	ATGCATGC	2	0.2	0.6547
>3	ATGCATGC	3	0.0	nan
>4	ATGCATGC	4	0.2	0.6547
>5	GGGGGGGGG	5	0.2	0.6547
	ATGCATGC	6	1.6	0.6547
	ATGCATGC	7	1.8	0.6290
	GGGGGGGGG	8	1.8	0.6290

Calc comp bias per site

Compositional bias per site

>1
GGGGGCC
>2
ATGCATCC
>3
ATGCATCC
>4
ATGCATCC
>5
GGGGGGCG

→
Calc comp
bias per
site

	chi-square	p-val	Multi-test corrected p-val
1	0.2	0.6547	0.6547
2	0.2	0.6547	0.6547
3	0.0	nan	nan
4	0.2	0.6547	0.6547
5	0.2	0.6547	0.6547
6	1.6	0.6547	0.4493
7	1.8	0.6290	0.1797
8	1.8	0.6290	0.1797

Evolutionary rate per site

>1
GGGGGCC

>2
ATGCATGC

>3
ATGCATGC

>4
ATGCATGC

>5
GGGGGGGG

Evolutionary rate per site

>1 GGGGGCC		1	0.48
>2 ATGCATGC		2	0.48
>3 ATGCATGC		3	0.0
>4 ATGCATGC	→	4	0.48
>5 GGGGGGGG	Calc evo rate per site	5	0.48
		6	0.56
		7	0.32
		8	0.32

Evolutionary rate per site

>1	GGGGGCCCC	1	0.48
>2	ATGCATGC	2	0.48
>3	ATGCATGC	3	0.0
>4	ATGCATGC	4	0.48
>5	GGGGGGGG	5	0.48
		6	0.56
		7	0.32
		8	0.32

Calc evo rate per site

Evolutionary rate per site

>1	GGGGGCC	1	0.48
>2	ATGCATC	2	0.48
>3	ATGCATC	3	0.0
>4	ATGCATC	4	0.48
>5	GGGGGGG	5	0.48
		6	0.56
		7	0.32
		8	0.32

Calc evo rate per site

Evolutionary rate per site

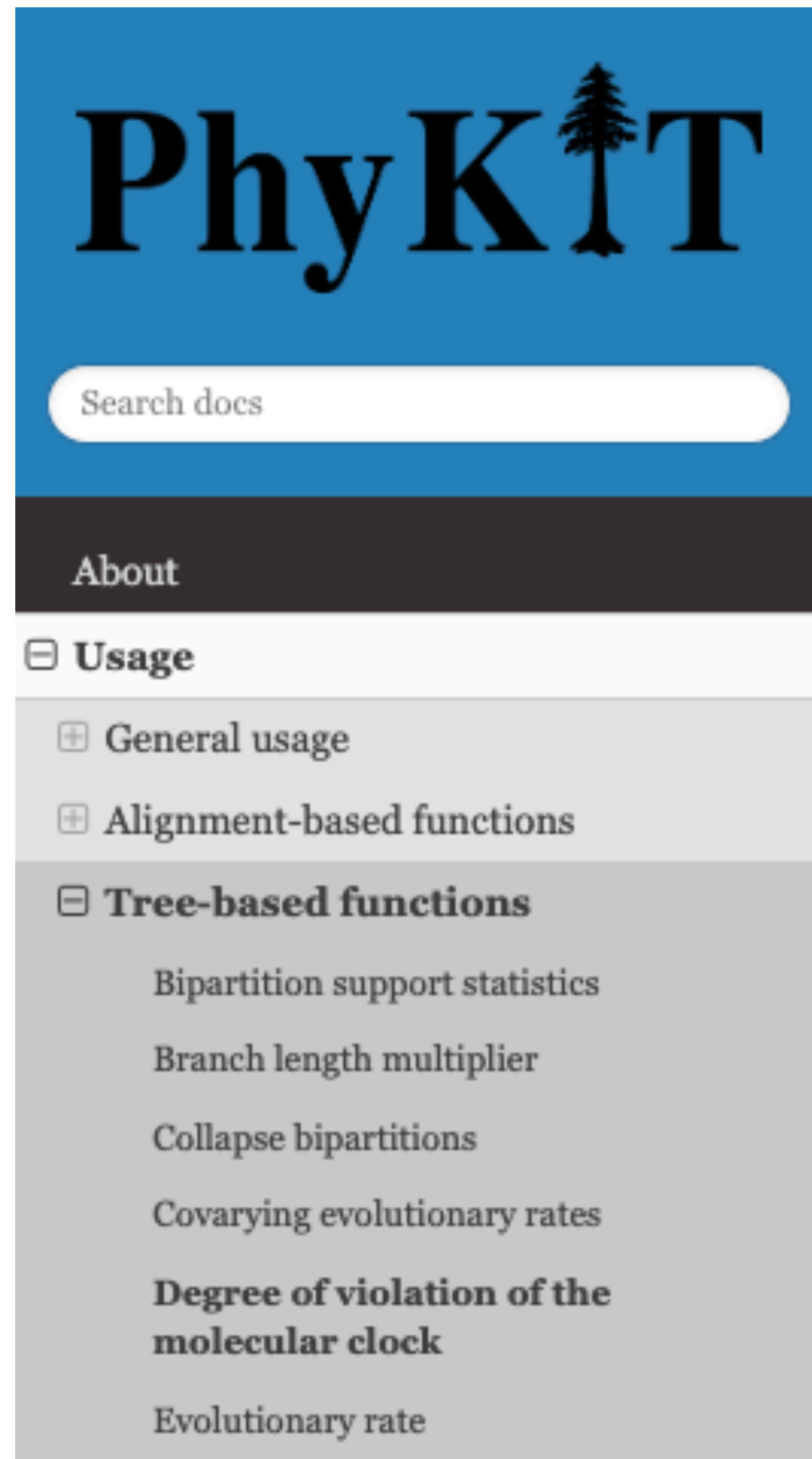
>1	GGGGCC	1	0.48
>2	ATGCATCC	2	0.48
>3	ATGCATCC	3	0.0
>4	ATGCATCC	4	0.48
>5	GGGGGGG	5	0.48
		6	0.56
		7	0.32
		8	0.32

Calc evo rate per site

So many metrics, so many details

1. Alignment length - **higher better**
2. Alignment length with no gaps - **higher better**
3. GC content (for NTs) - **lower better**
4. Pairwise identity - **depends**
5. # of parsimony informative sites - **higher better**
6. # of variable sites - **higher better**
7. Relative composition variability - **lower better**
8. Average bootstrap support value - **higher better**
9. Degree of violation of a molecular clock - **lower better**
10. Evolutionary rate - **depends**
11. Long branch score - **lower better**
12. Treeness - **higher better**
13. Saturation - **lower better**
14. Treeness / RCV - **higher better**
15. RCVT - **lower better**
16. Compositional bias per site - **lower better**
17. Evolutionary rate per site - **depends**

Where known, PhyKIT documentation will say



Degree of violation of the molecular clock

Function names: `degree_of_violation_of_a_molecular_clock`, `dvmc`

Command line interface: `pk_degree_of_violation_of_a_molecular_clock`, `pk_dvmc`

Calculate degree of violation of a molecular clock (or DVMC) in a phylogeny.

Lower DVMC values are thought to be desirable because they are indicative of a lower degree of violation in the molecular clock assumption.

Typically, outgroup taxa are not included in molecular clock analysis. Thus, prior to calculating DVMC from a single gene tree, users may want to prune outgroup taxa from the phylogeny. To prune tips from a phylogeny, see the `prune_tree` function.

Calculate DVMC in a tree following Liu et al., PNAS (2017), doi: 10.1073/pnas.1616744114.

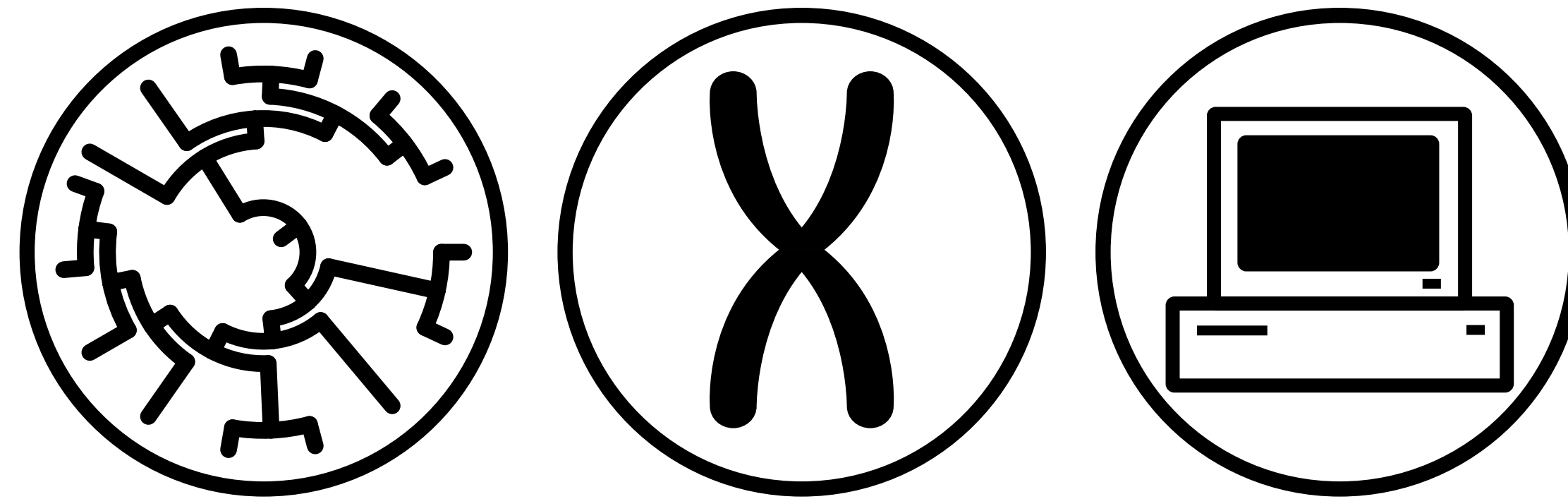
```
phykit degree_of_violation_of_a_molecular_clock <tree>
```

Options:

`<tree>`: input file tree name

<https://jlsteenwyk.com/PhyKIT>

Outline



- Introduction
- Inferring genetic networks from phylogenies
- Phylogenomic subsampling
- **Misc. notes before the tutorial**

Misc. notes on the tutorial

- There are steps in the tutorial for plotting
 - These steps are for the sake of completeness
 - But exporting figures in the container is a little complicated
 - Feel free to skip executing these steps
 - But please read and understand them

Misc. notes on the tutorial

- There are steps in the tutorial for plotting
 - These steps are for the sake of completeness
 - But exporting figures in the container is a little complicated
 - Feel free to skip executing these steps
 - But please read and understand them
- Gemma will have an easier time helping you than me

Misc. notes on the tutorial

- There are steps in the tutorial for plotting
 - These steps are for the sake of completeness
 - But exporting figures in the container is a little complicated
 - Feel free to skip executing these steps
 - But please read and understand them
- Gemma will have an easier time helping you than me
- Curious about career or something not related to the workshop?
 - Feel free to ask!

Thank you for your time and attention!

King Lab

Becca Arruda

Chrisa Staikou

Alain G. De Las Bayonas

Maxwell C. Coyle

Josean Reyes-Rivera

Michael Carver

Stefany Gonzalez



hhmi

Howard Hughes
Medical Institute



Life Sciences
RESEARCH FOUNDATION



King, N



Coyle, M



Buida, J



Li, Y