

Network Evolution

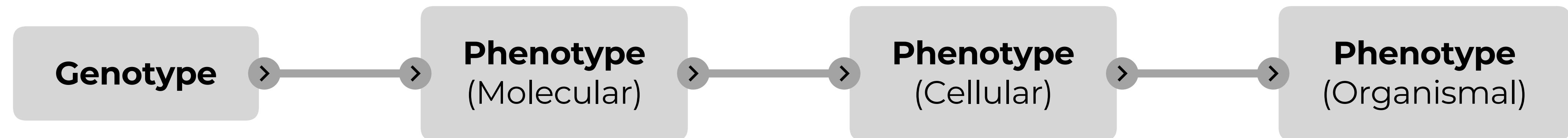
Evolution of Protein Networks and Genomic Conservation

2025/11/21
Chinmay P. Rele

Departmental Seminar Series
Reed Lab — Biological Sciences

Regulation is Complex

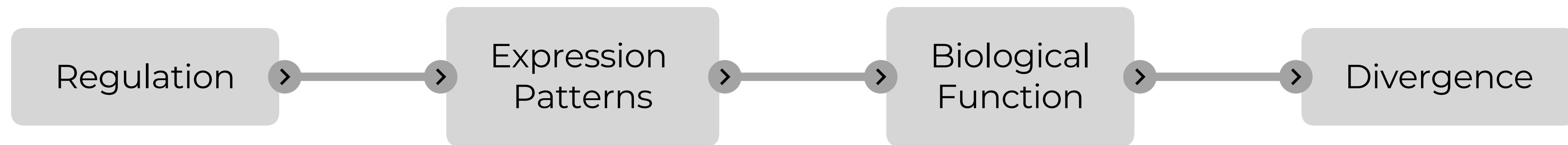
Why



- Genotype can affect phenotype at different scales.
- Evolutionary forces act on these different phenotypes differently.

Regulation is Complex

Why



Network Evolution

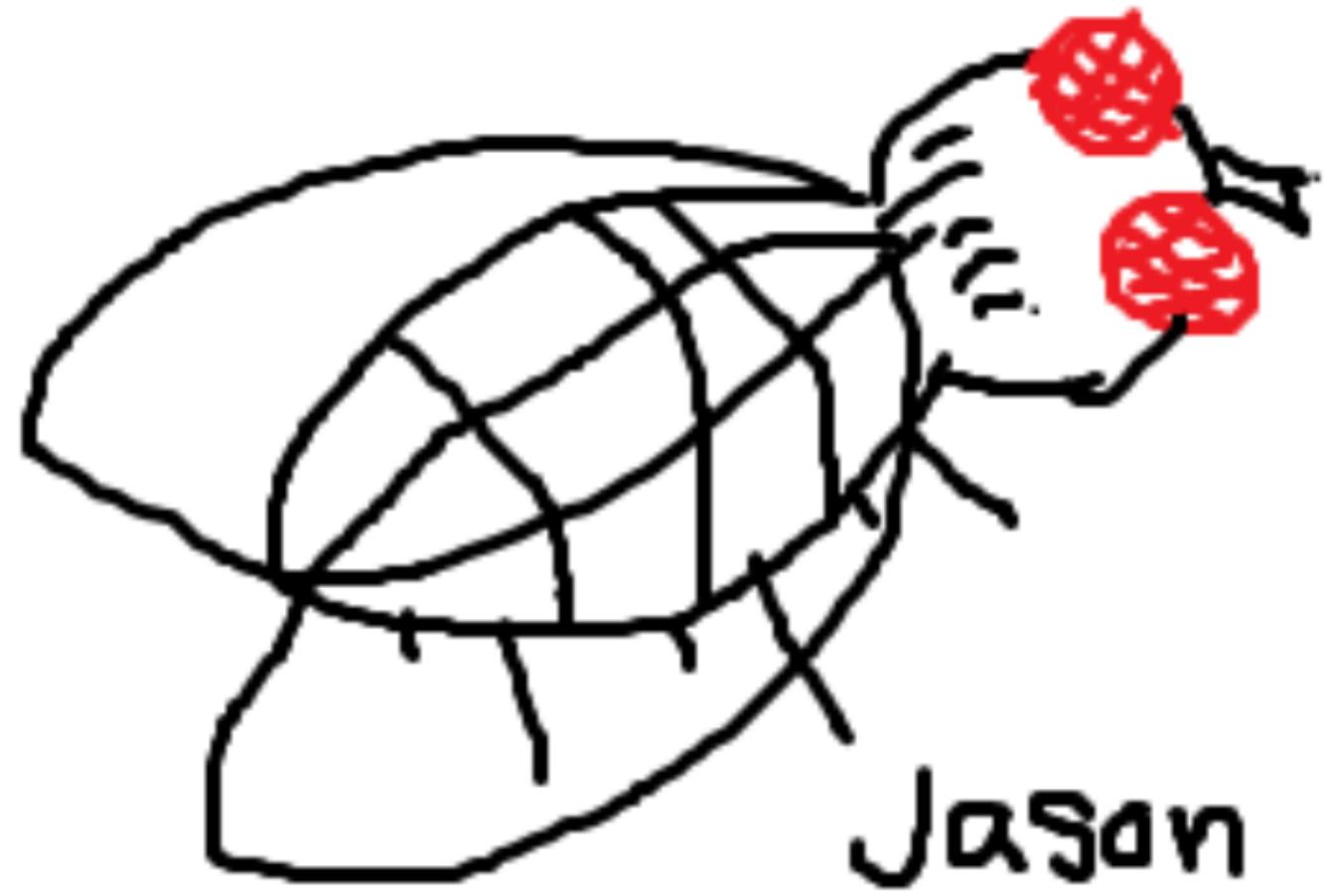
Why

- Predicts proteome rewiring
- Reveals functional adaptation
- Allows for understanding of disease mechanisms
 - ▶ Cancer
 - ▶ Neurogenerative
 - ▶ Infectious (viruses/bacteria)

Network Evolution

What

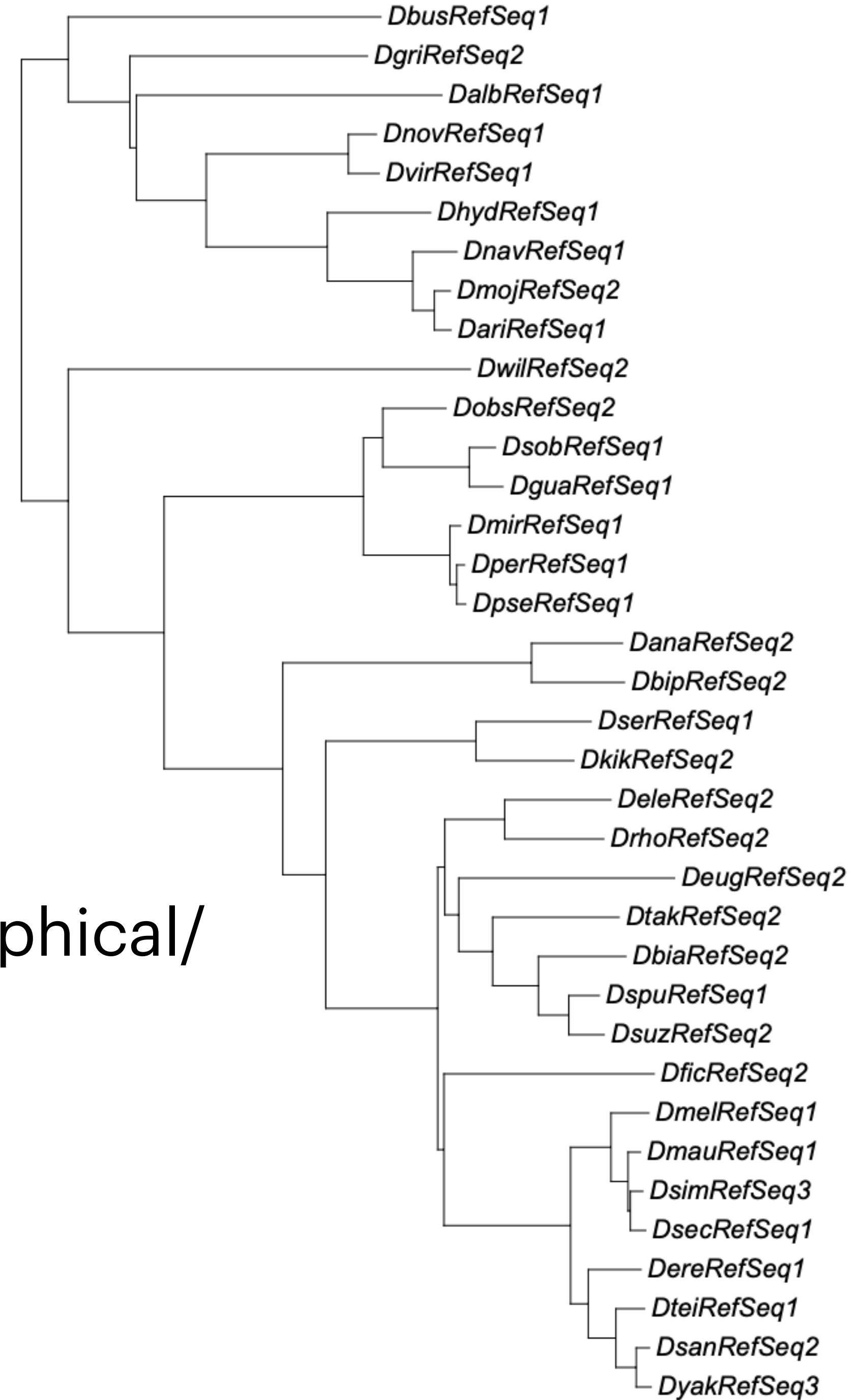
- 1990s → Yeast 2 Hybrid
- 2000/2003 → High throughput screening
- 2010s → Genomics and disease
- NOW → Machine Learning/Evolutionary protein-protein interaction mapping



Drosophila

Why

- Well annotated across clades by NCBI
- Evolutionary relationships known through geographical/morphological/genetic assays
- “Small” genomes easy to compute

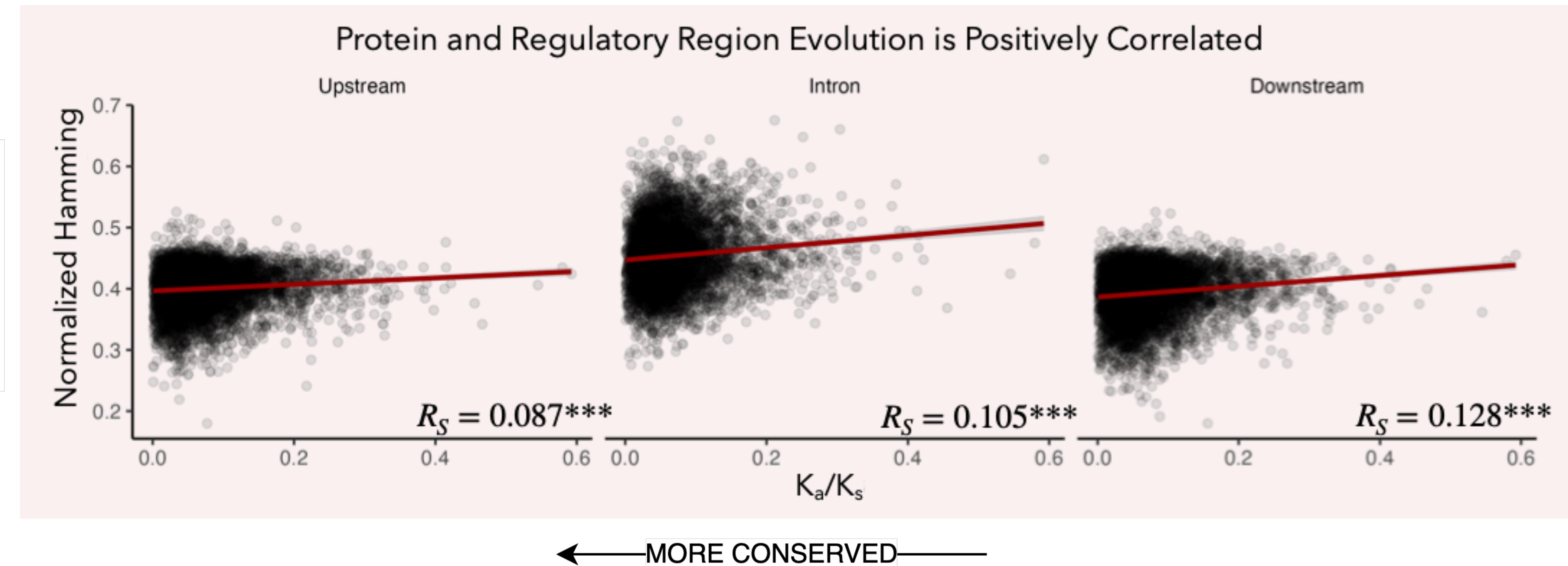


Chapters

1. Evolution of protein coding sequence to regions that regulate their expression
2. Evolution of protein coding sequence and regions that regulate their expression relative to network architecture
3. How the network architecture evolves across the clade

Sequence Conservation is Correlated

Coding Sequence vs. Regulatory Region



Sequence Conservation is Correlated

Measures of Conservation

$$\frac{K_a}{K_s} = \begin{cases} < 1 & \rightarrow \text{purifying selection/constraint} \\ 1 & \rightarrow \text{neutral evolution} \\ > 1 & \rightarrow \text{adaptive evolution} \end{cases}$$

$$K_a = \frac{2}{5} = 0.4$$

$$K_s = \frac{1}{5} = 0.2$$

$$\frac{K_a}{K_s} = \frac{0.4}{0.2} = 2$$

- Rate of synonymous/non-synonymous mutations (accounting for neutral evolution)
- Only for coding sequences

Met	Asp	Thr	Ala	Val
ATG	GAC	ACA	GCG	GTT
ATG	GCC	ACT	TCG	GTT
Met	Ala	Thr	Ser	Val

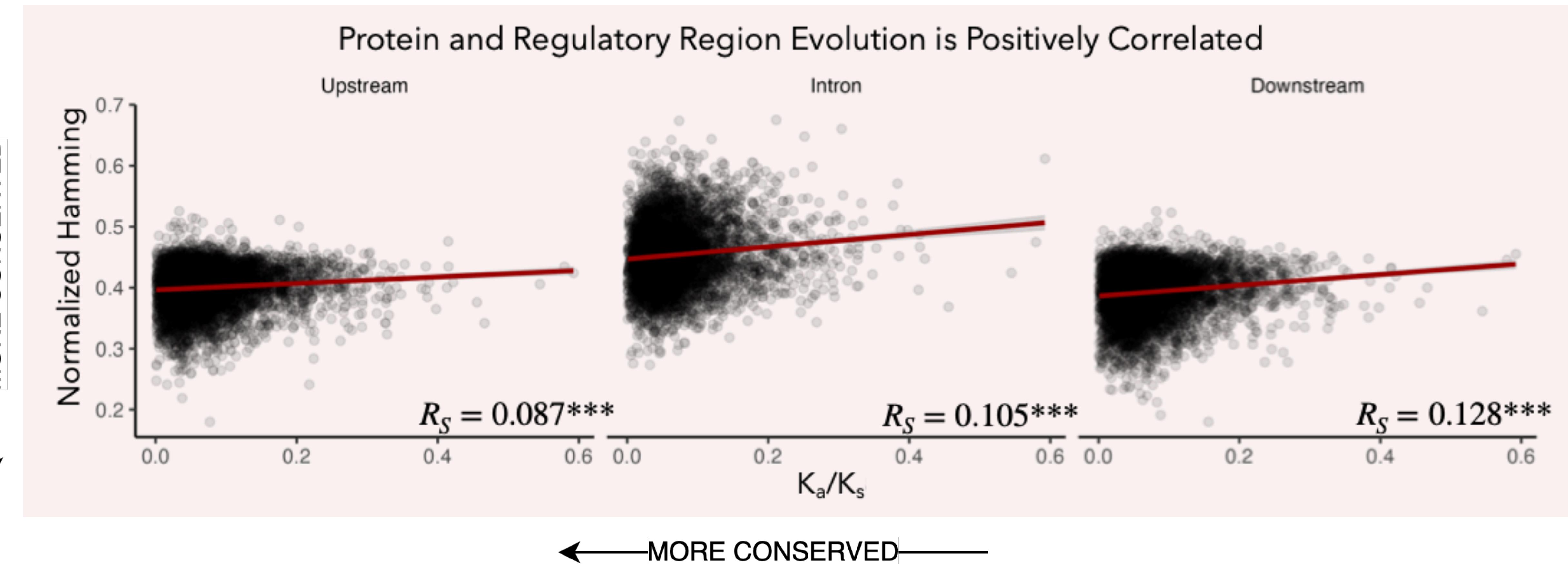
Sequence Conservation is Correlated

Measures of Conservation

Hamming	<ul style="list-style-type: none">"karolin" and "kathrin" is 3."karolin" and "kerstin" is 3."kathrin" and "kerstin" is 4.0000 and 1111 is 4.2173896 and 2233796 is 3.	$Hamming_{normalized} = \frac{Hamming}{Length} = \frac{4}{7}$
Total number of changes in the sequence of identical length	Normalized Hamming takes into account the hamming score relative to the total length of the sequences.	

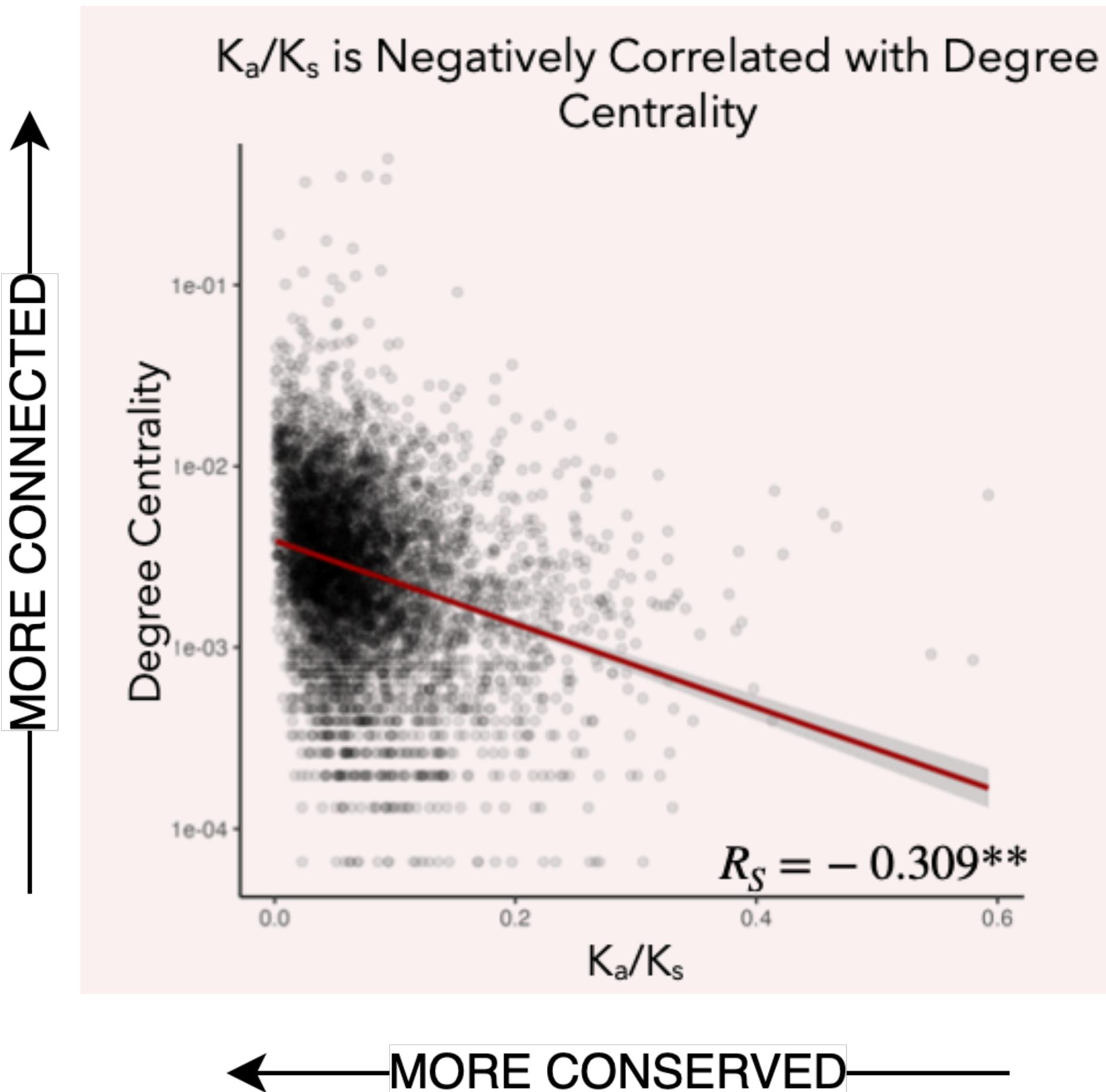
Sequence Conservation is Correlated

Coding Sequence vs. Regulatory Region



Network Position is Correlated with Sequence Divergence

Coding Sequence/Regulatory Region vs. Network

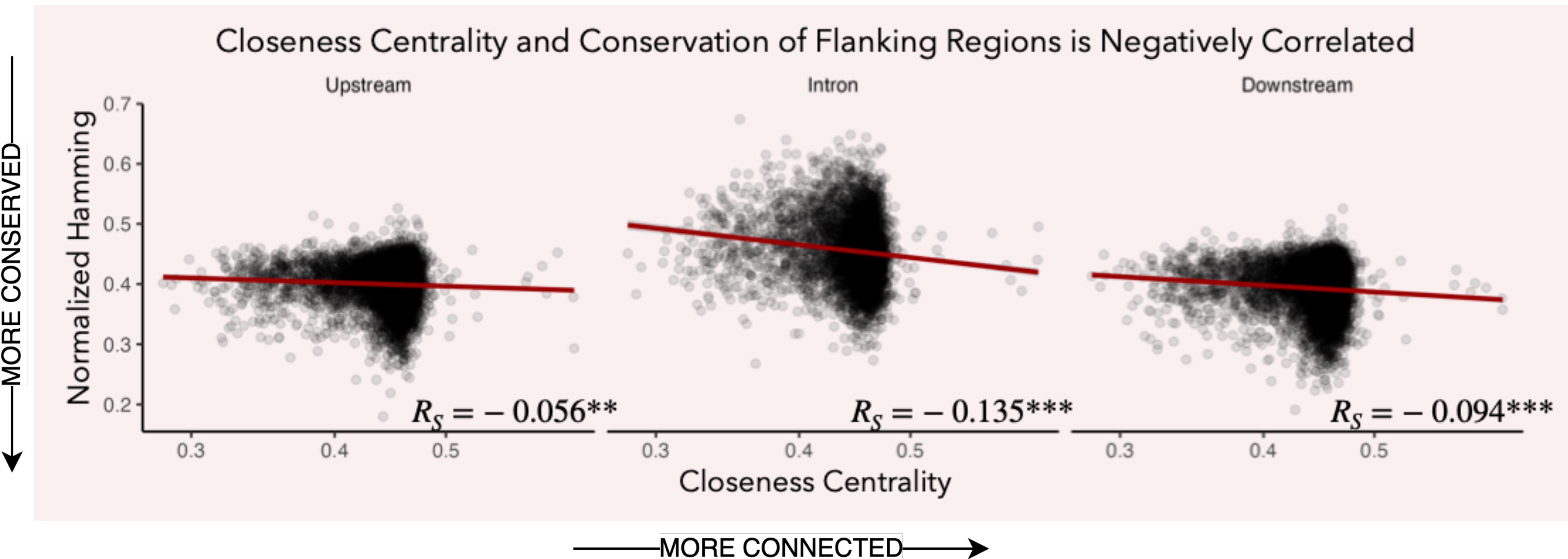


A higher degree centrality for a node indicates that it is more connected to other nodes in the network.

Genes with more conserved coding sequences tend to be connected with more genes.

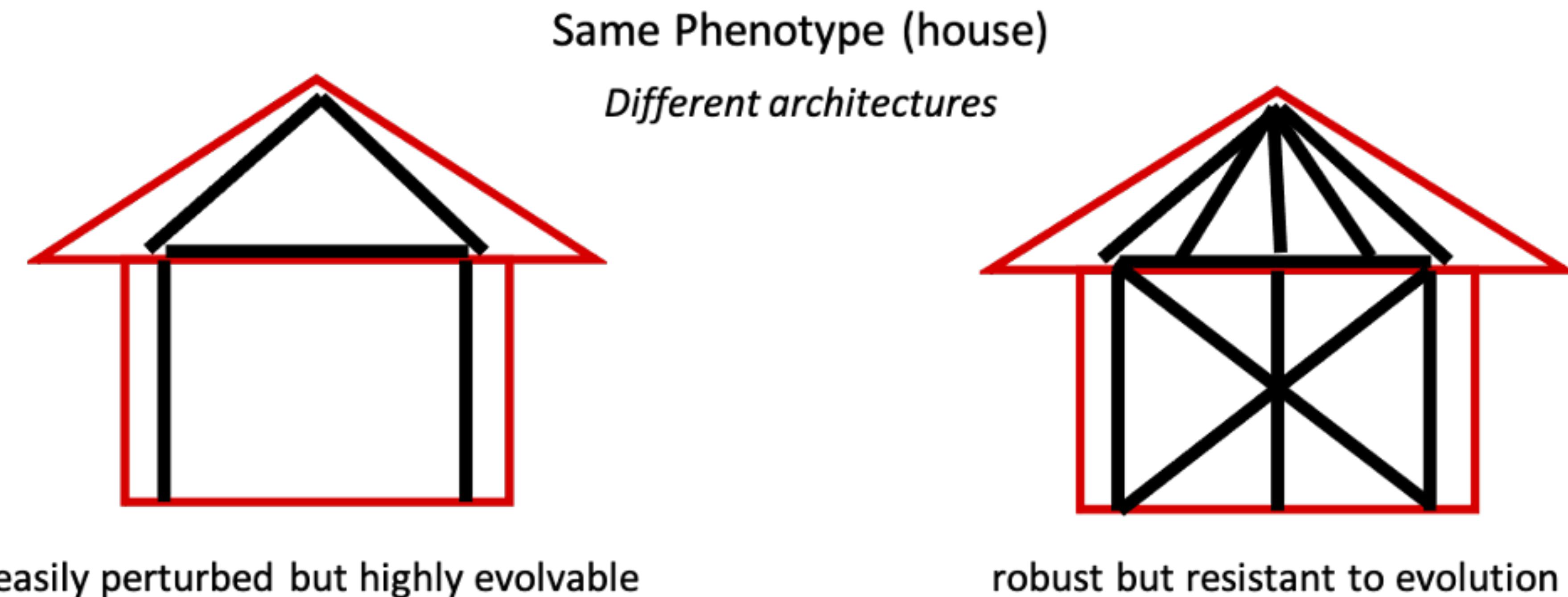
Network Position is Correlated with Sequence Divergence

Coding Sequence/Regulatory Region vs. Network



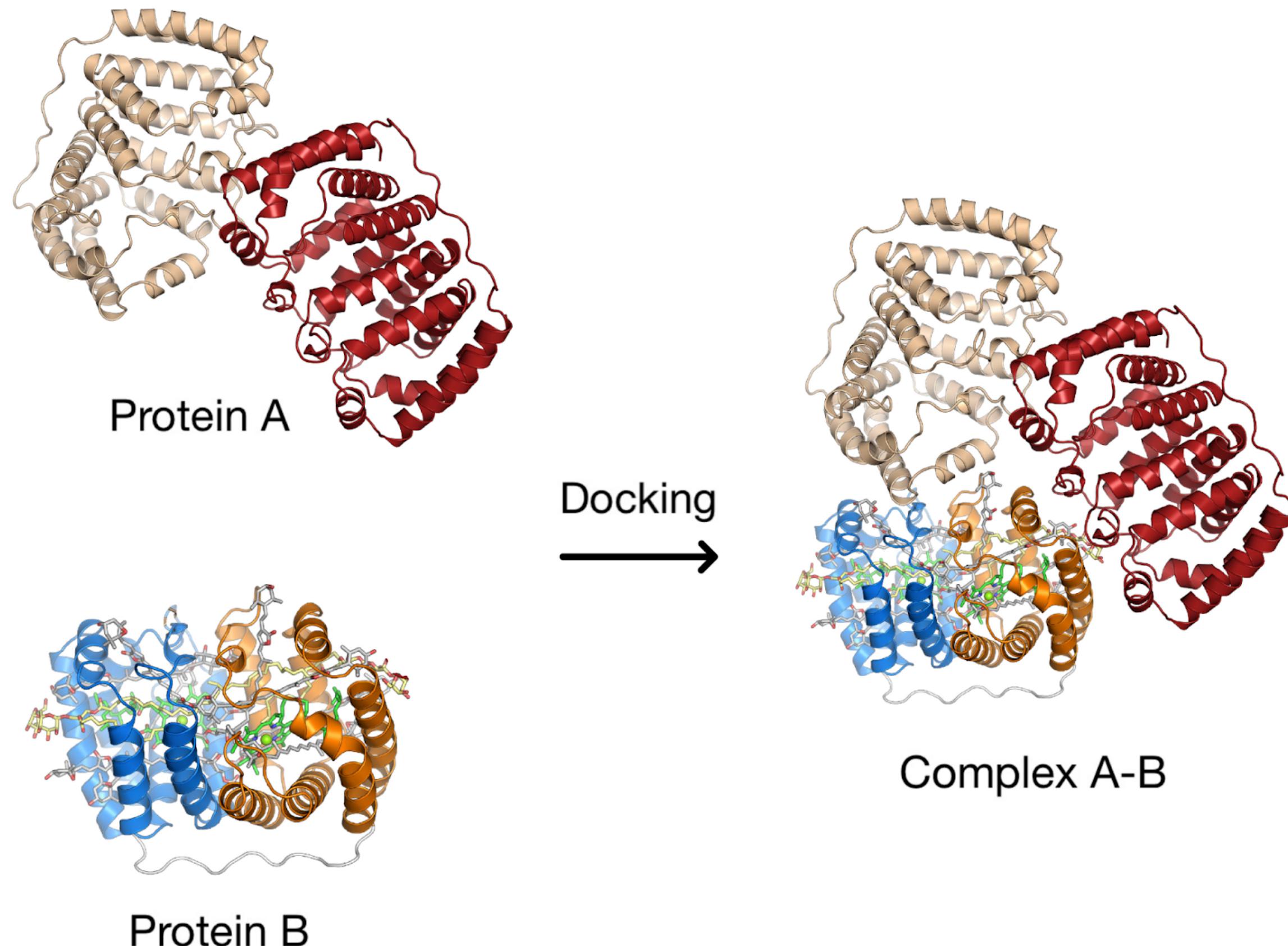
Network Architecture

How it Affects Evolution



Protein-Protein Interactions Inference

Connections



A connection/interaction is inferred if two proteins can interact with each other to form a stable complex.

Protein-Protein Interactions Inference

Interaction Assay Types

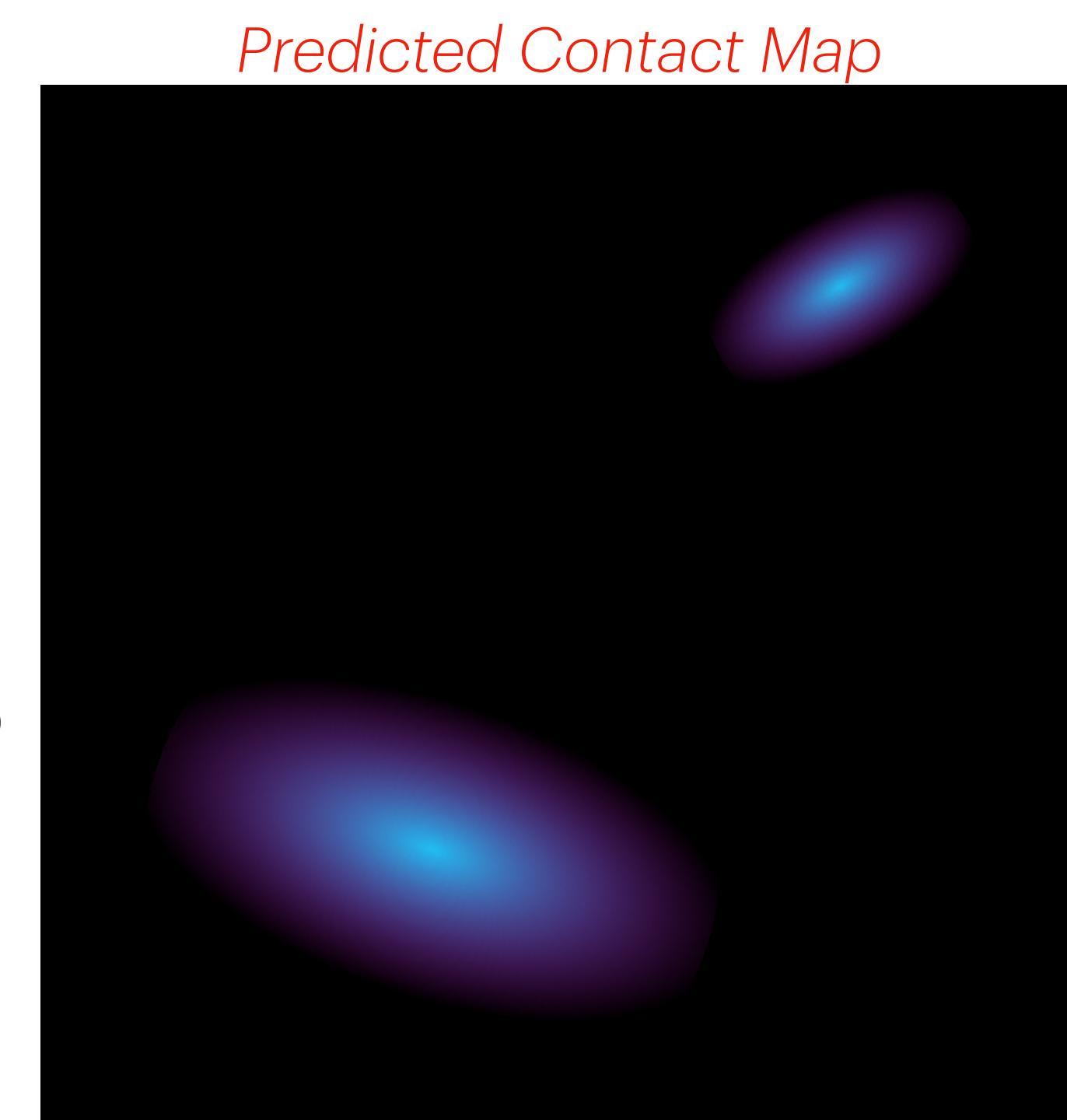
Empirical Studies	<i>in silico</i> Assays
Biologically-derived interactions	Generally sequence-derived interactions
Confirms real interactions	Predicts probabilities of interactions
Protein of interest to see all interacting partners	Can take any proteins and identify probability of interaction
High Confidence Biological	Can be run on multiple genes/species easily

PHILHARMONIC/D-SCRIPT

in silico Network Prediction



Physical Position
along Protein 2



Physical Position
along Protein 1

Interaction
Probability

PHILHARMONIC/D-SCRIPT

in silico Network Prediction

Gene A

GENE A.a

GENE A.b

GENE A.c

Gene B

GENE B.a

GENE B.b

GENE B.c

Gene C

GENE C.a

GENE C.b

PHILHARMONIC/D-SCRIPT

in silico Network Prediction

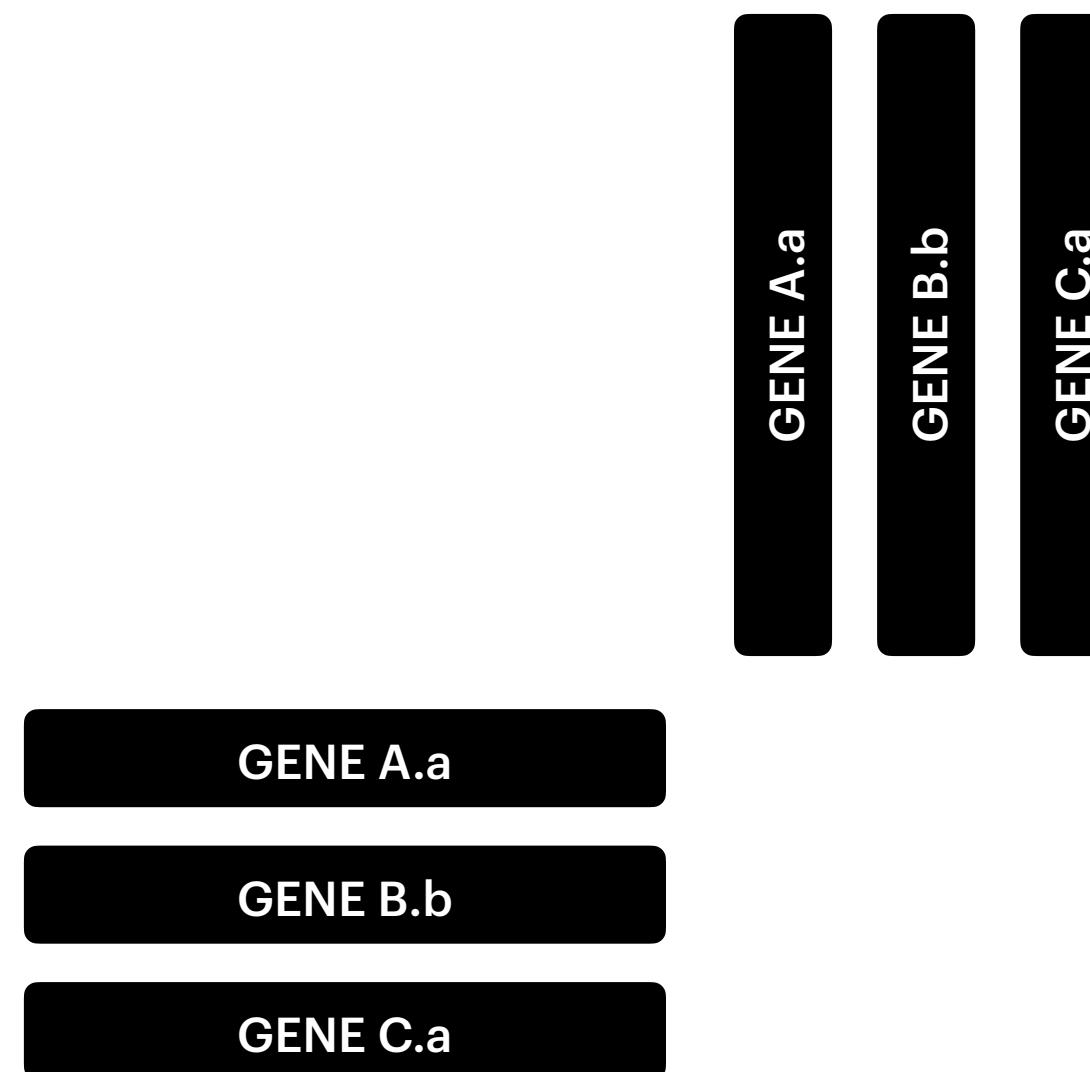
GENE A.a

GENE B.b

GENE C.a

PHILHARMONIC/D-SCRIPT

in silico Network Prediction



PHILHARMONIC/D-SCRIPT

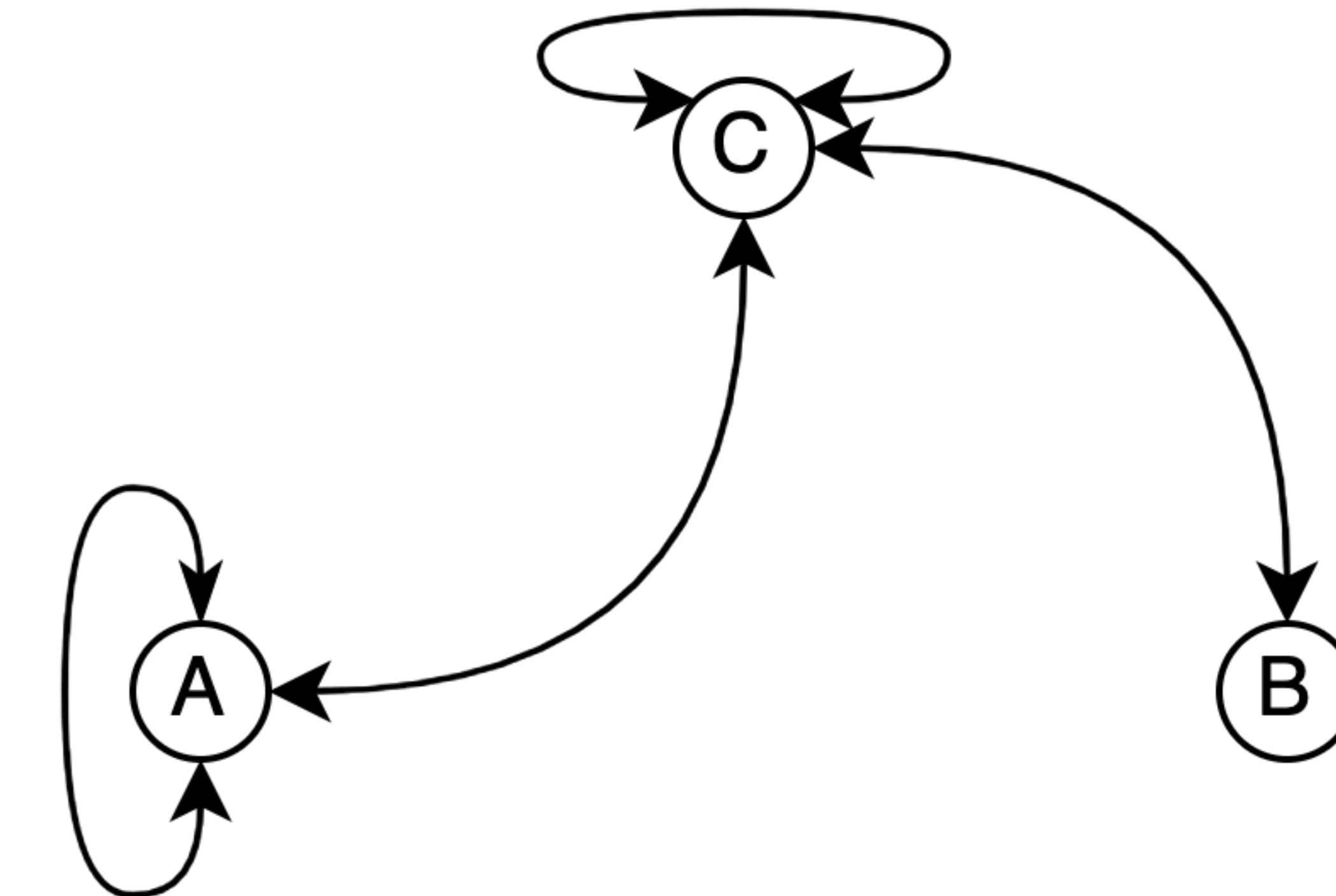
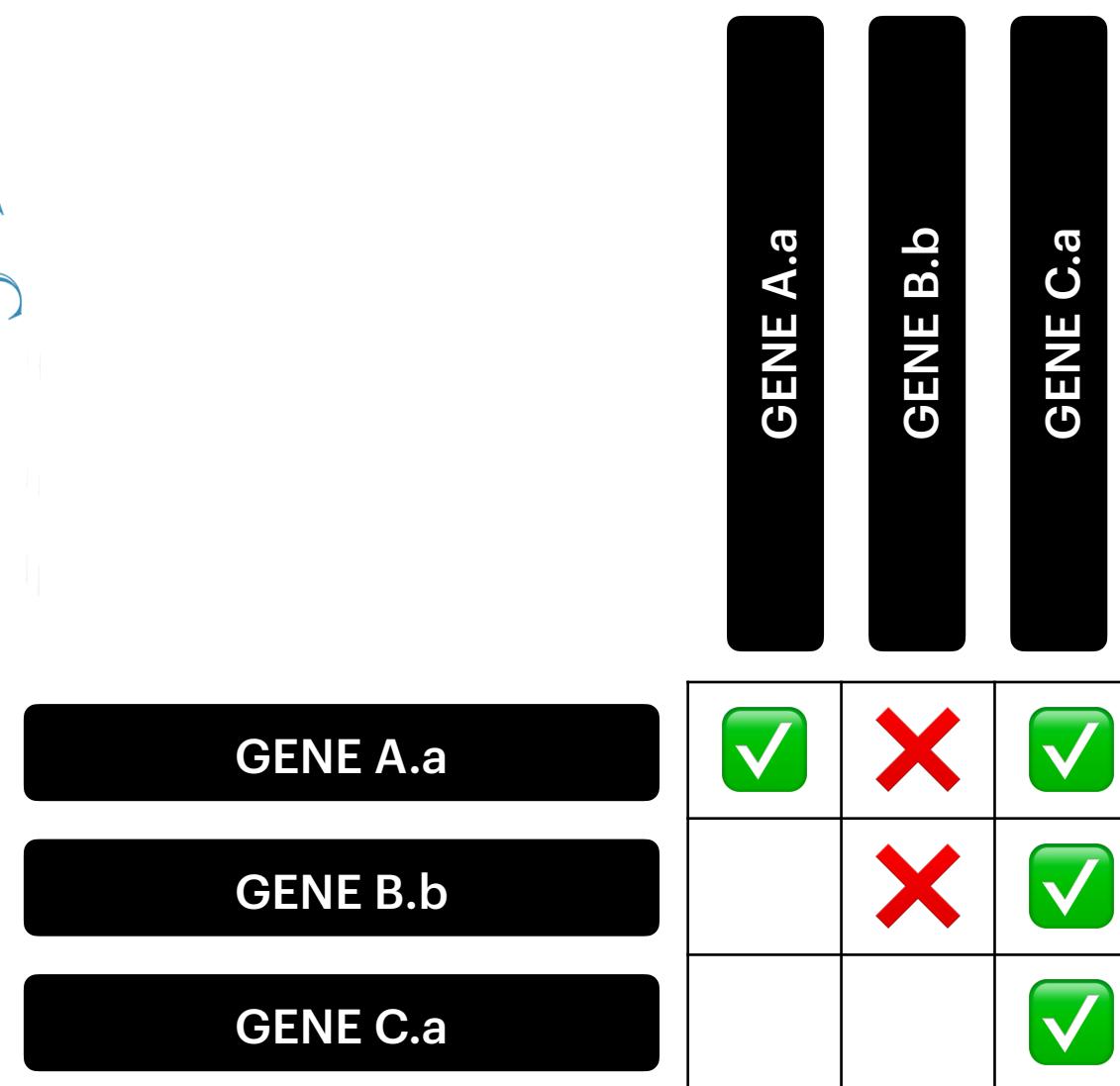
in silico Network Prediction



GENE A.a	GENE B.b	GENE C.a	
GENE A.a	✓	✗	✓
GENE B.b		✗	✓
GENE C.a			✓

PHILHARMONIC/D-SCRIPT

in silico Network Prediction



Network Summaries

Network Architecture

Each species have
similar _____
in their networks as
each other.

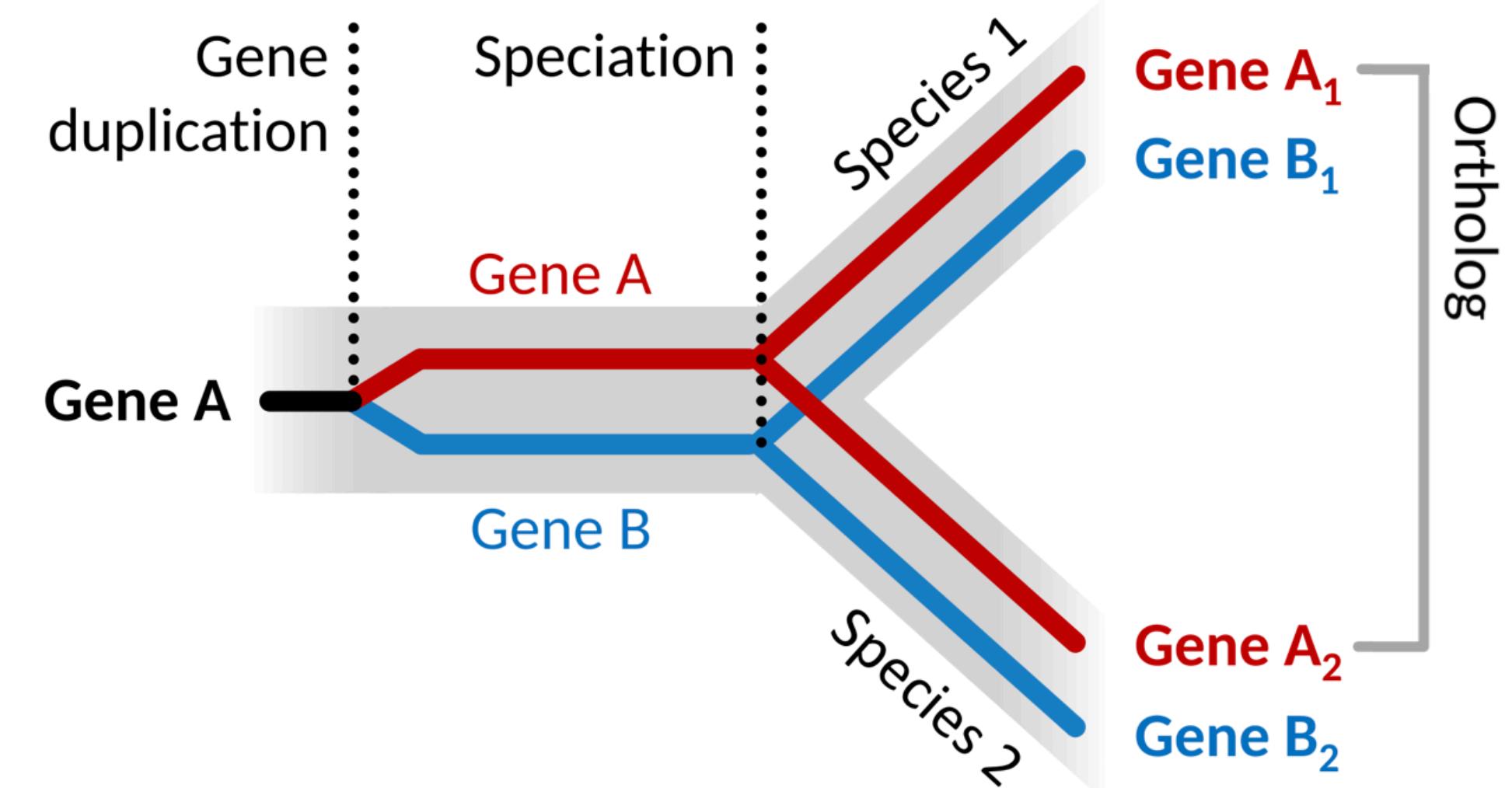
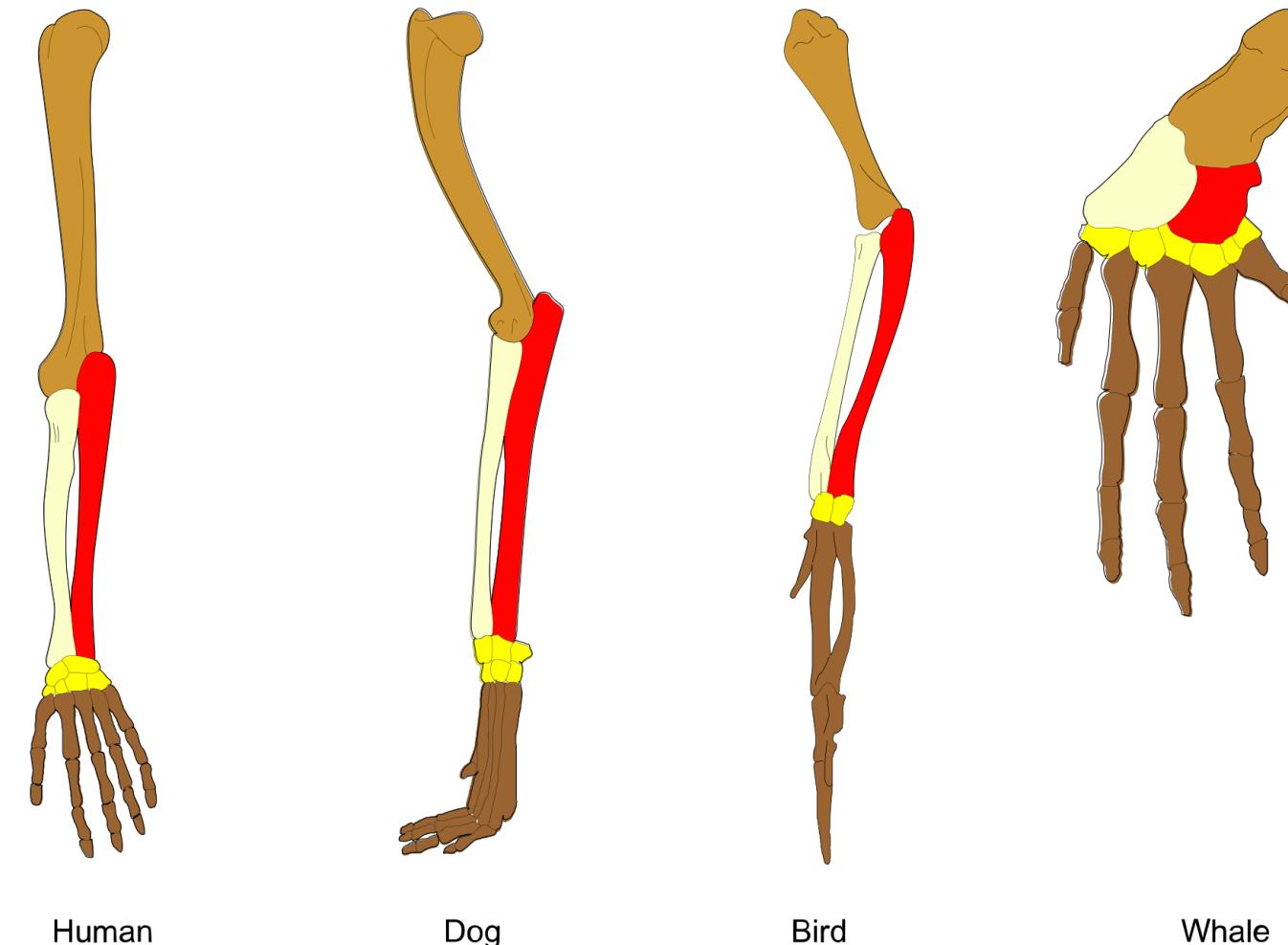
- number of proteins
- number of interactions
- density $\rho = \frac{\text{number of present edges}}{\text{number of possible edges}}$

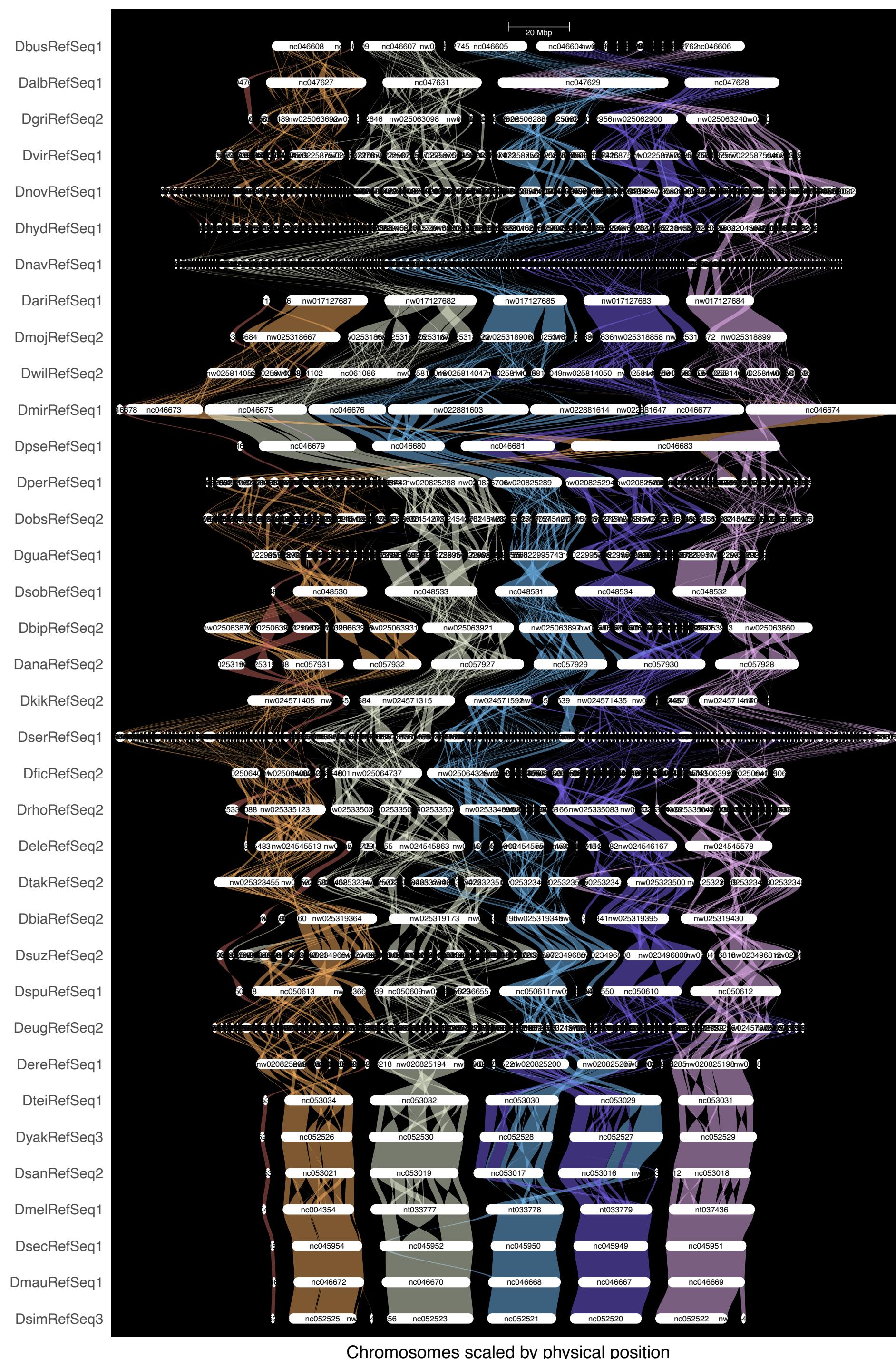
Neighborhood Consistency

Orthogroup Assignment

Orthology is a specialized version of homology.

Any gene that has the same ancestry in 2+ species.





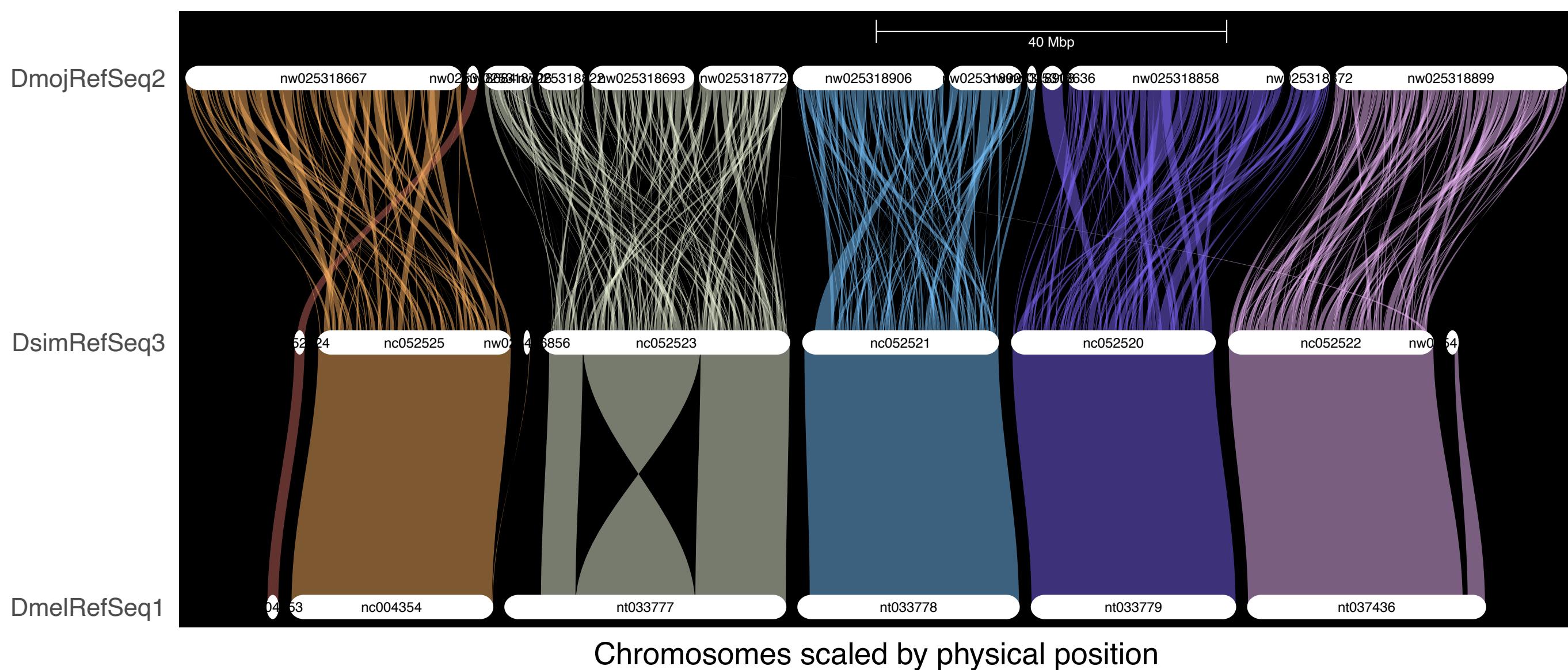
GENESPACE

Neighborhood Consistency

Specialization within Clade

Uses sequence identity to identify orthology.

Refines orthology using position within genome.



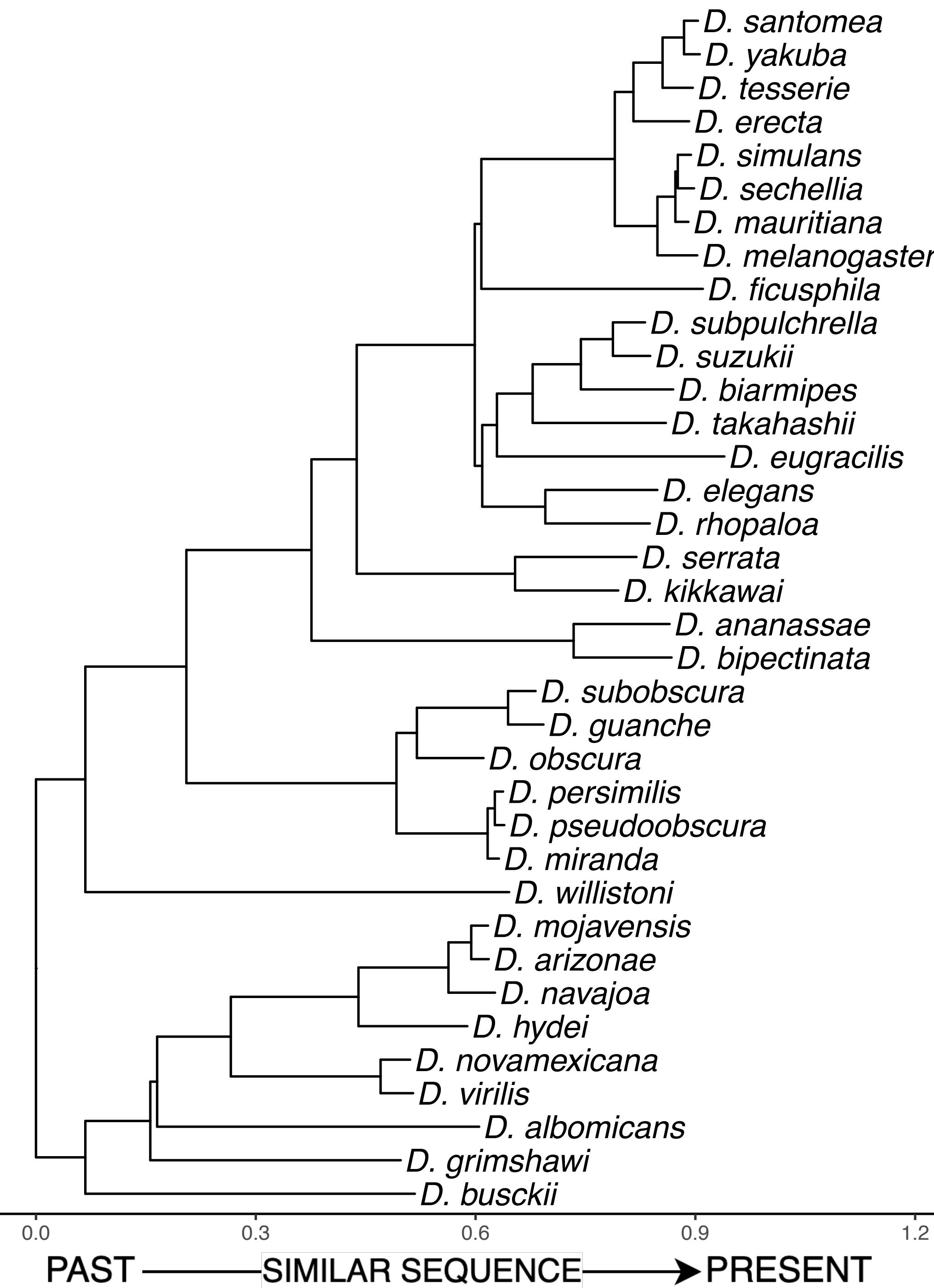
Chromosomes scaled by physical position

Neighborhood Consistency

Specialization within Clade

Closely related species have similar protein sequence.

Species with distant ancestry have more dissimilar sequence.

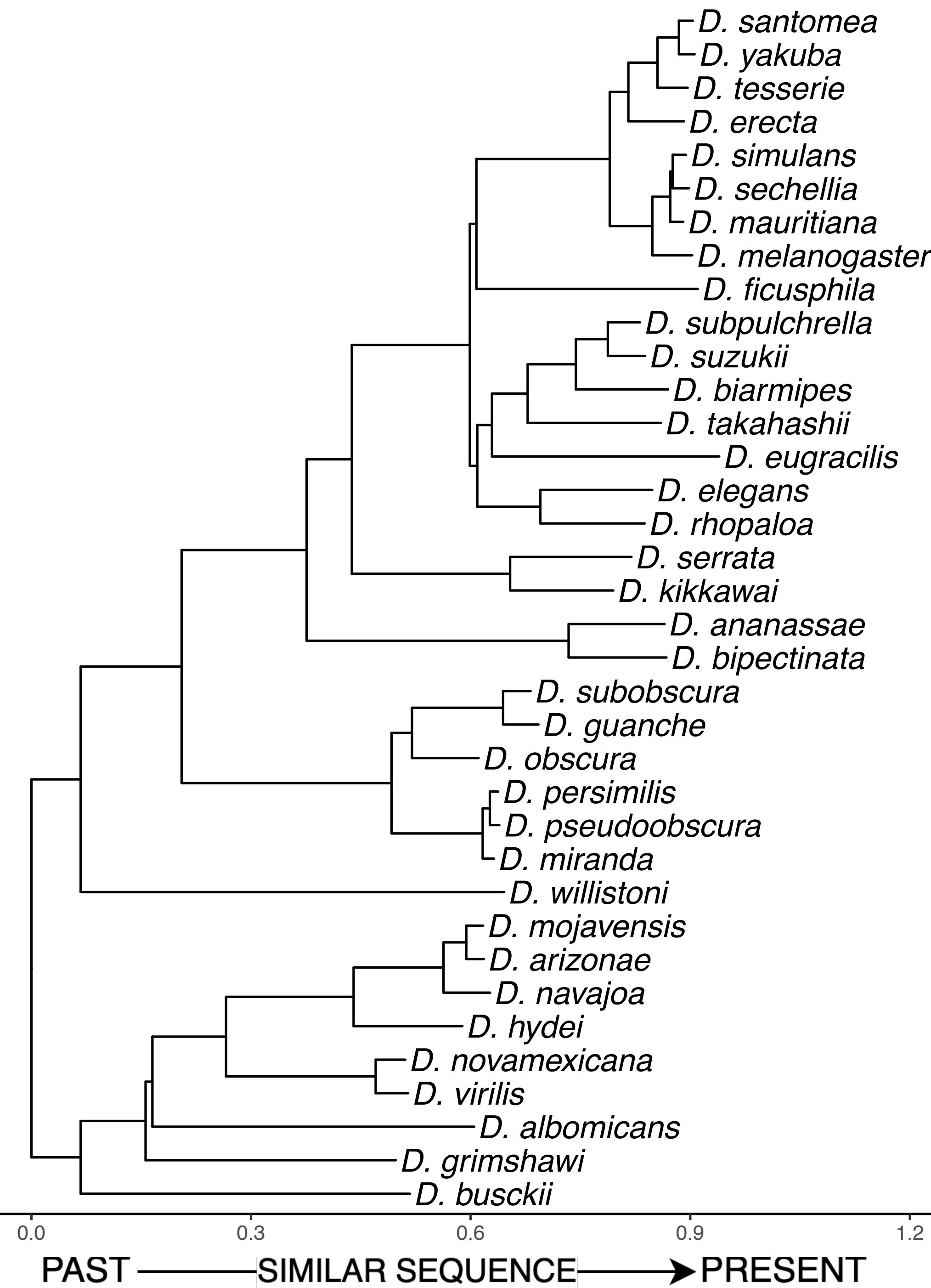


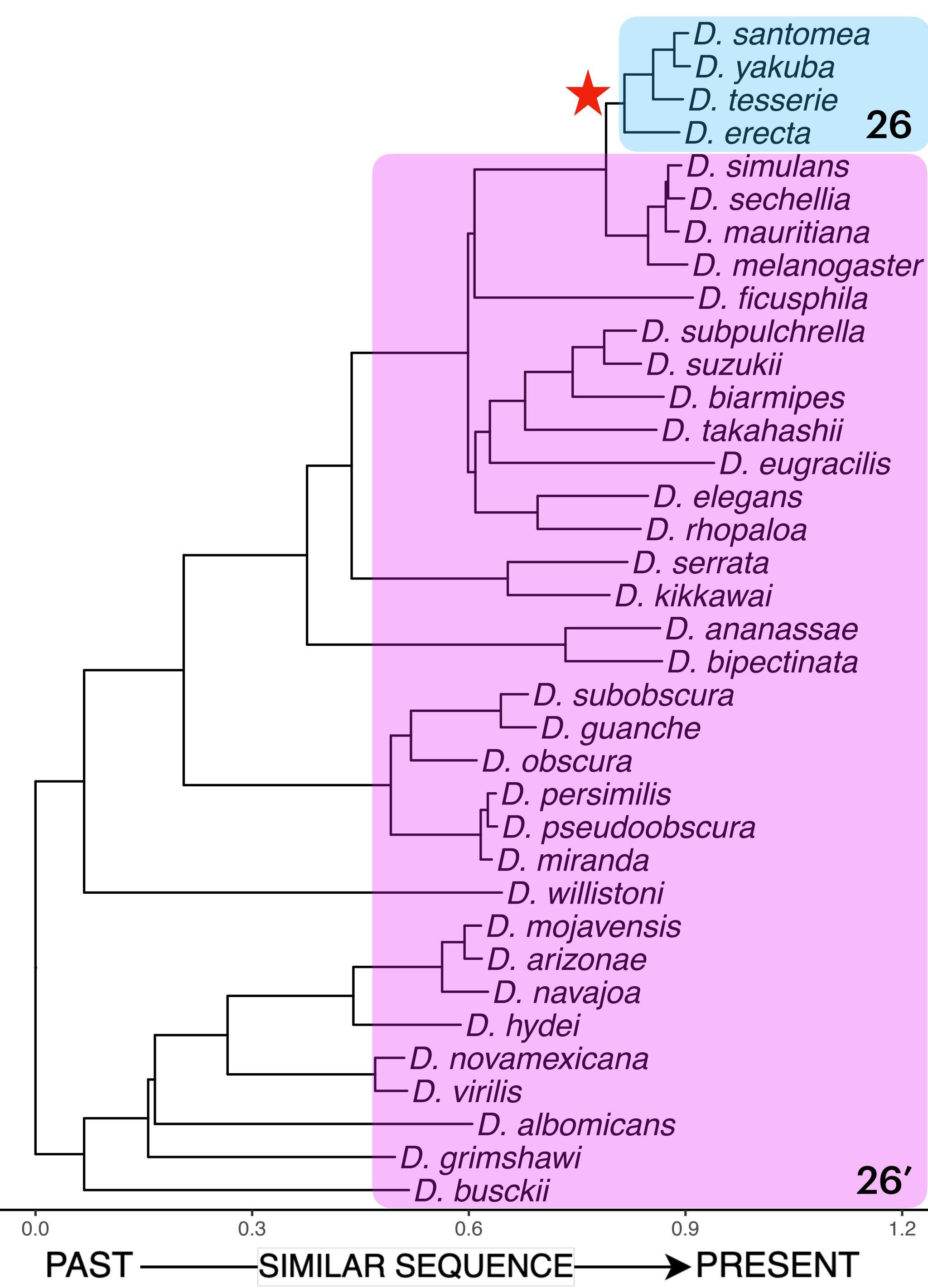
Neighborhood Consistency

Specialization within Clade

Sequence variation has already been studied across clade.

Variation in function is derived through variation in connections to other proteins.



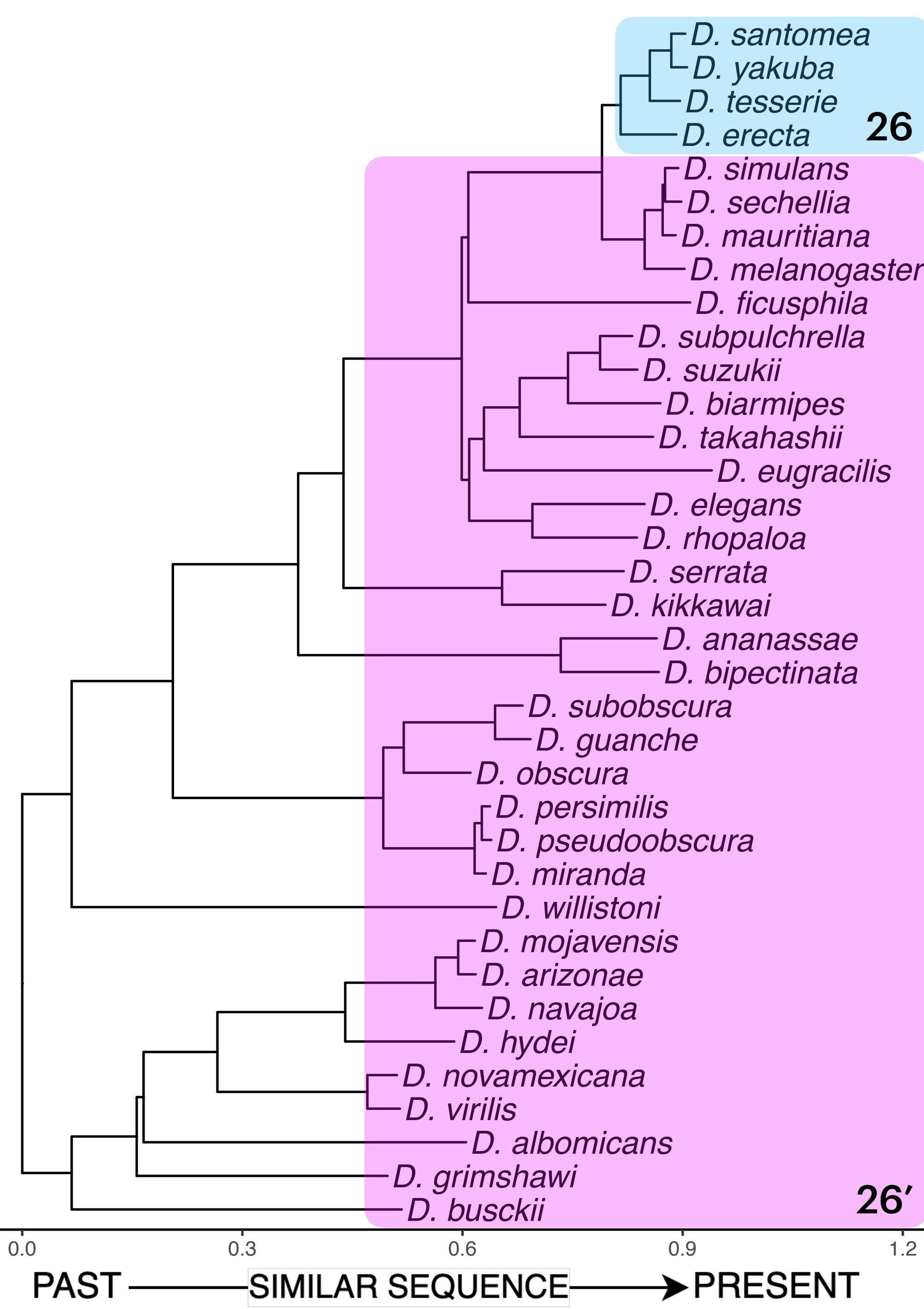


Neighborhood Consistency

Specialization within Clade

How similar the neighborhood of a specific protein within a node.

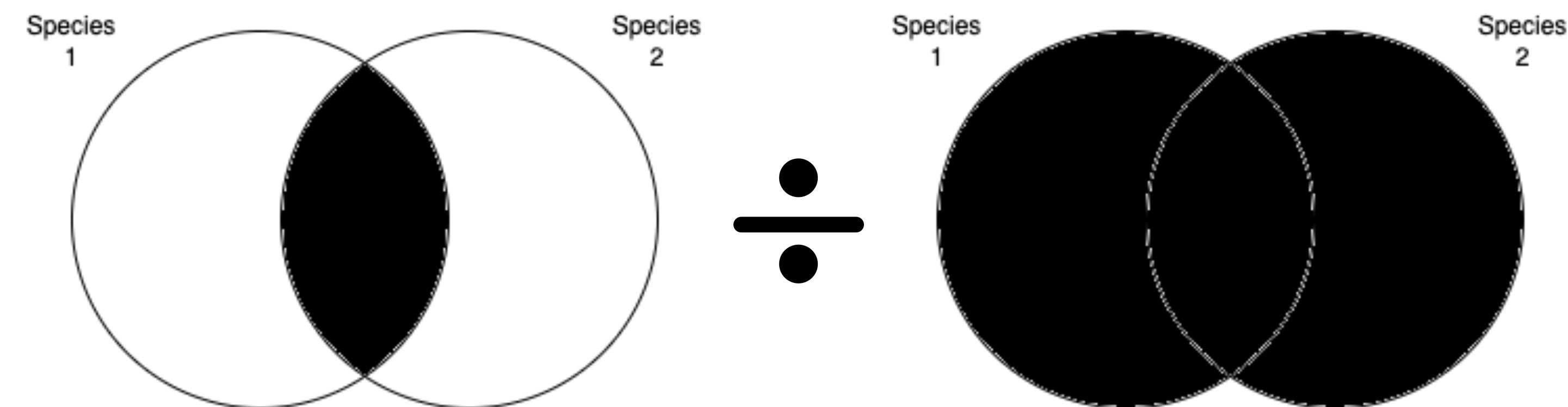
Compare for all species in one node (26) vs all species not in that node (26').



Neighborhood Consistency

Specialization within Clade

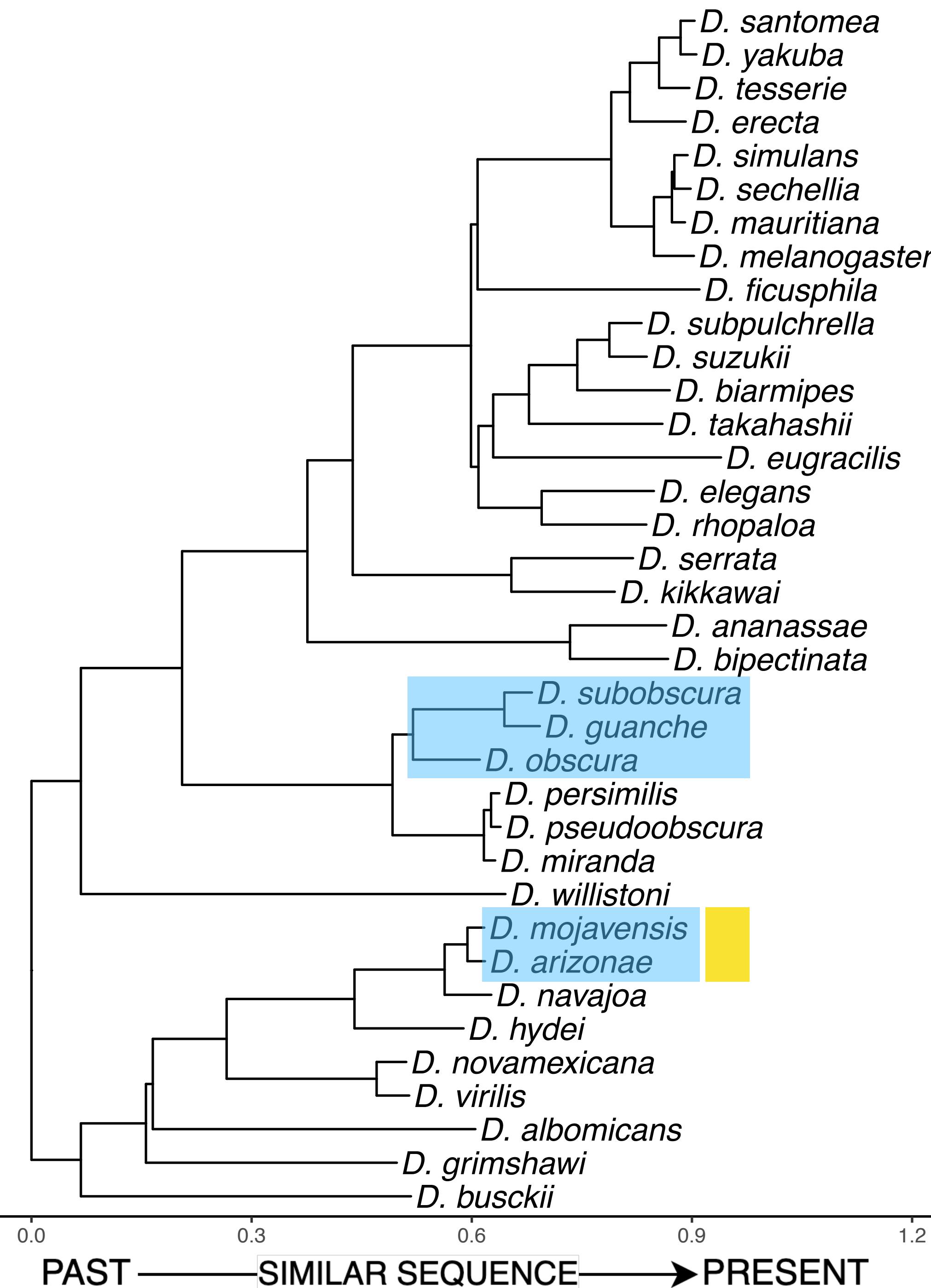
How similar the neighborhood of a specific protein within a node.



Variation of Jaccard Index

Neighborhood Consistency

Specialization within Clade

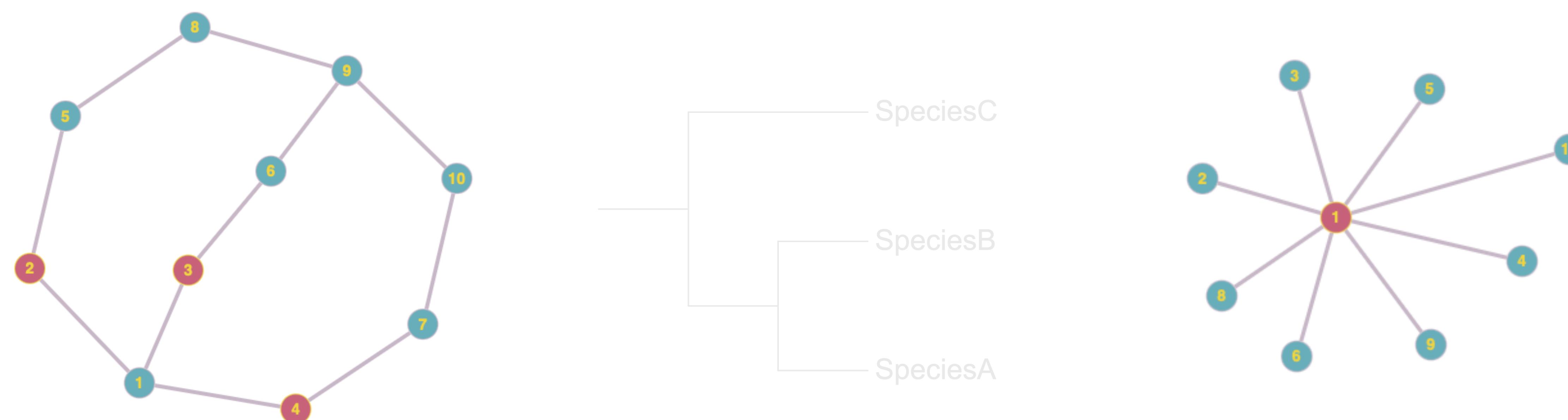


Node	Protein	Δ consistency
24	<i>Tankyrase</i>	-0.607
10		-0.531
24	<i>oskyddad</i>	-0.685

- Some genes are evolving more in certain lineages than others
- Pattern need further individual investigation

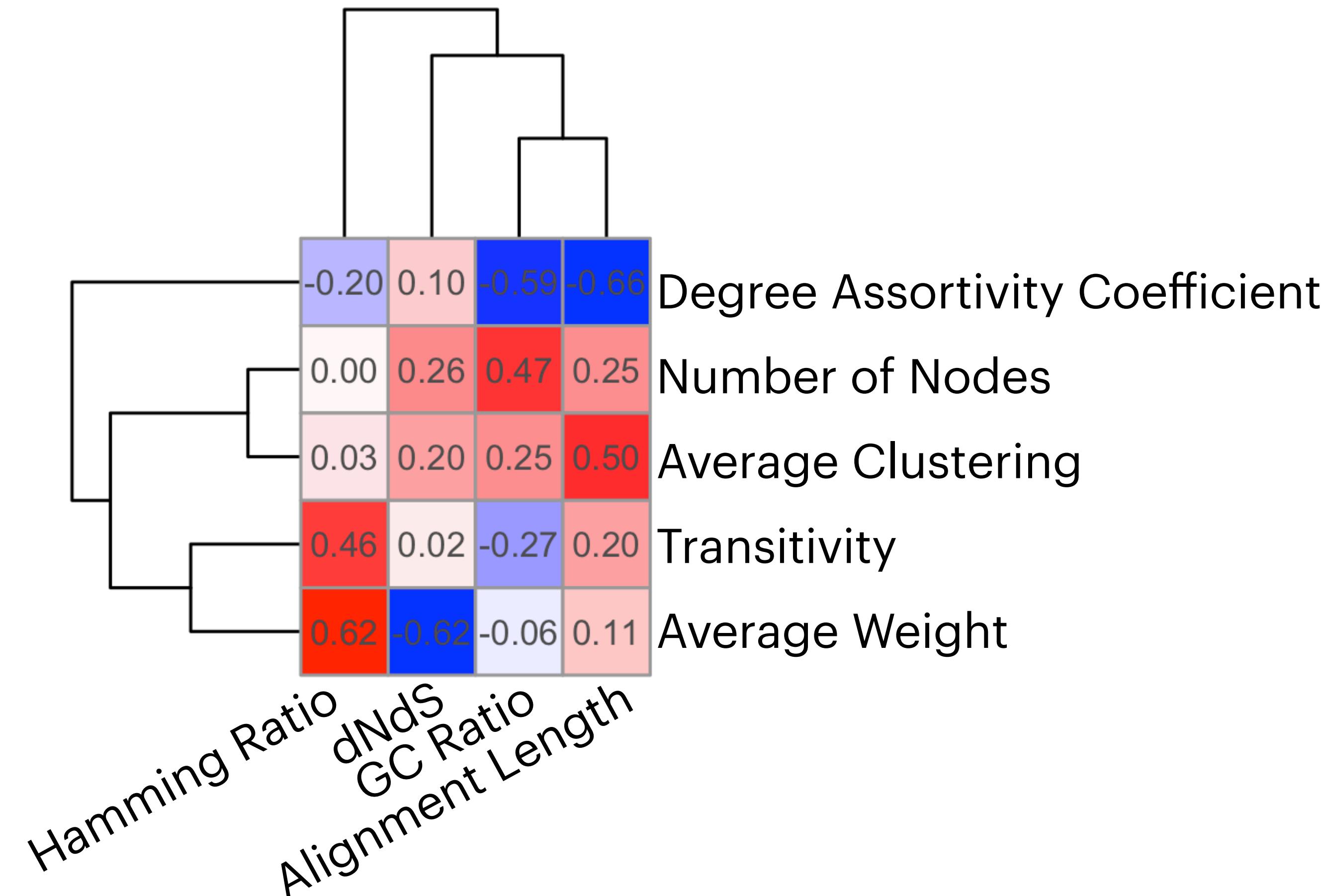
Neighborhood Consistency

Correlation to Genomic Sequence



Neighborhood Consistency

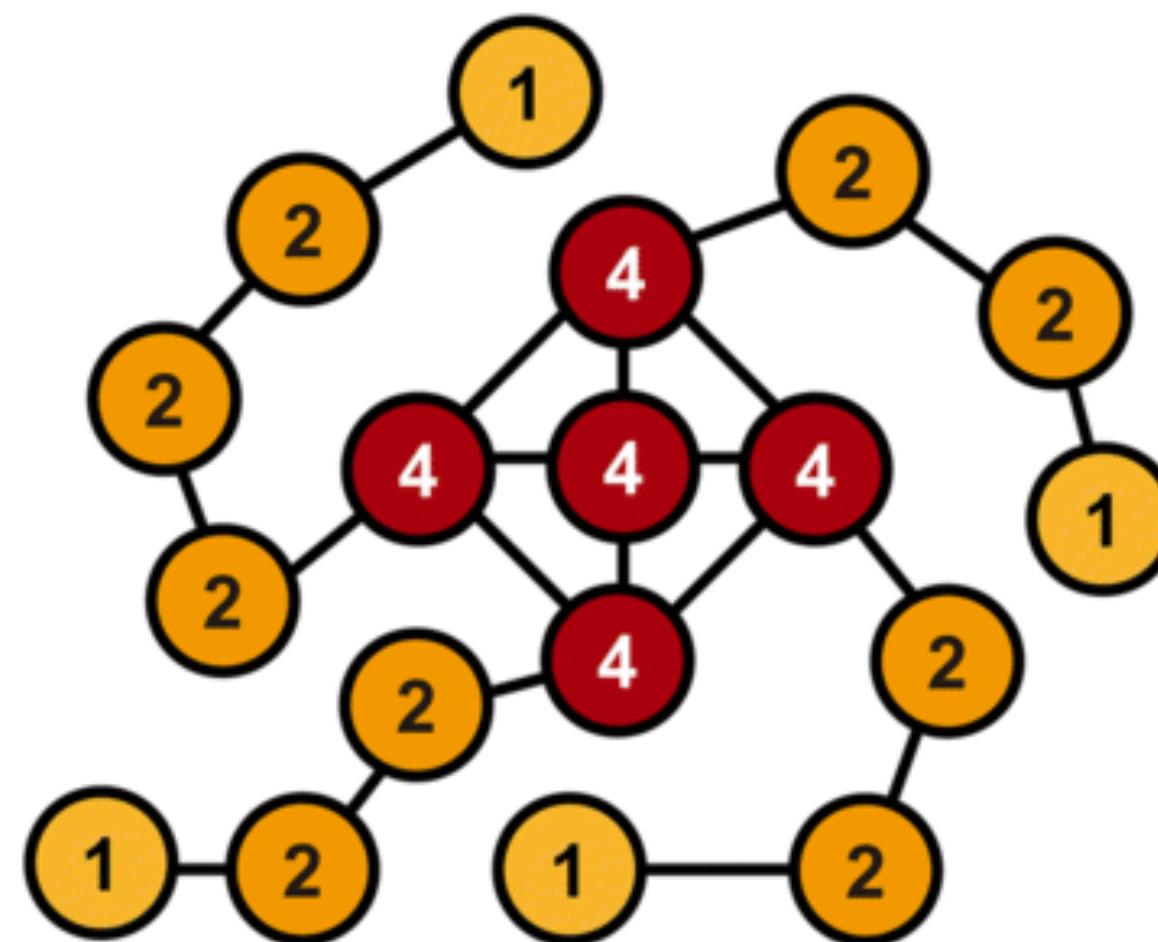
Correlation to Genomic Sequence



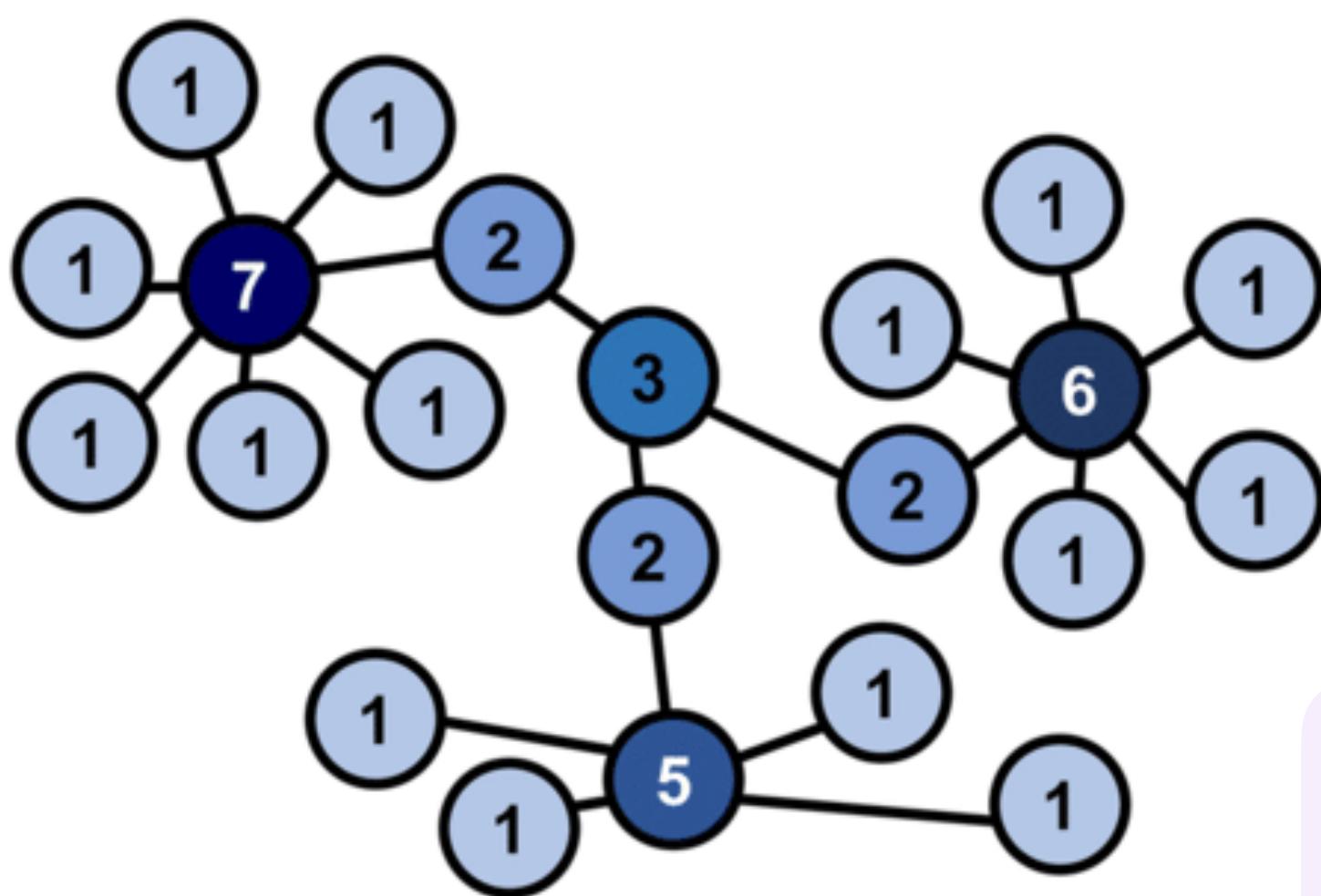
Neighborhood Consistency

Assortativity vs. Regulatory Region GC Ratio

A Assortative network



B Disassortative network



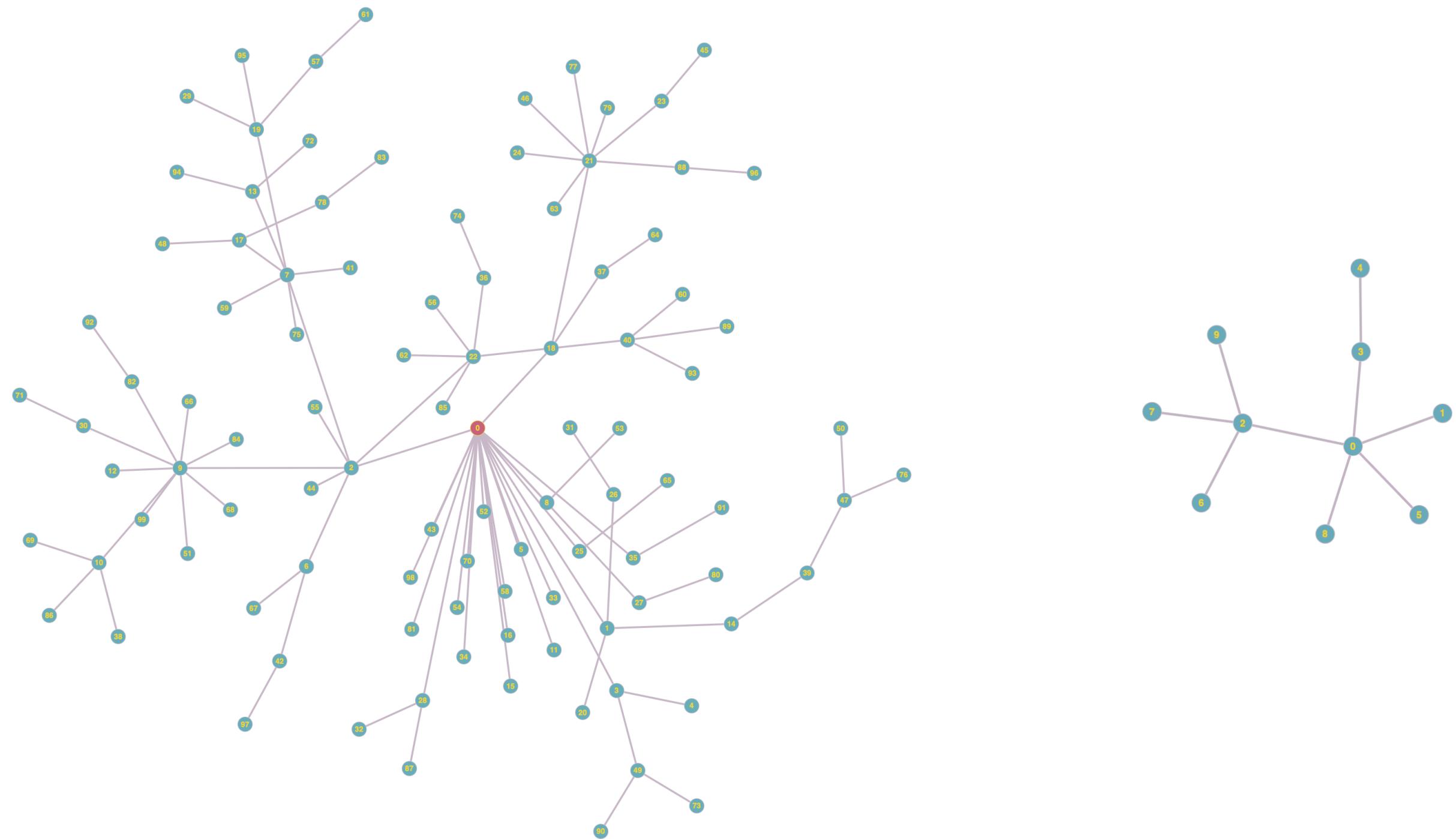
High NEGATIVE correlation with GC ratio.

- Low assortativity means that highly connected nodes are connected to lowly connected nodes.
- More **hub and spoke** architecture

Networks with more hub/spoke architecture had more GC content.

Neighborhood Consistency

Number of Nodes vs. Regulatory Region GC Ratio



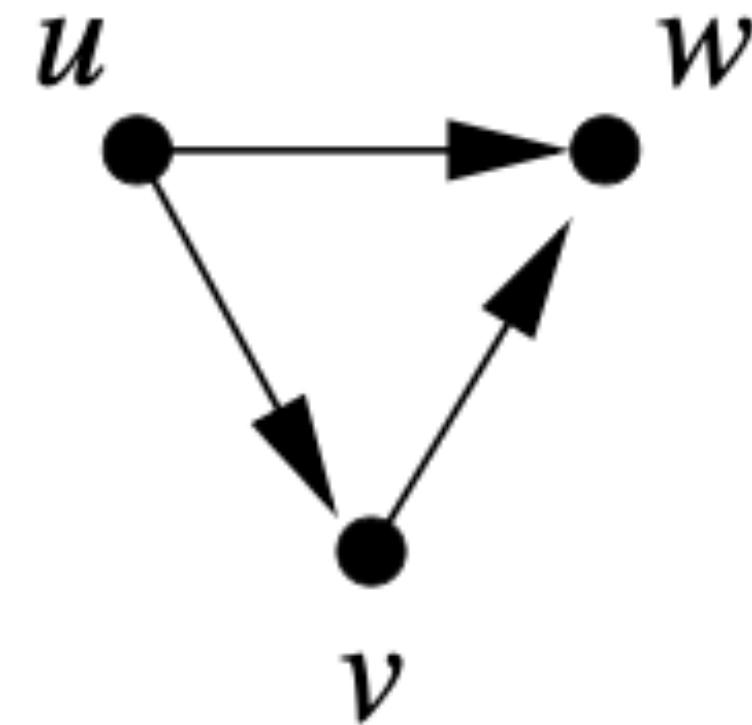
- Number of unique interacting partners across the entire clade

High **POSITIVE** correlation with regulatory region GC ratio.

The more the interacting partners, the higher the GC content.

Neighborhood Consistency

Transitivity vs. Regulatory Region Conservation



A transitive triple of nodes
in a directed network.

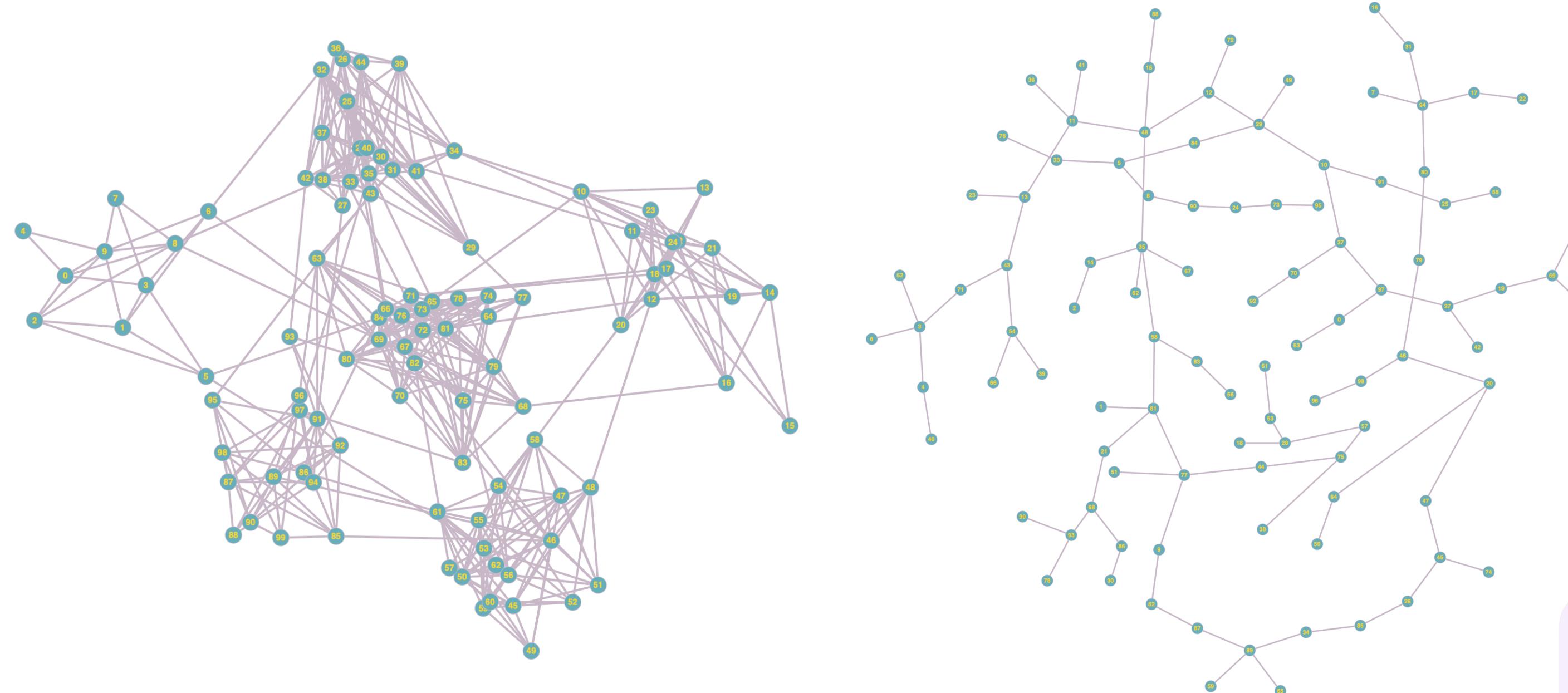
- Transitivity is how often two connected nodes have the same neighbor
- Global metric

High **POSITIVE** correlation with
Hamming.

*More closed triplets meant more
change in regulatory region
sequence.*

Neighborhood Consistency

Clustering vs. Regulatory Region Alignment Length



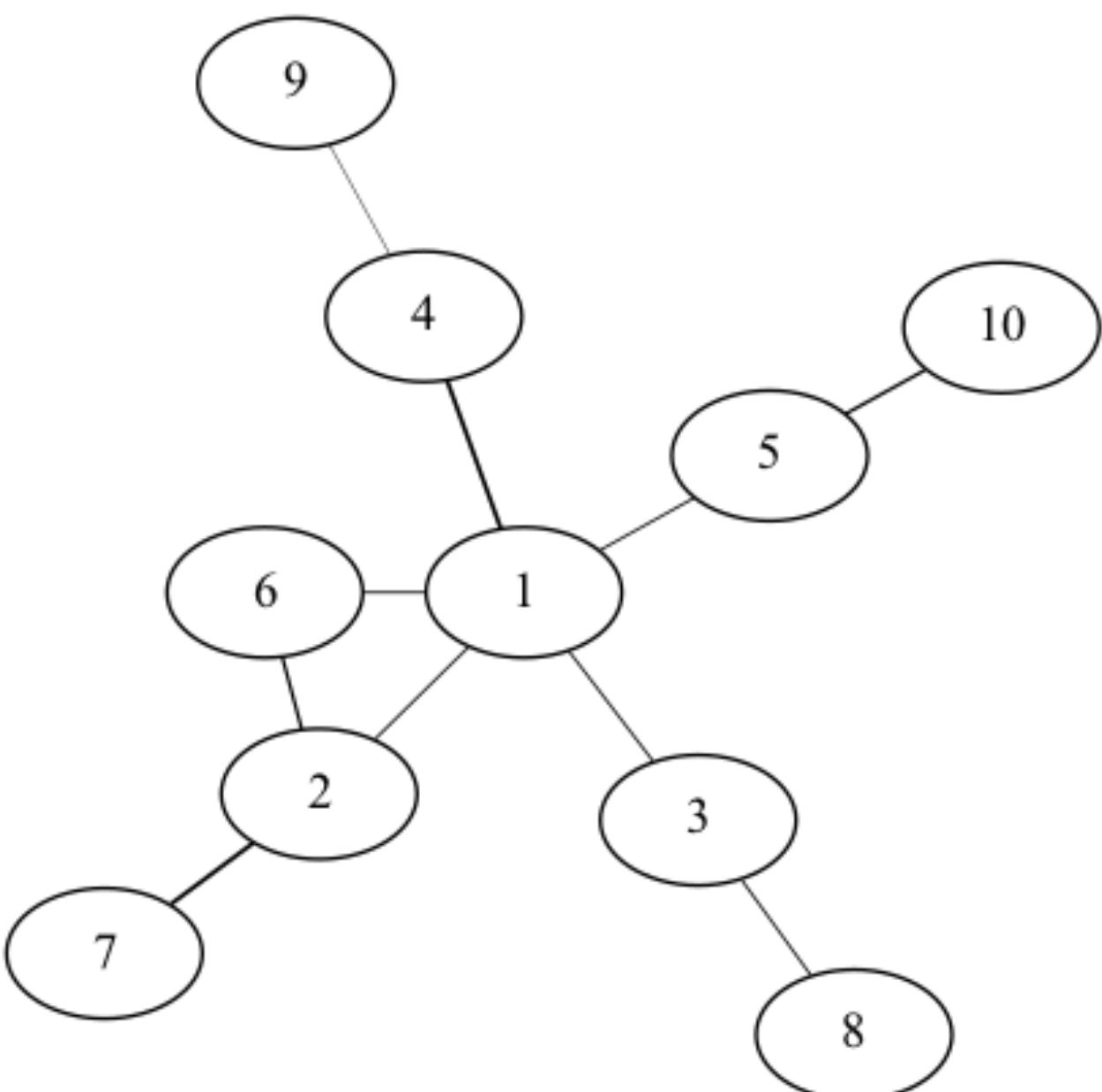
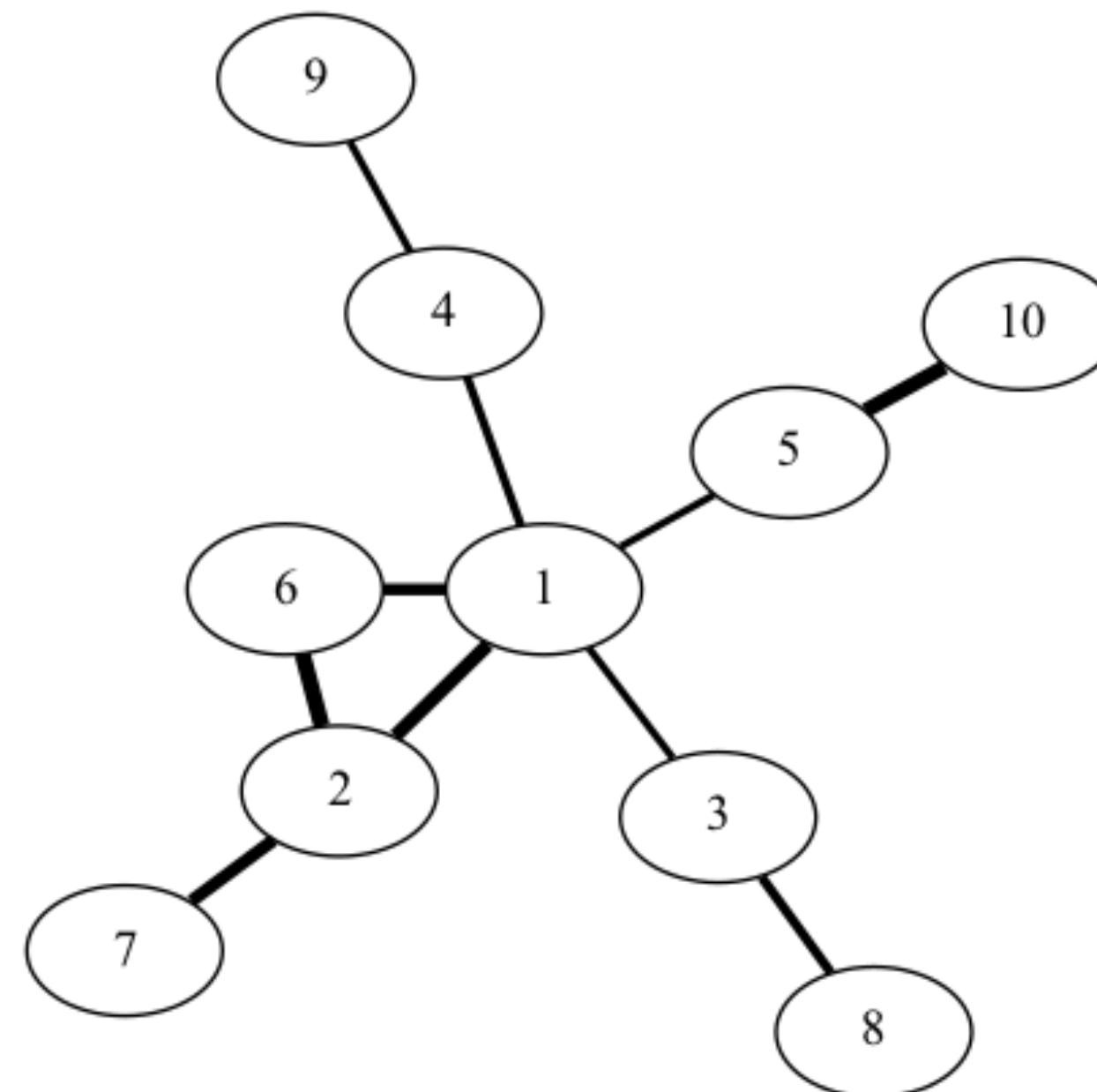
High **POSITIVE** correlation with regulatory region alignment length.

- How clustered are neighborhoods
- Local metric

Higher the clustering, the longer the multiple sequence alignment.

Neighborhood Consistency

Average Edge Weight vs. Conservation



Higher average edge weight indicates more conserved neighborhood across clade

High NEGATIVE correlation with dNdS.

High POSITIVE correlation with Hamming.

More conserved protein-coding sequence means more conserved connections. But why positive correlation with Hamming?

Network Evolution

Further Work

Perform **permutation testing** on edges of networks to see if patterns re-emerge

- Using scale-free model of preferential attachment

Further explore **motif distribution** in regulatory regions using MEME

- Using motif entropy as an analog of change in motif distribution

Network Evolution

Conclusion

- Coding sequence and regulatory regions show positively correlated correlation.
- Genomic conservation correlates with network connectivity.
- Changes in networks are correlated with protein evolution.
 - Interesting potential for exploring relationship between network architecture.



Laura K. Reed

Paige F. Ferguson

Kevin M. Kocot

Michael R. McKain

John H. Yoder

Réka Albert ext.

Daryl Lam

Wilson Leung

Samuel Sledzieski

Faith Ocitti

Joshua Lotfi

Richard Adkins

Reed Lab Members

Logan T. Cohen

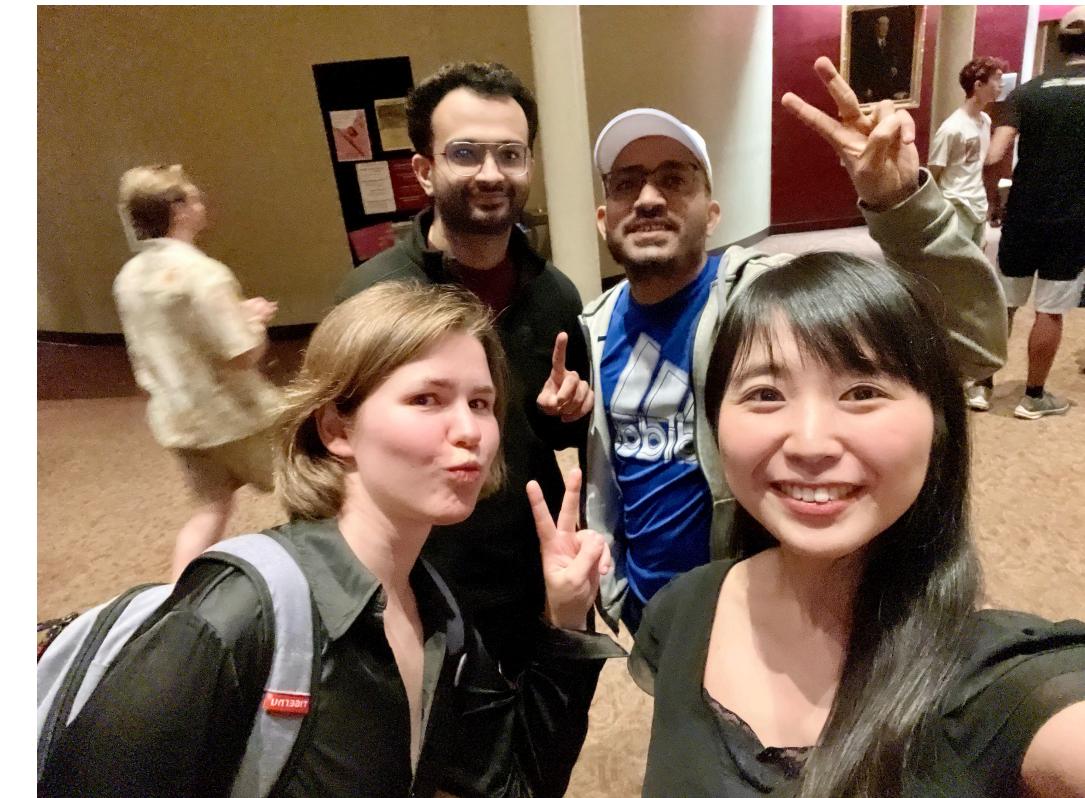
Tolulope R. Kolapo-Mabayoje

Sahar Abdollahi

Bishnu Adhikari

David "Kyle" Breault

Acknowledgements



Anna Parul
Torin A. Alter
German A. Lopez Otero
Hrishikesh S. Tupkar
Friends and Family

AA ANM BNG DH
DMP DP EM ER GK
HA HG HN ID IN JC
JV LAJ MH MKW
MM MS NR PC PS
RJ RN SA SK SP SS



NCTTA/USATT/Senior Games



fin