

# Sparsification of Phylogenetic Covariance Matrices of $k$ -Regular Trees

Sean Svihla

Department of Applied Mathematics  
University of Colorado - Boulder

AofA 2024

# k-Regular Trees

$T = (V, E)$ :

- **planted**
- **ordered**
- **unlabelled**
- **$k$ -regular**,  $k \geq 2$ .

**Leaves** labelled as encountered in DFS.

**Interior nodes** and **edges** labelled by leaf descendants.

**Edge length function**

$\ell : E \rightarrow \mathbb{R}_+$ .

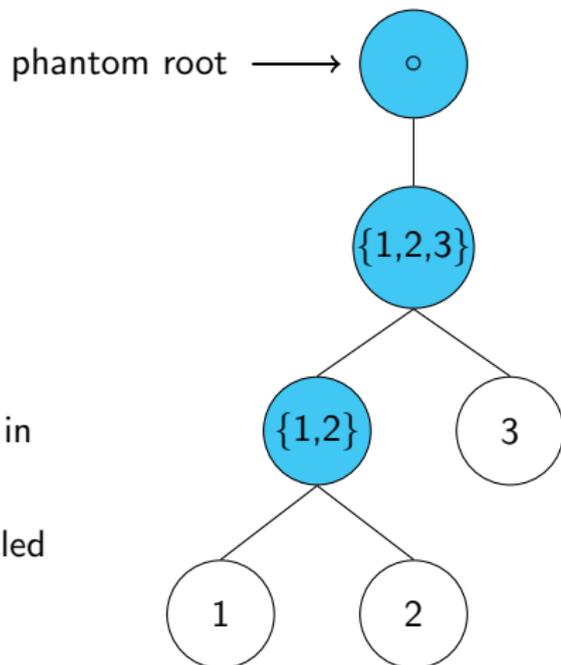


Figure: 2-regular tree example.

# Phylogenetic Covariance Matrices

(a.k.a. Cophenetic Matrices)

**Phylogenetic covariance** is a function of shared path length:

$$C(i, v) = \ell(e_1);$$

$$C(i, j) = \ell(e_2) + \ell(e_1);$$

$$C(i, i) = \ell(e_3) + \ell(e_2) + \ell(e_1).$$

**Phylogenetic covariance matrix:**

$$C = \left( C(i, j) \right)_{i, j \in L}.$$

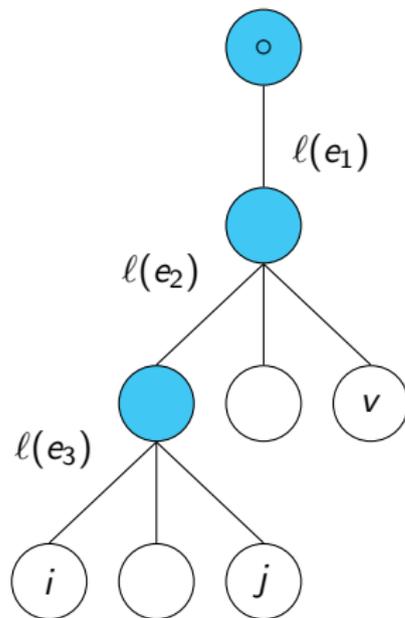


Figure: 3-regular tree with interior nodes colored cyan.

# Phylogenetic Covariance

[Harmon (2019)]

If traits evolve as a **Brownian Motion** (BM) along each edge, i.e.:

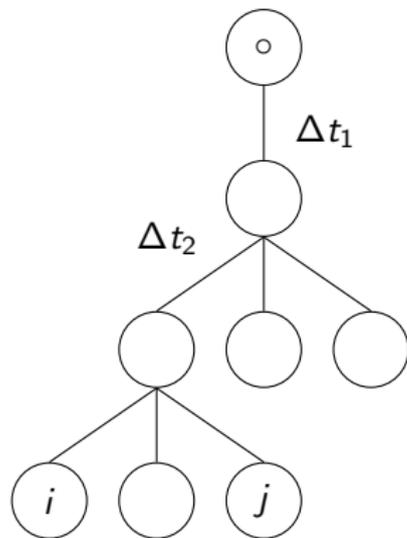


Figure: BM representation of phylogenetic covariance matrices.

$X_v :=$  trait value at a leaf  $v$

$U, V :=$  Gaussian variables

then  $U, V$ , and  $B(\Delta t_1 + \Delta t_2)$  are independent, hence

$$X_i = B(\Delta t_1 + \Delta t_2) + U$$

$$X_j = B(\Delta t_1 + \Delta t_2) + V$$

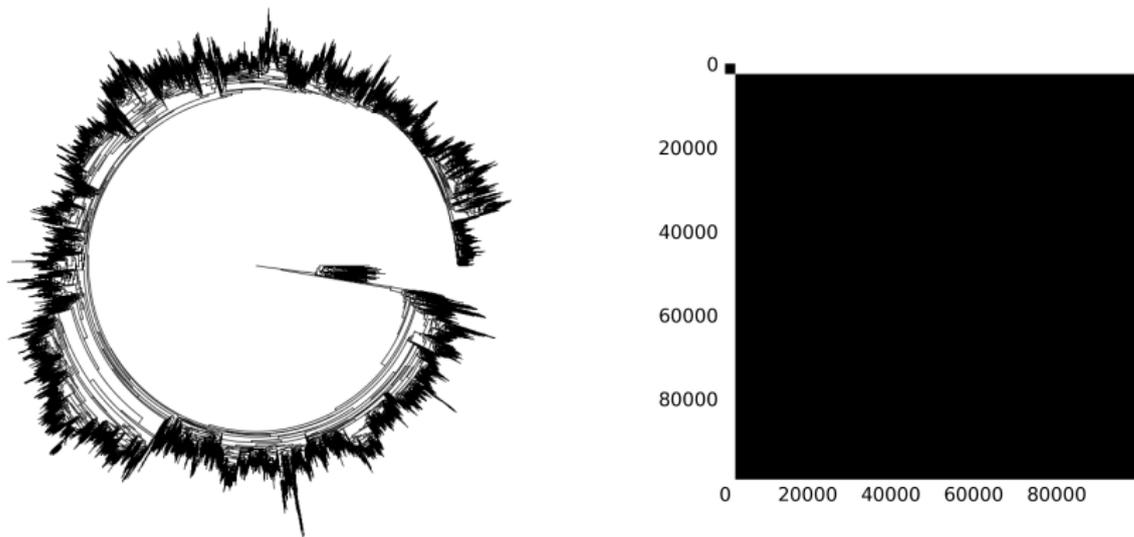
$\Downarrow$

$$\text{Cov}(X_i, X_j) = \text{Var}(B(\Delta t_1 + \Delta t_2))$$

$$= (\Delta t_1 + \Delta t_2) \cdot \sigma^2$$

$$= \ell(e_1) + \ell(e_2).$$

# Phylogenetic covariance matrices are typically dense

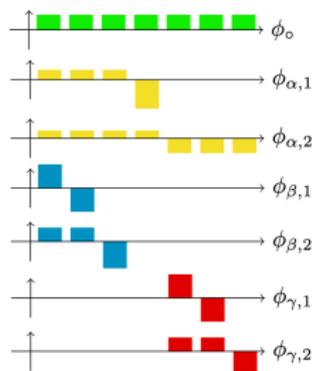
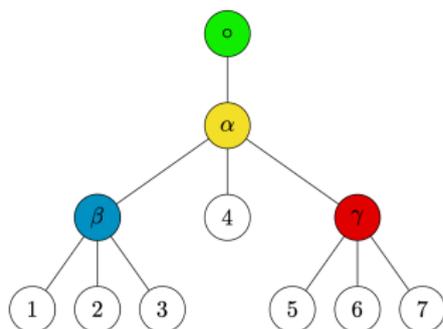


**Figure:** (Left) Circular layout of reference binary phylogenetic tree with  $\approx 100,000$  leaves. (Right) Heatmap of associated covariance matrix. 94% of its  $\approx 10$  billion entries are non-zero.<sup>a</sup>

<sup>a</sup>Figures from [Gorman & Lladser (2023)].

# The Haar-like Wavelets

[Gavish, Nadler & Coifman (2010)]



- **Orthonormal basis** for the linear space of functions  $L \rightarrow \mathbb{R}$ .
- As many wavelets as leaves.
- The wavelet associated with the root is **constant**.

For each  $v \in I \setminus \{o\}$

- there are associated wavelets  $\phi_{v,j}$ ,  $j = 1, 2, \dots, |\text{children}(v)| - 1$ ,
- $\forall j$ ,  $\text{supp}(\phi_{v,j}) \subset L(v)$ , and
- $\forall j$ ,  $\phi_{v,j}$  takes a single positive and single negative value.

# The Haar-like Wavelets

[Gavish, Nadler & Coifman (2010)]

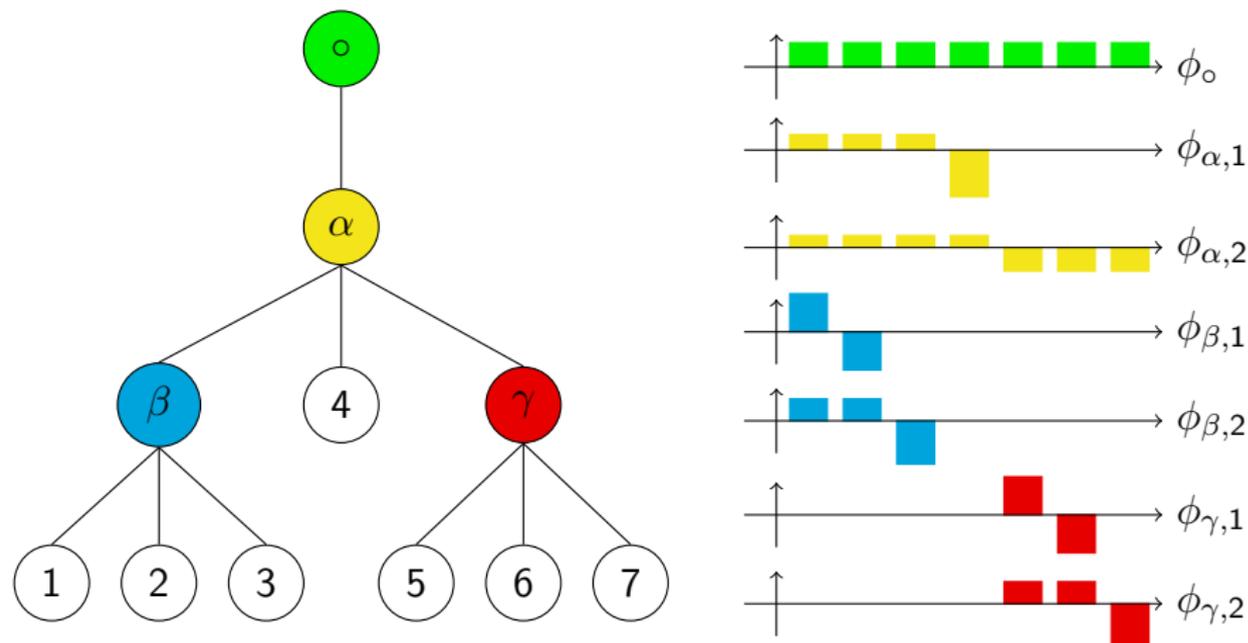


Figure: Haar-like wavelets associated with a 3-regular tree.

# Sparsification of phylogenetic covariance matrices

**Trace branch length** is defined  $\forall e \in E$ ,  $\ell^*(e) := |L(e)| \cdot \ell(e)$ .

**Theorem.**<sup>a</sup> Let  $\Phi$  be the matrix with Haar-like wavelets as columns. If  $u, v \in I$  then

$$(\Phi' C \Phi)(u, v) = \sum_{i \in L} \varphi_u(i) \cdot \ell^*(i, v) \cdot \varphi_v(i).$$

In particular,  $L(u) \cap L(v) = \emptyset \implies (\Phi' C \Phi)(u, v) = 0$ .

---

<sup>a</sup>[Gorman & Lladser (2023)], [Svihla & Lladser (2024)].

## Many wavelet pairs have disjoint support

**Recall:**  $\forall \varphi_u$  wavelet associated with  $u \in I$ ,  $\text{supp}(\varphi_u) \subset L(u)$ .

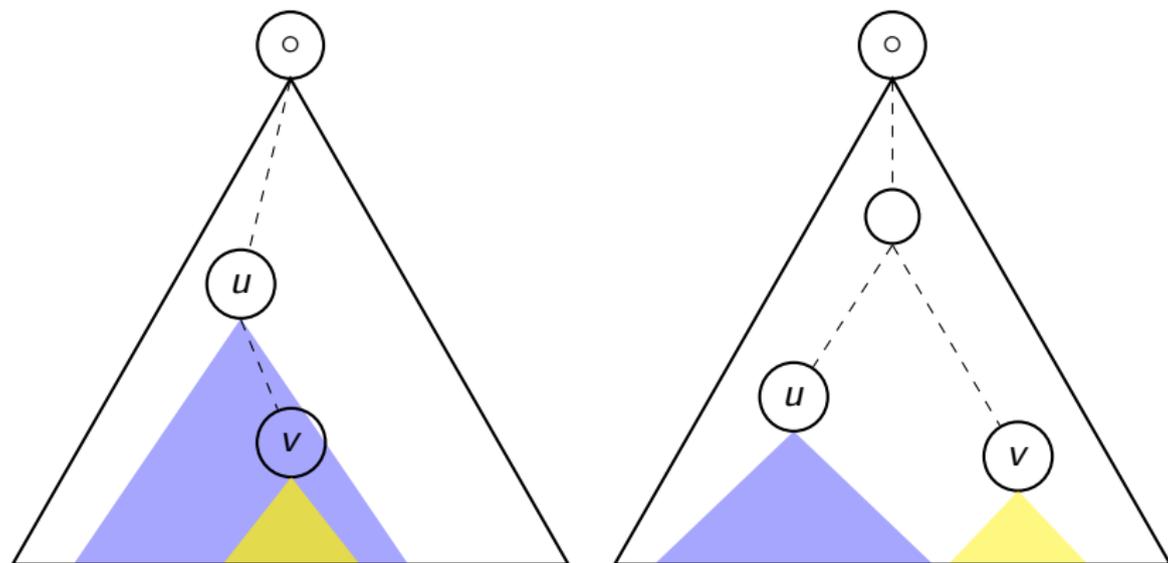
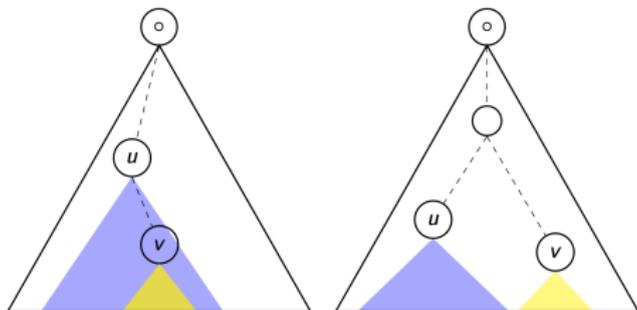


Figure:  $L(u) \cap L(v) \neq \emptyset \iff$  “ $u$  is an ancestor of  $v$  or vice versa”.

# Sparsification of phylogenetic covariance matrices



**Theorem.<sup>b</sup>** For a  $k$ -regular tree  $T$ , if  $\zeta$  denotes the fraction of vanishing entries of  $\Phi' C \Phi$ , then

$$(1 - \zeta) \leq \frac{(k - 1)^2}{|I|} + 2(k - 1)^2 \frac{\text{IPL}(T)}{|I|^2}.$$

In particular, if  $\text{IPL}(T) \ll |I|^2$  then  $\zeta = 1 - o(1)$ .

---

<sup>b</sup>[Svihla & Lladser (2024)].

There's no reason why  $\text{IPL}(T) \ll |I|^2$  should happen!

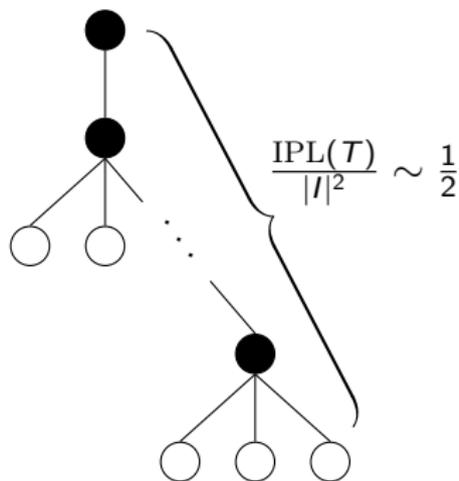


Figure: 3-regular caterpillar tree.

**What about a typical  $k$ -regular tree?**

# Generating function for $k$ -regular trees

Let  $\mathbb{T}_n$  be a **uniformly at random  $k$ -regular tree** of size  $n$ .

**Goal:** Find asymptotic formulas for  $\mathbb{E}[\text{IPL}(\mathbb{T}_n)]$  and  $\mathbb{V}[\text{IPL}(\mathbb{T}_n)]$ .

**Definition.** Let  $Q(z, u)$  be the (bivariate) generating function of the class of  $k$ -regular trees, where  $z$  marks the **size** and  $u$  marks the **internal path length** of each tree.

To address the goal, we need to understand the singularities of

$$Q(z) := Q(z, u) \Big|_{u=1};$$

$$Q_u(z) := \frac{\partial Q}{\partial u}(z, u) \Big|_{u=1}; \quad Q_{uu}(z) := \frac{\partial^2 Q}{\partial u^2}(z, u) \Big|_{u=1}.$$

# Radius of convergence of $Q(z)$

By the **symbolic method**:

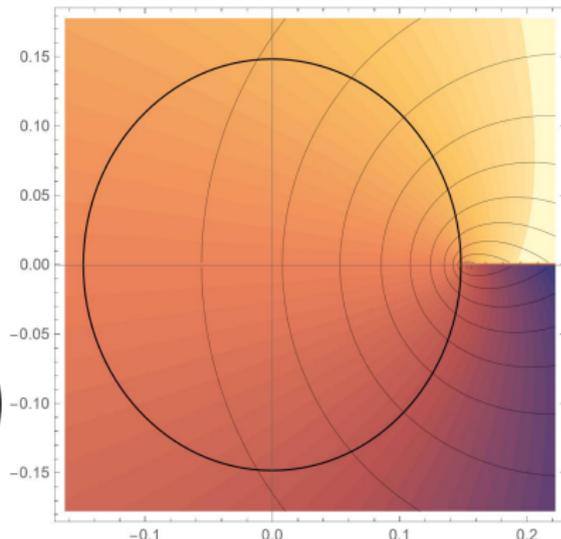
$$Q(z) = 1 + z\{Q(z)\}^k$$

By **Lagrange inversion**:

$$[z^n]Q(z) = \frac{1}{(k-1)n+1} \binom{kn}{(k-1)n}$$

So,  $Q(z)$  has **radius of convergence** and a **singularity** at:

$$z_k := \frac{(k-1)^{k-1}}{k^k}.$$



**Figure:** Plot of  $Q(z)$  and boundary of disk of convergence  $|z| < z_k$  when  $k = 3$ . Plot is colored by argument; contour lines denote modulus.

## Interlude: Hypergeometric functions

**Definition.**  $F(z) = \sum_{n=0}^{\infty} f_n z^n$  is called **hypergeometric** if there are  $a_1, \dots, a_p, b_1, \dots, b_q \in \mathbb{R}$  such that

$$\frac{f_{n+1}}{f_n} = \frac{(n+a_1) \cdots (n+a_p)}{(n+1) \cdot (n+b_1) \cdots (n+b_q)}.$$

In this case, we write:

$$F(z) = {}_pF_q \left[ \begin{matrix} a_1, \dots, a_p \\ b_1, \dots, b_q \end{matrix}; z \right].$$

**Theorem<sup>†</sup>** If the **balance**  $s := \left( \sum_{j=1}^q b_j - \sum_{j=1}^p a_j \right) > 0$  then  ${}_pF_q(a_1, \dots, a_p; b_1, \dots, b_q; z)$  converges at  $z = 1$ .

---

<sup>†</sup>[Evans & Stanton (1984)].

# $Q(z)$ is hypergeometric

## Proposition.<sup>c</sup>

$$Q(z) = {}_{k-1}F_{k-2} \left[ \begin{matrix} \frac{1}{k}, \frac{2}{k}, \dots, \frac{k-1}{k} \\ \frac{2}{k-1}, \frac{3}{k-1}, \dots, \frac{k-2}{k-1}, \frac{k}{k-1} \end{matrix} ; \frac{z}{z_k} \right]$$

and it has balance  $s > 0$ . Furthermore, if we define

$$p(t) := t(1-t)^{k-1},$$

then

$$Q(p(t)) = \frac{1}{1-t}, \text{ for all } 0 \leq t \leq \frac{1}{k}.$$

---

<sup>c</sup>[Weisstein (2023)], [Svihla & Lladser (2024)].

In particular, since  $z_k = p\left(\frac{1}{k}\right)$ ,  $Q(z_k) = \frac{k}{k-1}$ .

$Q(z)$  fits the smooth implicit function schema

$$Q(z_k) = \frac{k}{k-1}$$

**Lemma.**<sup>d</sup>  $z = z_k$  is the only singularity of  $Q(z)$  on  $|z| \leq z_k$ , and

$$Q(z) = 1 + g(z) - h(z) \cdot \sqrt{1 - \frac{z}{z_k}},$$

locally about  $z_k$ , with  $g(z)$  and  $h(z)$  analytic nearby. Additionally,

$$[z^n] Q(z) = \sqrt{\frac{k}{2\pi n^3(k-1)^3}} \cdot z_k^{-n} (1 + O(n^{-1})).$$

---

<sup>d</sup>[Drmota (2009)], [Svihla & Lladser (2024)].

# Partial derivatives of $Q(z, u)$ at $u = 1$

Again due to the **symbolic method**:

$$Q(z, u) = 1 + z \cdot \{Q(zu, u)\}^k.$$

**Implicit differentiation** then gives that  $Q_u(z)$  and  $Q_{uu}(z)$  are **linear combinations** of generating functions of the form

$$\frac{f(z)\{Q(z)\}^a}{(1 - kz\{Q(z)\}^{k-1})^b}.$$

**Lemma.**<sup>e</sup> The equation  $kz\{Q(z)\}^{k-1} = 1$ , with  $|z| \leq z_k$ , has only  $z_k$  as a solution.

---

<sup>e</sup>[Svihla & Lladser (2024)].

## Partial derivatives of $Q(z, u)$ at $u = 1$

**Lemma.**<sup>f</sup> If  $f : \mathbb{C} \rightarrow \mathbb{C}$  is an entire analytic function such that  $f(z_k) \neq 0$ , and  $a \geq 0$  and  $b \geq 1$  are integers, then

$$\begin{aligned} [z^n] \frac{f(z)\{Q(z)\}^a}{(1-kz\{Q(z)\}^{k-1})^b} \\ = \frac{f(z_k)}{2^{b/2}\Gamma(b/2)} \left(\frac{k}{k-1}\right)^{a+b/2} n^{(b-2)/2} z_k^{-n} (1 + O(n^{-1/2})). \end{aligned}$$

<sup>f</sup>[Svihla & Lladser (2024)].

Lastly, for our final and main result are the **well-known formulas**:

$$\mathbb{E}[\text{IPL}(\mathbb{T}_n)] = \frac{[z^n]Q_u(z)}{[z^n]Q(z)};$$

$$\mathbb{V}[\text{IPL}(\mathbb{T}_n)] = \frac{[z^n]Q_{uu}(z)}{[z^n]Q(z)} + \frac{[z^n]Q_u(z)}{[z^n]Q(z)} - \left(\frac{[z^n]Q_u(z)}{[z^n]Q(z)}\right)^2.$$

# Sparsification of a random $k$ -regular tree

**Theorem.**<sup>§</sup> If  $\mathbb{T}_n$  is a uniformly at random  $k$ -regular tree of size  $n$  then

$$\mathbb{E}[\text{IPL}(\mathbb{T}_n)] \sim \sqrt{\frac{\pi k}{2(k-1)}} n^{3/2} \quad \text{and}$$
$$\mathbb{V}[\text{IPL}(\mathbb{T}_n)] \sim \frac{k}{2(k-1)} \left( \frac{10}{3} - \pi \right) n^3.$$

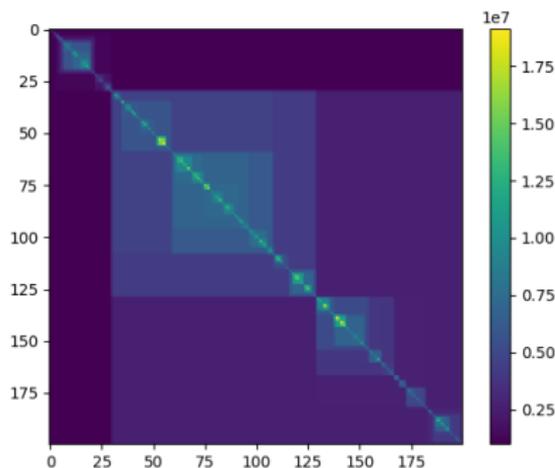
In particular, if  $C_n$  and  $\Phi_n$  are the phylogenetic covariance and Haar-like matrices associated with  $\mathbb{T}_n$ , respectively, and  $\zeta_n$  denotes fraction of vanishing entries in  $\Phi_n' C_n \Phi_n$  then

$$\lim_{n \rightarrow \infty} \zeta_n \stackrel{p}{=} 1.$$

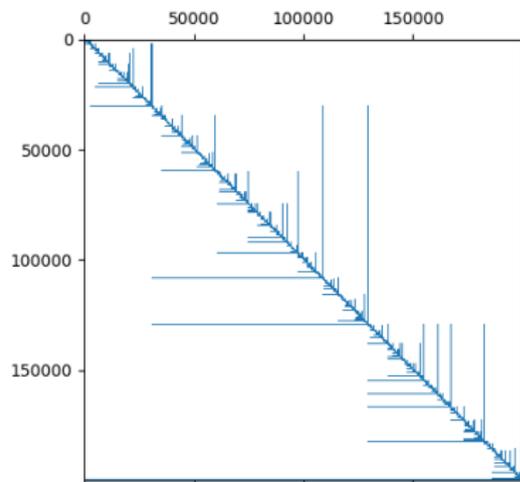
---

<sup>§</sup>[Svihla & Lladser (2024)].

# Does the sparsification work in practice?



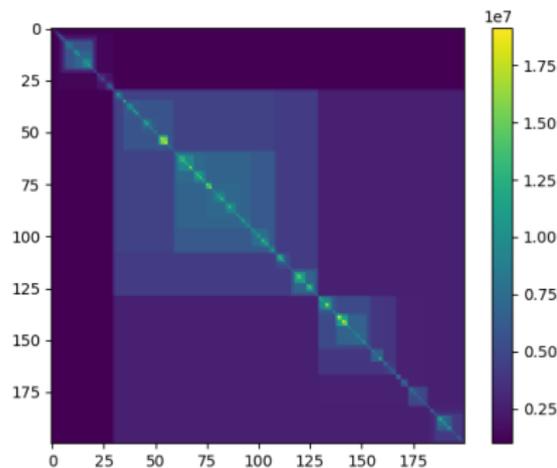
(a)  $200,000 \times 200,000$  dense matrix.



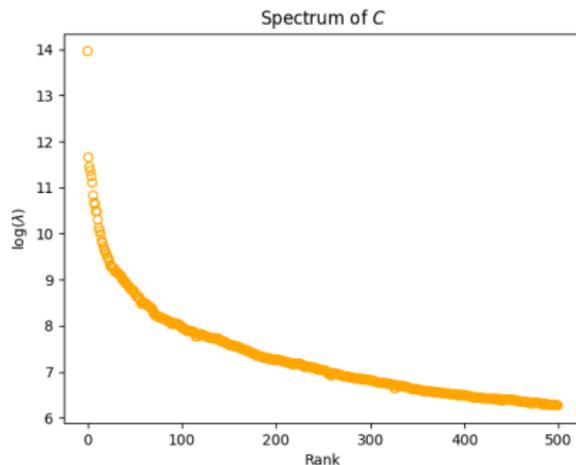
(b) Sparsified matrix.

**Figure:** Density/sparsity pattern of the phylogenetic covariance matrix of a random 3-regular tree with  $\approx 200,000$  leaves. The dense matrix has over 40 billion non-zero entries, but 99.97% of these vanish after changing basis to the Haar-like wavelets of the tree.

# Sparsification enables manipulating large dense matrices



(a)  $200,000 \times 200,000$  dense matrix.



(b) Matrix spectrum.

**Figure:** With over 40 billion non-zero entries, everyday software cannot manipulate the dense matrix on the left. Nevertheless, it can compute the 500 largest eigenvalues from its sparsified version, as shown in the plot on the right.

# Summary

- The Haar-like wavelets can be used to sparsify phylogenetic covariance matrices.
- We can derive a lower-bound on the number of vanishing entries after changing to the Haar-like basis by counting wavelets with disjoint support.
- With high probability, a large random  $k$ -regular tree has a covariance matrix which is highly sparsified by its Haar-like basis.
- This enables manipulating large and dense phylogenetic covariance matrices from their sparsified representation.

# Acknowledgements

- Manuel Lladser (Ph.D. research advisor).
- Lladser Research Group.
- Work partially funded by the NSF grant No. 1836914.

**Thank you!**