

**Document Version**

Final published version

**Licence**

CC BY

**Citation (APA)**

Holtgreffe, N., Huber, K. T., van Iersel, L., Jones, M., Martin, S., & Moulton, V. (2025). Squirrel: Reconstructing semi-directed phylogenetic level-1 networks from four-leaved networks or sequence alignments. *Molecular Biology and Evolution*, 42(4), Article msaf067. <https://doi.org/10.1093/molbev/msaf067>

**Important note**

To cite this publication, please use the final published version (if applicable). Please check the document version above.

**Copyright**

In case the licence states "Dutch Copyright Act (Article 25fa)", this publication was made available Green Open Access via the TU Delft Institutional Repository pursuant to Dutch Copyright Act (Article 25fa, the Taverne amendment). This provision does not affect copyright ownership. Unless copyright is transferred by contract or statute, it remains with the copyright holder.

**Sharing and reuse**

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

**Takedown policy**

Please contact us and provide details if you believe this document breaches copyrights. We will remove access to the work immediately and investigate your claim.

# SQUIRREL: Reconstructing Semi-directed Phylogenetic Level-1 Networks from Four-Leaved Networks or Sequence Alignments

Niels Holtgreve <sup>1</sup>, Katharina T. Huber <sup>2</sup>, Leo van Iersel <sup>1</sup>, Mark Jones <sup>1</sup>, Samuel Martin <sup>3</sup>, Vincent Moulton <sup>2,\*</sup>

<sup>1</sup>Delft Institute of Applied Mathematics, Delft University of Technology, Mekelweg 4, Delft 2628 CD, The Netherlands

<sup>2</sup>School of Computing Sciences, University of East Anglia, Norwich NR4 7TJ, UK

<sup>3</sup>European Bioinformatics Institute, Hinxton CB10 1SD, UK

\*Corresponding author: E-mail: [v.moulton@uea.ac.uk](mailto:v.moulton@uea.ac.uk).

Associate editor: Yoko Satta

## Abstract

With the increasing availability of genomic data, biologists aim to find more accurate descriptions of evolutionary histories influenced by secondary contact, where diverging lineages reconnect before diverging again. Such reticulate evolutionary events can be more accurately represented in phylogenetic networks than in phylogenetic trees. Since the root location of phylogenetic networks cannot be inferred from biological data under several evolutionary models, we consider semi-directed (phylogenetic) networks: partially directed graphs without a root in which the directed edges represent reticulate evolutionary events. By specifying a known outgroup, the rooted topology can be recovered from such networks. We introduce the algorithm SQUIRREL (Semi-directed Quarnet-based Inference to Reconstruct Level-1 Networks) which constructs a semi-directed level-1 network from a full set of quarnets (four-leaf semi-directed networks). Our method also includes a heuristic to construct such a quarnet set directly from sequence alignments. We demonstrate SQUIRREL's performance through simulations and on real sequence data sets, the largest of which contains 29 aligned sequences close to 1.7 Mb long. The resulting networks are obtained on a standard laptop within a few minutes. Lastly, we prove that SQUIRREL is combinatorially consistent: given a full set of quarnets coming from a triangle-free semi-directed level-1 network, it is guaranteed to reconstruct the original network. SQUIRREL is implemented in Python, has an easy-to-use graphical user interface that takes sequence alignments or quarnets as input, and is freely available at <https://github.com/nholtgreve/squirrel>.

**Keywords:** semi-directed phylogenetic network, rooted phylogenetic network, quarnet, traveling salesman problem, sequence alignment, network reconstruction.

## Introduction

Secondary contact, where diverging lineages come into contact and hybridize before continuing to diverge, is commonplace in evolution. This process is poorly described by most phylogenetic reconstruction methods which generally assume a bifurcating tree model. Secondary contact has been widely documented for diverse sets of taxa, including viruses (e.g. HIV and SARS-CoV-2, see [Worobey et al. 2008](#); [Pekar et al. 2021](#); [Jiao et al. 2024](#)), bacteria (e.g. [Diop et al. 2022](#)), plants (e.g. [Ehrendorfer 1959](#); [Rieseberg et al. 2003](#)), birds (e.g. [Taylor and Larson 2019](#)), fish (e.g. [Meier et al. 2019](#); [Du et al. 2024](#)), invertebrates (e.g. [Zhang et al. 2016](#)), and primates, including humans (e.g. [Patterson et al. 2006](#); [Green et al. 2010](#)). Through secondary contact, introgression—the exchange of genetic material between hybridizing lineages—may occur by means of complex processes, often involving multiple rounds of backcrossing.

Evolutionary histories shaped by secondary contact can be more accurately represented by rooted phylogenetic level-1 networks than by strictly bifurcating rooted phylogenetic trees. Rooted phylogenetic level-1 networks are directed acyclic graphs that are largely tree-like in structure, describing patterns of divergence, but include localized reticulations where

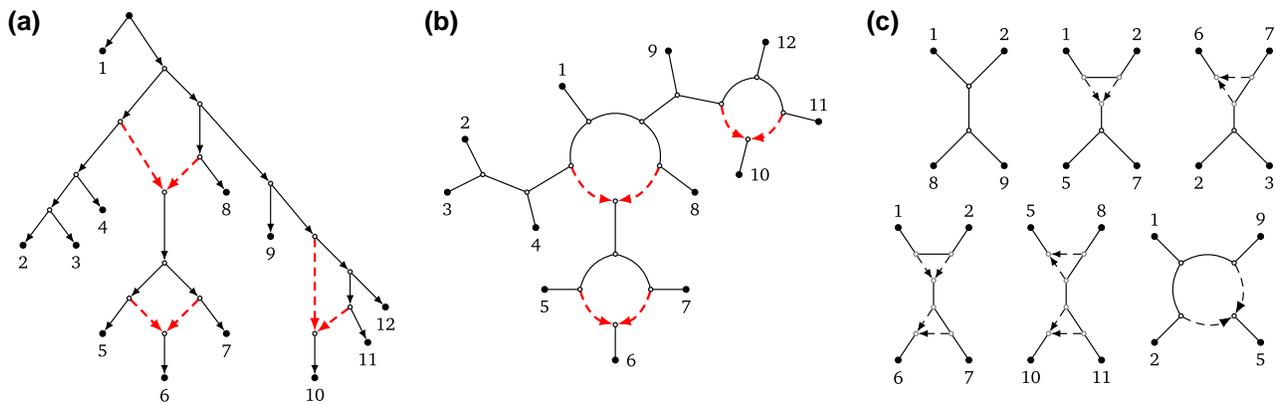
lineages have merged through reticulate events (see, e.g. [Fig. 1a](#) and see the Materials and Methods for a more formal definition). Application of these networks is highly desirable, but their construction is computationally intensive, and their use has remained out of reach for most biologists. Results reported here, including an efficient algorithm and software, address the challenge of building phylogenetic level-1 networks, thus offering the possibility of finding a more realistic description of biological diversity.

Our results are achieved by considering *semi-directed (phylogenetic) networks* ([Solís-Lemus and Ané 2016](#)), in which there is no root and only branches representing reticulate events carry information about direction (see the Materials and Methods for a more formal definition). These networks have gained considerable interest recently (see, e.g. [Solís-Lemus and Ané 2016](#); [Allman et al. 2019](#); [Kong et al. 2024](#); [Warnow et al. 2024](#); [Wu and Solís-Lemus 2024](#); [Frohn et al. 2025](#)), as it has been shown that under certain models of evolution it is theoretically impossible to infer the root of a rooted phylogenetic network directly from data ([Baños 2019](#); [Gross et al. 2021](#); [Xu and Ané 2023](#)). For an example of a semi-directed level-1 network, see [Fig. 1b](#). In case an outgroup is available, this can be used to root the semi-

Received: November 7, 2024. Revised: January 21, 2025. Accepted: March 4, 2025

© The Author(s) 2025. Published by Oxford University Press on behalf of Society for Molecular Biology and Evolution.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.



**Fig. 1.** a) A rooted phylogenetic level-1 network on 12 taxa represented by numbers 1–12, with the dashed reticulation edges pointing towards reticulation vertices which represent reticulate events. b) The semi-directed topology of the rooted network, which is a triangle-free semi-directed level-1 network on 12 leaves, again with the reticulation edges dashed. This network uniquely determines the rooted network by specifying leaf 1 as an outgroup. c) Some of the quarnets induced by the semi-directed network. When ignoring the leaf labels, these are all six possible level-1 quarnet shapes. The top left quarnet is a quartet tree, the bottom right quarnet is the only one that contains a cycle of length 4 (4-cycle), and the other four quarnets contain one or two triangles (3-cycles). The tf-quarnets (triangle-free quarnets) can be obtained from the quarnets by contracting each of the triangles to a single node. The quartet tree and quarnet with a 4-cycle are both already triangle-free.

directed network (Solís-Lemus and Ané 2016), as illustrated in Fig. 1a and b. Several identifiability results have been recently proven for semi-directed level-1 networks. In particular, it was shown that such networks can be theoretically recovered from data under various models of evolution (Baños 2019; Gross et al. 2021; Xu and Ané 2023). By focusing on semi-directed networks, we offer a tractable way for reconstructing phylogenetic level-1 networks.

Recently, two algebraic approaches have been introduced to construct semi-directed level-1, four-leaved networks, or *quarnets* (see Fig. 1c): QNR-SVM (Barton et al. 2022) and an algorithm in Martin et al. (2023). These methods take as input sequence data and both employ algebraic invariants to infer quarnets under the Jukes–Cantor model (Barton et al. 2022; Martin et al. 2023) and the Kimura 2-parameter model (Martin et al. 2023). To infer evolutionary relationships for larger data sets, methods are therefore required to puzzle together such quarnets into larger networks (see, e.g. Schmidt et al. 2002 and Oldman et al. 2016 for two of the earliest algorithms where this approach was used for trees and rooted networks, respectively). It is known that the quarnets coming from a semi-directed level-1 network uniquely characterize the network (Huber et al. 2024) and that theoretically they can be puzzled together efficiently to reconstruct the network (Frohn et al. 2025). However, a set of quarnets stemming from real data will unavoidably contain erroneous quarnets, thus creating the need for a more robust algorithm.

In this paper, we introduce SQUIRREL (Semi-directed Quarnet-based Inference to Reconstruct Level-1 Networks): an efficient software tool and algorithm that builds a semi-directed level-1 network from a given full set of quarnets (that is, a *dense* set that contains one quarnet for each subset of four taxa). We complement SQUIRREL with a fast heuristic method to construct quarnets from sequence data: the  $\delta$ -heuristic (see the Materials and Methods for a formal description). Note that various existing algorithms and programs can be used to infer level-1 networks (both rooted and semi-directed) from biological data that are based on alternative approaches. For example, PHYLONET (Than et al. 2008; Yu and Nakhleh 2015), SNAQ (Solís-Lemus and Ané 2016; Solís-Lemus et al. 2017), and PHYNEST (Kong et al. 2024) are all software tools using likelihood-based algorithms operating under a coalescent model.

SNAQ builds semi-directed networks, whereas both PHYNEST and PHYLONET focus on rooted networks. These methods assume an upper bound on the number of reticulate events and either take gene trees (PHYLONET and SNAQ) or sequence data (PHYNEST) as input, after which they perform a potentially time-consuming search through the space of networks to optimize a likelihood criterion. On the other hand, NANUQ (Allman et al. 2019) and the recent extension NANUQ<sup>+</sup> (Allman et al. 2024b) do not employ a likelihood framework and instead use concordance factors on four-taxon subsets to produce a semi-directed level-1 network up to contracting triangles (3-cycles) and identifying the locations of reticulations in 4-cycles. This approach is faster but requires other methods to compute the input gene trees first, which itself can be a challenging step (Chifman and Kubatko 2014; Simmons and Gatesy 2015; Zhang and Mirarab 2022; Steenwyk et al. 2023). Other approaches use Bayesian methodology to construct rooted networks [e.g. SPECIESNETWORK {Zhang et al. 2018a}] but are not yet able to scale to larger data sets. Lastly, LEV1ATHAN (Huber et al. 2010) and TRILONET (Oldman et al. 2016) take a combinatorial stance towards the network construction problem; they take as input a set of rooted three-leaf trees (LEV1ATHAN) or rooted three-leaf networks (TRILONET) and output a rooted level-1 network, with TRILONET including a heuristic to generate rooted three-leaf networks from sequence data.

We now present a brief overview of how SQUIRREL works; a formal description of the algorithm (plus supporting figures) is given in the Materials and Methods section. As with NANUQ and to a lesser extent SNAQ, SQUIRREL constructs networks up to the contraction of triangles (see Fig. 1b), thus resulting in a binary triangle-free semi-directed level-1 network (i.e. a network with no cycles that contain just three vertices). Since triangles are relatively difficult to infer correctly (Gross et al. 2021), SQUIRREL does not use the location of any triangles in the quarnets and instead only employs *tf-quarnets* (triangle-free quarnets; see Fig. 1c). As shown in Frohn et al. (2025), by considering *tf-quarnets*, we still maintain enough information to theoretically construct the complete semi-directed level-1 network up to contracting its triangles. If quarnets with triangles are given in the input, *tf-quarnets* are obtained by contracting the triangles. Hence, each *tf-quarnet* is either a quartet tree or contains a 4-cycle.

Given a dense set of weighted tf-quarnets, SQUIRREL first uses all of the tf-quarnets that are quartet trees to build a sequence of nonbinary phylogenetic trees, using an algorithm from [Berry and Gascuel \(2000\)](#) and employing techniques from the QUARTETJOINING algorithm ([Grünwald et al. 2009](#)) that constructs phylogenetic trees from quartet trees. Within each of the nonbinary phylogenetic trees in the sequence, the internal vertices with high degree are replaced by a suitable cycle. In particular, SQUIRREL repeatedly solves the TRAVELING SALESMAN PROBLEM (TSP, see, e.g. [Bellman 1962](#); [Held and Karp 1962](#)) with suitably defined distances to create a cyclic ordering of the subnetworks around the cycles. This results in a sequence of candidate level-1 networks, from which SQUIRREL returns the one that agrees, in a well-defined sense, with most of the original tf-quarnets. If an outgroup is specified, this network can in turn be transformed into a rooted network.

We emphasize that any method that is able to create a dense set of tf-quarnets from biological data (possibly incorporating, e.g. incomplete lineage sorting) could be used to generate input for SQUIRREL. Furthermore, SQUIRREL takes into account weights the tf-quarnets might have, which can be used to model confidence or bootstrap support. Reassuringly, SQUIRREL is consistent in the sense that it will reconstruct the correct network if all tf-quarnets are derived from a triangle-free semi-directed level-1 network, a fact that we prove in [Theorem 1](#) in the Materials and Methods section.

## Results

### Simulation Study

Following the simulation studies for LEV1ATHAN ([Huber et al. 2010](#)) and TRILONET ([Oldman et al. 2016](#)), we analyze what effect noise in a set of tf-quarnets has on the performance of SQUIRREL. To this end, we generate 100 random triangle-free semi-directed level-1 networks for every number  $n \in \{10, 15, 20, 25, 30, 35\}$  of leaves (see [supplementary material Section B, Supplementary Material](#) online for the generating algorithm). For each network  $\mathcal{N}$ , the reticulation number  $r(\mathcal{N})$  (i.e. the number of reticulations) is chosen uniformly at random from  $\{0, \dots, \lfloor n/3 \rfloor\}$ . This results in a set of 600 random networks  $\mathcal{N}$ , each inducing a set  $\mathcal{Q}(\mathcal{N})$  of tf-quarnets. For each network  $\mathcal{N}$  and each perturbation ratio  $\varepsilon \in \{0, 0.01, 0.02, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5\}$ , we create a noisy set of tf-quarnets  $\mathcal{Q}_\varepsilon(\mathcal{N})$  by changing the undirected underlying topology of a fraction of the tf-quarnets uniformly at random which is given by  $\varepsilon$ . Then, if this creates a 4-cycle, we pick a random location for the reticulation. We use this scheme for the creation of noise to prevent 4-cycles from only changing their reticulation and keeping their circular ordering. Such a perturbation will barely influence the output of the algorithm, since reticulations of 4-cycle tf-quarnets are only used to determine the location of reticulations in 4-cycles of the final networks. The resulting  $5,400 = 600 \times 9$  sets of unweighted tf-quarnets  $\mathcal{Q}_\varepsilon(\mathcal{N})$  are used as input for SQUIRREL. The average computation times ranged from below a second for the networks with the fewest leaves to below two minutes for the networks with 35 leaves.

To measure how well SQUIRREL reconstructs the original networks from these noisy tf-quarnet sets, we compute two similarity scores for every input network  $\mathcal{N}$  and output network  $\mathcal{M}$ . The first score is the *tf-quarnet consistency score* (modeled

after a similar score in [Huber et al. 2010](#) and [Oldman et al. 2016](#)) which is defined as

$$C(\mathcal{N}, \mathcal{M}) = \frac{|\mathcal{Q}(\mathcal{N}) \cap \mathcal{Q}(\mathcal{M})|}{|\mathcal{Q}(\mathcal{N})|}. \quad (1)$$

This score measures what fraction of the tf-quarnets induced by  $\mathcal{N}$  are also induced by the constructed network  $\mathcal{M}$ . We also consider its symmetric counterpart: the *tf-quarnet symmetric consistency score*, defined as

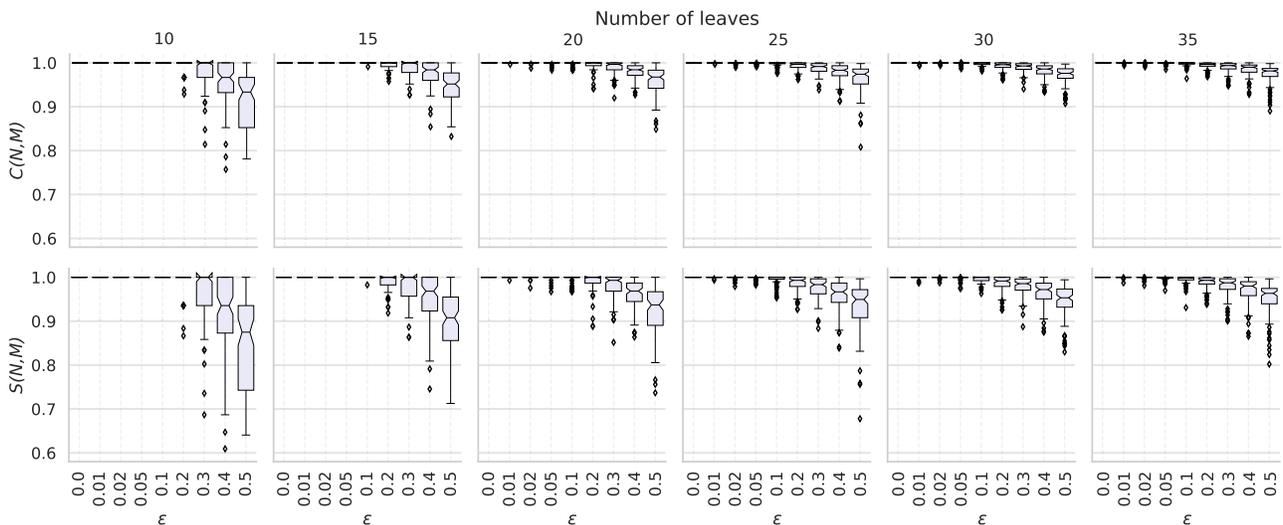
$$S(\mathcal{N}, \mathcal{M}) = \frac{|\mathcal{Q}(\mathcal{N}) \cap \mathcal{Q}(\mathcal{M})|}{|\mathcal{Q}(\mathcal{N}) \cup \mathcal{Q}(\mathcal{M})|}. \quad (2)$$

Both scores are always in the interval  $[0, 1]$  and attain a value of 1 if and only if  $\mathcal{N} = \mathcal{M}$ , which follows from [Frohn et al. \(2025\)](#). The boxplots in [Fig. 2](#) show the distribution of the two scores for different perturbation ratios  $\varepsilon$  and leaf set sizes  $n$ . As expected, both scores decrease for larger values of  $\varepsilon$ . However, the decrease seems fairly limited, with both consistency scores averaging above 0.91 even for sets containing only 50% of the original tf-quarnets.

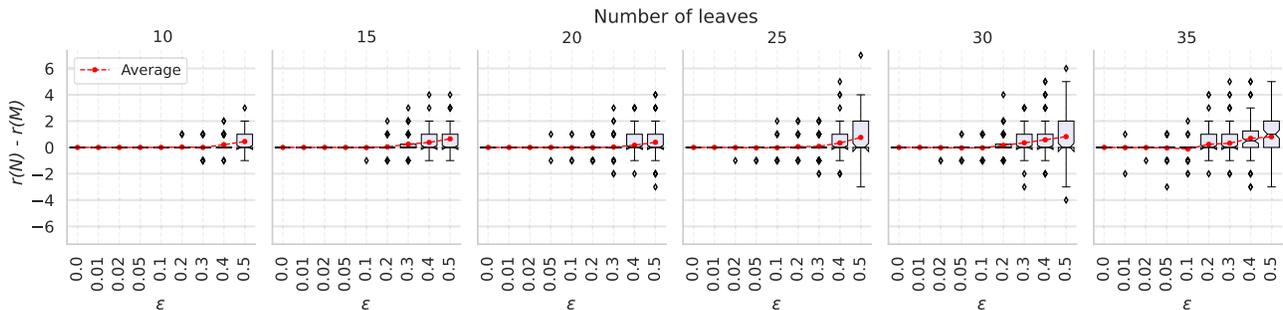
To investigate in what way noise in a set of tf-quarnets influences the structure of the reconstructed networks, we compute the difference in the reticulation numbers  $r(\mathcal{N}) - r(\mathcal{M})$  between the input networks  $\mathcal{N}$  and output networks  $\mathcal{M}$ . The boxplots in [Fig. 3](#) show the result of this experiment, again for different values of  $\varepsilon$  and  $n$ . Up to a value of  $\varepsilon = 0.1$ , SQUIRREL reconstructs networks with the correct reticulation number in almost all cases. For higher values, the differences are more spread out, while the average difference slowly becomes positive. Thus, it seems that SQUIRREL slightly favors networks with fewer reticulations for high values of  $\varepsilon$ , although the average absolute differences remain below a reasonably small 1.5. A possible explanation could be that by not considering triangles in the quarnets, the signal in the data indicating reticulate events is weakened.

We also perform a study with simulated nucleotide sequences to test the performance of the  $\delta$ -heuristic combined with SQUIRREL, using a similar approach to the simulations presented in [Holland et al. \(2002\)](#) and [Oldman et al. \(2016\)](#). For each of our 600 previously generated networks, we simulate one multiple sequence alignment (MSA) for every sequence length  $k \in \{1, 10, 100 \text{ kb}, 1 \text{ Mb}\}$  as follows. Briefly, we first root every semi-directed network  $\mathcal{N}$  uniformly at random on some edge (making sure that it is a valid root location) to create a rooted phylogenetic network. We then use the software tool SEQ-GEN ([Rambaut and Grass 1997](#)) to simulate MSAs of equal length along all displayed trees of the rooted phylogenetic network under the K2P model with transition-transversion bias 4 (as in [Holland et al. 2002](#); [Oldman et al. 2016](#)). The MSAs of the displayed trees are then concatenated to create one MSA with the desired length  $k$ . Since our  $\delta$ -heuristic treats every site of the MSA independently, this way of generating MSAs is asymptotically equivalent to generating MSAs under the K2P network-based Markov model with reticulation parameters of 0.5 (see, e.g. [Gross et al. 2021](#)).

The branch lengths (i.e. the expected number of substitutions along each edge) that are used for the simulations are determined as follows. Given an edge  $(u, v)$  of one of the rooted phylogenetic networks, we let  $p_{(u,v)}$  be the average length (in terms of number of edges) of all unique paths from the root to any leaf that contain the edge  $(u, v)$ . Then, we assign the edge  $(u, v)$  a branch length of  $0.3/p_{(u,v)}$ , which ensures that



**Fig. 2.** Boxplots showing the spread of  $C$ - and  $S$ -scores between the input network  $\mathcal{N}$  and output network  $\mathcal{M}$ , when applying SQUIRREL to sets of tf-quarnets with leaf set sizes  $n$  and perturbation ratios  $\epsilon$ . The boxplots show the quartiles of the data and its outliers. A single outlier in the case of  $n = 10$  and  $\epsilon = 0.5$  has a  $C$ - and  $S$ -score below 0.6 and is omitted from the figure for clarity.



**Fig. 3.** Boxplots showing the variation of the difference in reticulation number  $r(\mathcal{N}) - r(\mathcal{M})$  of the input network  $\mathcal{N}$  and output network  $\mathcal{M}$ , when applying SQUIRREL to sets of tf-quarnets with leaf set sizes  $n$  and perturbation ratios  $\epsilon$ . The boxplots show the quartiles of the data and its outliers.

every path in the network from a root to a leaf roughly has a total length of 0.3, as is the case in the simulations by [Holland et al. \(2002\)](#) and [Oldman et al. \(2016\)](#).

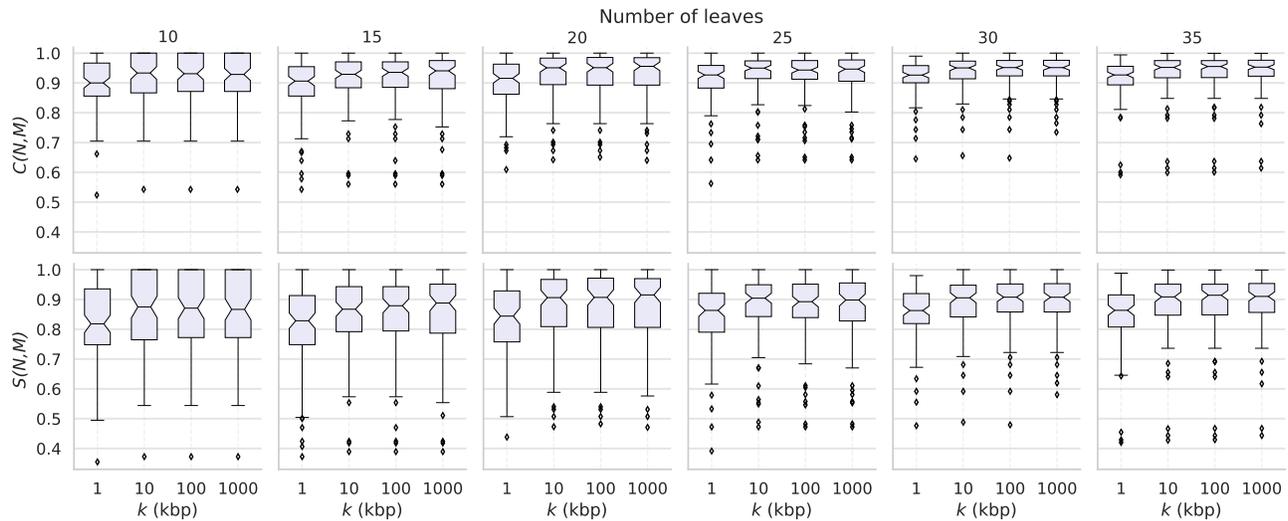
We then use the  $2,400 = 600 \times 4$  simulated MSAs as input for our  $\delta$ -heuristic to construct dense sets of weighted tf-quarnets, which are in turn used to construct semi-directed networks with SQUIRREL. As before, we compare every constructed semi-directed network  $\mathcal{M}$  with the original semi-directed network  $\mathcal{N}$  in terms of  $C$ -score,  $S$ -score and difference in reticulation number  $r(\mathcal{N}) - r(\mathcal{M})$ . The results are depicted in [Figs. 4](#) and [5](#), respectively. We observe that both consistency scores increase as the sequence length changes from 1 to 10 kb. Additionally, both the average and the variation of the difference in reticulation number decrease. Interestingly, the increase of the sequence length from 10 to 100 kb or 1 Mb does not seem to have much further effect. As was the case in our previous experiment, an increase in the number of leaves  $n$  of the original semi-directed network improves the two considered consistency scores, yet also results in a greater spread of the difference in reticulation number between the original and constructed network. The latter point can be explained by the fact that smaller networks simply allow for fewer reticulations, thus also bounding the largest possible difference in reticulation number.

## Biological Data

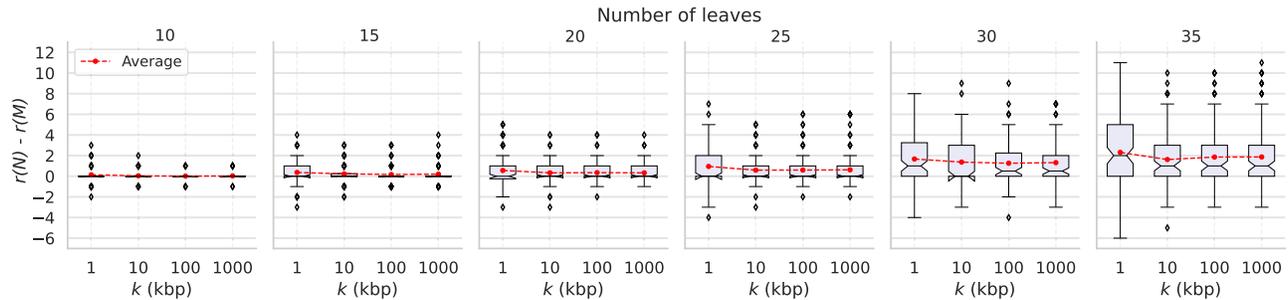
To illustrate the applicability of SQUIRREL to biological data, we consider three data sets on groups of taxa with evidence of secondary contact in their evolutionary histories: a large set of tf-quarnets generated with the MML algorithm from [Martin et al. \(2023\)](#) (named after the authors), a short MSA on few taxa from [Salemi and Vandamme \(2003\)](#), and a long MSA on many taxa from [Vanderpool et al. \(2020\)](#).

**Xiphophorus.** We first test the applicability of SQUIRREL to a set of tf-quarnets that was generated with the MML algorithm ([Martin et al. 2023](#)). For each four-taxon subset, this algorithm creates a ranking of the possible 4-cycles according to some scoring criterion (with the lowest score being the best). Based on the scores, it either detects a quartet tree (which we give a weight of 1), or it chooses the best 4-cycle, which we give a weight of  $\min(1, s_2/s_1 - 1)$ , where  $s_1, s_2$  are the two lowest (and thus best) scores. In this manner, we take into account how close the scores for the two best scoring 4-cycles are.

The data set we consider contains 14,950 weighted tf-quarnets on a set of 25 swordtail fish and platyfish (genus *Xiphophorus*) and the single outgroup *Pseudoxiphophorus*



**Fig. 4.** Boxplots showing the spread of  $C$ - and  $S$ -scores between the input network  $\mathcal{N}$  and output network  $\mathcal{M}$ , when applying the  $\delta$ -heuristic and SQUIRREL to MSAs with leaf set sizes  $n$  and sequence lengths  $k$ . The boxplots show the quartiles of the data and its outliers.



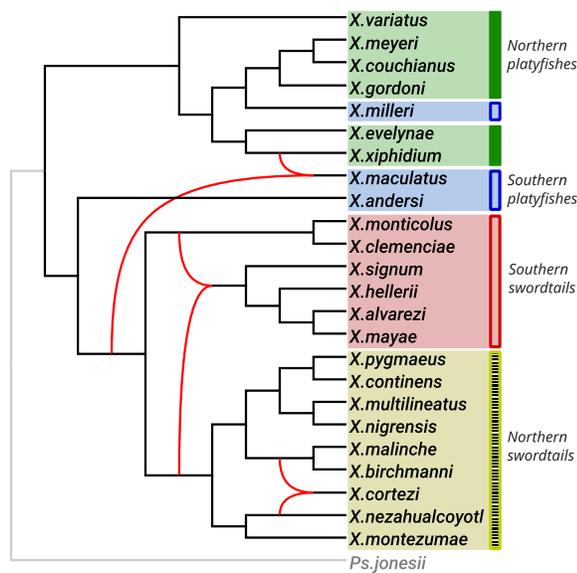
**Fig. 5.** Boxplots showing the variation of the difference in reticulation number  $r(\mathcal{N}) - r(\mathcal{M})$  of the input network  $\mathcal{N}$  and output network  $\mathcal{M}$ , when applying the  $\delta$ -heuristic and SQUIRREL to MSAs with leaf set sizes  $n$  and sequence lengths  $k$ . The boxplots show the quartiles of the data and its outliers.

*jonesii*. This genus has been widely studied and much evidence has been presented for widespread hybridization within the genus (see, e.g. Rosenthal et al. 2003; Culumber et al. 2011; Cui et al. 2013; Kang et al. 2013; Schumer et al. 2013; Solís-Lemus and Ané 2016, and the references therein), making it difficult to capture the full evolutionary history. Traditionally, the genus is divided into four major lineages: northern swordtails, southern swordtails, northern platyfishes, and southern platyfishes (Meyer et al. 2006; Cui et al. 2013). The best network generated by SQUIRREL (taking less than two minutes) had a weighted tf-quarnet consistency score of 0.974 and is shown in Fig. 6. However, many of the other candidate networks had scores that were very close to the score of the best scoring network.

Since the weighted tf-quarnet consistency score measures how consistent the network is with the tf-quarnets, taking their weights into account [see equation (3) in the Materials and Methods], it should be noted that a weighted consistency score close to 1 does not necessarily imply a close to 100% level of confidence that the network is correct. Instead, it reflects whether the quarnets with high weight (i.e. high confidence in their correctness) are consistent with the constructed network, making it most useful as a relative measure to assess if there is a clear best network or if multiple networks perform similarly well. In contrast, the unweighted consistency score [see equation (1)] can be more easily interpreted as an absolute measure of performance, but it may discard useful information about

quarnet confidence if such information is available. A more statistically sound way to generate weights for the tf-quarnets inferred with the MML algorithm from Martin et al. (2023) (similar to the bootstrap support in Barton et al. 2022) would possibly increase the confidence of SQUIRREL in a single best network. Hence, we would welcome further research efforts into computing confidence scores for inferred tf-quarnets which can be used as input weights for SQUIRREL.

The constructed network clearly divides the three major *Xiphophorus* clades (northern swordtails, southern swordtails, and platyfishes) but similar to other studies (Meyer et al. 2006; Cui et al. 2013) intertwines northern and southern platyfishes. Our network has one reticulation edge involving an ancestor of both the northern and the southern swordtails. Another reticulate event places the northern swordtail *Xiphophorus cortezi* both as a sibling of *Xiphophorus nezahualcoyotl* and of the clade (*Xiphophorus malinche*, *Xiphophorus birchmanni*). This reticulate event aligns with previous work in Cui et al. (2013), where the precise placement of *X. cortezi* within this subset of the species (including *Xiphophorus montezumae*) was also uncertain and depended on the inference methods used. Furthermore, one of the subtrees displayed in our network for this subset of the species (i.e. the subtree that includes *X. montezumae*) is the same as the subtree of the network inferred by SNAQ (Solís-Lemus and Ané 2016; Solís-Lemus et al. 2017). The last reticulate event involves the southern platyfish *Xiphophorus maculatus*,



**Fig. 6.** Phylogenetic network inferred by SQUIRREL from a dense set of weighted tf-quarnets on the genus *Xiphophorus* (generated from a MSA with the MML algorithm from Martin et al. 2023). The four major lineages are indicated by the different shaded areas. The reticulation edges are curved, while the edges leading to the outgroup *Pseudoxiphophorus jonesii* are in grey.

for which Cui et al. (2013) report difficulties placing it in the mitochondrial DNA tree. Judging from the many different inferred networks and possible reticulate events (see again Rosenthal et al. 2003; Culumber et al. 2011; Cui et al. 2013; Kang et al. 2013; Schumer et al. 2013; Solís-Lemus and Ané 2016), capturing the evolutionary history of the complete genus as a level-1 network might be too much to ask for because the truth may not be level-1. As an example, evolutionary histories containing many hybridization events between more distantly related species (such as horizontal gene transfer) cannot always be captured well by a level-1 network, since such events often result in complex networks with many nested reticulation events (see, e.g. Soucy et al. 2015, Fig. 5).

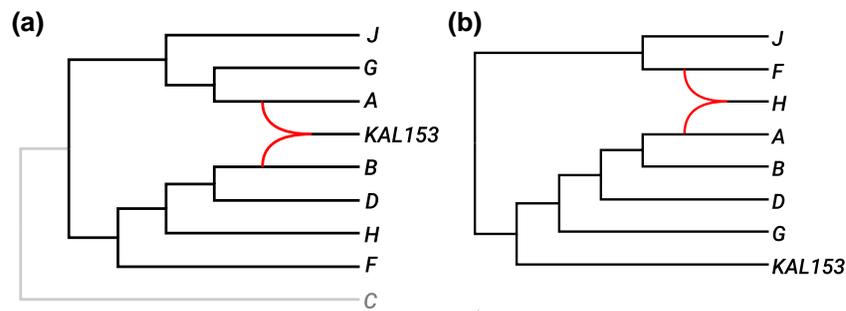
**HIV.** We now consider an MSA of the HIV-1 virus data set containing nine sequences of length 9,953 bp which first appeared in Salemi and Vandamme (2003). This data set is well-studied (Lemey et al. 2009; Huber et al. 2010; Oldman et al. 2016) and contains sequences of the HIV-1 M-group subtypes A, B, C, D, F, G, H, and J as well as a sequence for *KAL153* which is believed to be a recombinant of subtypes A and B (see Lemey et al. 2009, Ch. 16). We use our  $\delta$ -heuristic (formally described in the Materials and Methods) to obtain a weighted set of tf-quarnets from the MSA and then apply SQUIRREL to construct a network, which we root using the outgroup C (as in Salemi and Vandamme 2003; Huber et al. 2010). The  $\delta$ -heuristic and SQUIRREL constructed a clear best scoring network (shown in Fig. 7a) with a weighted tf-quarnet consistency of 0.58 within 1 s.

Indeed, SQUIRREL, combined with the  $\delta$ -heuristic, is able to identify *KAL153* as a recombinant of subtypes A and B, agreeing with the analysis in (Lemey et al. 2009, Ch. 16). This compares favorably to TRILONET (Oldman et al. 2016), where the subtype H was identified as a recombinant (see the constructed network in Fig. 7b). LEV1ATHAN (Huber et al. 2010) was able to identify *KAL153* as a recombinant, but it relies on other algorithms to make the step from sequences to gene trees.

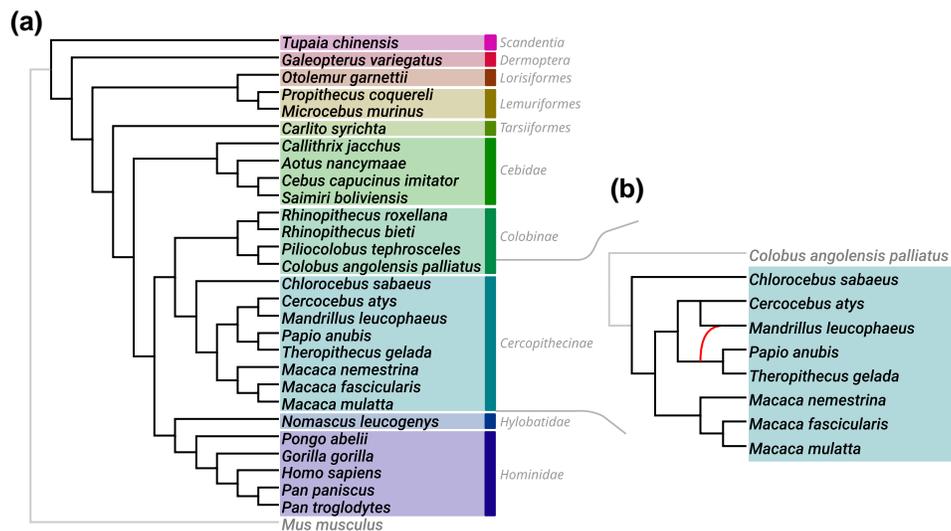
**Primates.** To investigate the performance of SQUIRREL and the  $\delta$ -heuristic on data sets with many taxa and long sequences, we consider an MSA from Vanderpool et al. (2020) of length 1,761,114 bp that contains concatenated sequences for 26 primate species, 2 closely related nonprimate species and the outgroup *Mus musculus*. We first apply the  $\delta$ -heuristic to the MSA to obtain a set of 23,751 weighted tf-quarnets. Subsequently, we use SQUIRREL (specifying *Mus musculus* as the outgroup to root it) and obtain the tree in Fig. 8a after a few minutes on a standard laptop. The tree coincides exactly with the species tree obtained in Vanderpool et al. (2020) using the gene tree-based algorithm ASTRAL III (Zhang et al. 2018b), while largely agreeing with two previously inferred phylogenies (Perelman et al. 2011; Springer et al. 2012). The weighted tf-quarnet consistency score of the tree is 0.995, but some of the other generated candidate networks (which contain reticulations) have scores within 0.003 from this best value, suggesting that reticulate events might have occurred.

We investigate this further by looking only at the eight primates in the *Cercopithecinae* subfamily, for which Vanderpool et al. (2020) have demonstrated possible reticulate events. Combining the  $\delta$ -heuristic and SQUIRREL we generated a set of candidate networks for these eight species and the outgroup *Colobus angolensis palliatus*. Two of the networks had a much higher score than the others and they only differed from each other by the addition of a reticulation edge. In particular, the second best scoring network (shown in Fig. 8b) had a score of 0.956, while the best scoring network was the subtree of the original network with score 0.974 (also shown in Fig. 8b, by ignoring the curved reticulation edge). The *blobtree* of the network (obtained by contracting the cycle into a single node) exactly matches one of the *blobtrees* inferred with TINNIK (Allman et al. 2024a). The particular reticulate event we found was not reported in Vanderpool et al. (2020). However, our reticulate event might be more probable since it is between species in the same continent (Africa), while the study by Vanderpool et al. (2020) mentions possible reticulate events between species on different continents (Asia and Africa). Lastly, Vanderpool et al. (2020) found evidence for a “complex pattern of ancient introgression” (p. 14) within the subfamily and state that roughly 40% of the species within the subfamily are known to hybridize (Tung and Barreiro 2017), which suggests that the true nature of the subfamily might not be well-represented by a level-1 network. This is further supported by the fact that the analysis done in Vanderpool et al. (2020) with PHYLONET (Than et al. 2008; Yu and Nakhleh 2015) and SNAQ (Solís-Lemus and Ané 2016; Solís-Lemus et al. 2017) also gave ambiguous results, while PHYNEST (Kong et al. 2024) yet again concludes with a different network.

The *Cercopithecinae* subfamily (again with outgroup *C. angolensis palliatus*) also featured in Barton et al. (2022) in the context of using the QNR-SVM algorithm for inferring quarnets from a data set. The reason for restricting to a subset was stated as the lack of an algorithm that puzzles together many quarnets. Instead, the authors puzzle them together by hand to obtain a network with a single reticulation that induces 81% of the well-supported quarnets. Using their quarnet weighting scheme, SQUIRREL was able to identify a tree inducing 85% of the well-supported quarnets. (Here, we used a variation of SQUIRREL that takes into account the triangles of the quarnets to choose the best scoring network, instead of the default of just focusing on the tf-quarnets.) Therefore,



**Fig. 7.** a) Phylogenetic network inferred by SQUIRREL (using the  $\delta$ -heuristic to create tf-quarnets) from an MSA of the HIV-1 data set under consideration. The reticulation edges are curved, while the edges leading to the outgroup C are in grey. b) Phylogenetic network inferred by TriLoNet (Oldman et al. 2016) on the same HIV-1 data set (without the outgroup C), again with curved reticulation edges.



**Fig. 8.** a) Phylogenetic tree inferred by SQUIRREL (using the  $\delta$ -heuristic to create tf-quarnets) from an MSA of the primate data set under consideration, with the edges leading to the outgroup *Mus musculus* in grey. The different shaded areas indicate different taxonomical groups as they appear in Vanderpool et al. (2020). The two nonprimate species are *Tupaia chinensis* and *Galeopterus variegatus*. b) In conjunction with the  $\delta$ -heuristic to create tf-quarnets, SQUIRREL inferred two networks with very close weighted tf-quarnet consistency scores from the considered MSA of the subfamily of *Cercopitheciinae* (using *Colobus angolensis palliatus* as outgroup). One of them is the depicted network and the other is the phylogenetic tree obtained from that network by ignoring the curved reticulation edge.

SQUIRREL might be a viable tool to puzzle together quarnets obtained with an algorithm such as QNR-SVM, while still being able to scale to larger data sets unfit for resolving conflicting quarnets by hand.

## Discussion

We have introduced SQUIRREL: a combinatorially consistent algorithm that can puzzle together a dense set of quarnets to create a semi-directed level-1 network. In addition, when combined with the model-based method QNR-SVM (Barton et al. 2022) or the MML algorithm (Martin et al. 2023) for inferring quarnets, SQUIRREL provides a method to create a level-1 network directly from sequence data. To the best of our knowledge, SQUIRREL is one of the first methods that allows the construction of semi-directed level-1 networks from biological data using collections of quarnets. The only other approaches we are aware of that use quarnet information are NANUQ (Allman et al. 2019) and the recently presented NANUQ<sup>+</sup> (Allman et al. 2024b). Although NANUQ<sup>+</sup> uses a

similar distance-based strategy to SQUIRREL to expand the cycles in a network, both NANUQ and NANUQ<sup>+</sup> take as input a collection of gene trees, rather than a dense set of quarnets or a sequence alignment.

Any method that creates a dense set of quarnets from biological data could be used as input for SQUIRREL. In particular, if such a method is statistically consistent under some model (possibly incorporating, e.g. incomplete lineage sorting), the combinatorial consistency of SQUIRREL ensures that the combined inference is consistent as well. Furthermore, SQUIRREL could in principle be combined with methods that may not scale well to larger taxa sets but are still able to construct partial semi-directed level-1 networks (containing some but not all of the studied taxa) from biological data. Indeed, as with supertree methods, partial networks on larger sets of taxa could be converted to quarnets for SQUIRREL by restricting those partial networks to four taxa. This would require a rule to decide what to do in case partial networks overlap on more than four taxa and they induce conflicting quarnets. Hence, a possible direction for future research would be

adapting SQUIRREL to work with nondense sets of quarnets which could contain any number of quarnets for each subset of four taxa.

Using the  $\delta$ -heuristic, SQUIRREL is able to quickly construct a level-1 network directly from sequence data. Our sequence simulations show that the  $\delta$ -heuristic is likely not statistically consistent under the tested K2P model. In particular, an increase in sequence length beyond 10 kb does not give a visible improvement under our simulation settings, which one would expect for a statistically consistent quarnet inference method. Despite the lack of a statistical basis of the  $\delta$ -heuristic, it already shows promising similarity scores for MSAs with a length of 1 kb when combined with SQUIRREL. Furthermore, a major advantage is its speed. As an example, this approach was able to construct a network with 29 taxa from an MSA of length 1.7 Mb within a few minutes on a standard laptop (see the Results section). Hence, we do not see the  $\delta$ -heuristic (combined with SQUIRREL) as an alternative for known model-based methods, but rather as a complementary tool. For one, this approach can be used to generate reasonable starting networks for the time-intensive search through the network space of likelihood-based methods [such as PHYLONET {[Than et al. 2008](#); [Yu and Nakhleh 2015](#)}, SNAQ {[Solís-Lemus and Ané 2016](#); [Solís-Lemus et al. 2017](#)}, and PHYNEST {[Kong et al. 2024](#)}]. On the other hand, it can be used to quickly gain insight into sequence data without the need to first infer gene trees with a different tool, as is the case for NANUQ ([Allman et al. 2019](#)), which requires many accurate gene trees to make a good estimate of the concordance factors.

With the increasing availability of genome and transcriptome data, biologists are also likely to explore the reconstruction of separate phylogenetic networks for multiple sets of short orthologous sequences. Rapid construction of such networks for the same set of taxa across different sets of orthologues opens up the possibility for comparative analyses. A possible research direction in this area would be to combine SQUIRREL's speed for constructing semi-directed level-1 networks with the tf-quarnet consistency score or the recently introduced dissimilarity measure for semi-directed networks that generalizes the widely-used Robinson–Foulds distance for phylogenetic trees ([Maxfield et al. 2025](#)), which would permit the rapid comparison of networks computed for different sets of orthologues. It also leads to the interesting problem of finding a consensus of a collection of semi-directed networks, which to our best knowledge has not yet been addressed in the literature. One approach to this problem could be to treat it as a supernetwork question where all input networks have the same leaf set, and use the approach suggested earlier in this section.

Our simulations indicate that SQUIRREL can construct networks closely resembling an underlying network in terms of tf-quarnets, even if many of the tf-quarnets are wrongly inferred. In particular, both of the considered consistency scores average above 0.91 even for sets containing only 50% of the original tf-quarnets. This is a significant improvement compared to a similar experiment to the triplet/trinet-based LEV1ATHAN and TRILONET algorithms, where the *trinet consistency score* (the rooted three-leaf network analogue of our tf-quarnet consistency score) dropped below 0.5 for sets still containing 75% of the trinetts ([Oldman et al. 2016](#)). These results can be considered as evidence that SQUIRREL is able to

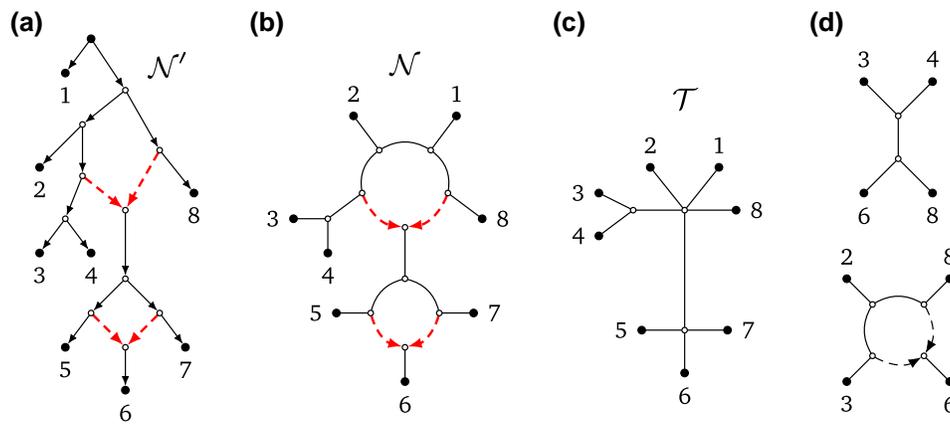
construct networks with a high topological resemblance to the original network in terms of tf-quarnets, even for a high percentage of incorrect tf-quarnets. As mentioned in the Result section, even though the tf-quarnets are theoretically enough to construct a triangle-free semi-directed level-1 network, in practice, contracting the triangles might somewhat weaken the signal of reticulation events. Note that theoretically (that is, when all quarnets come from a single network with  $n$  leaves) only  $\mathcal{O}(n \log n)$  tf-quarnets are required to reconstruct the network, instead of the full set of  $\mathcal{O}(n^4)$  tf-quarnets ([Frohn et al. 2025](#)). Thus, even sets with many incorrect tf-quarnets might still hold enough information to reconstruct the original network. This could also explain why a higher number of leaves seems to have a positive effect on the similarity score:  $\mathcal{O}(n \log n)$  grows slower than  $\mathcal{O}(n^4)$ , so the fraction of tf-quarnets necessary to reconstruct a network decreases when  $n$  grows.

Although several methods can construct semi-directed level-1 networks, the assumption that a network is level-1 might be too restrictive in many cases for biological data. A major breakthrough would be to develop a practical algorithm that is able to construct networks that are more complex than level-1 networks. Some theoretical results have already appeared towards tackling this problem. For example, it is known that semi-directed level-2 networks are uniquely encoded by the quarnets they induce ([Huber et al. 2024](#)). In addition, under several models, the circular ordering around the blobs of outerlabeled planar networks (a class of semi-directed networks more general than semi-directed level-1 networks) is also shown to be identifiable ([Rhodes et al. 2025](#)). Furthermore, the recently introduced TINNiK algorithm ([Allman et al. 2024a](#)) uses concordance factors computed from gene trees to construct the blobtree of networks with arbitrary level under the network multispecies coalescent model. Although such a blobtree still remains a tree, it does indicate in what areas of the underlying network reticulations may have occurred. It might also be worth looking for an extension of SQUIRREL to nonbinary networks, where high-degree vertices are allowed which do not necessarily represent reticulate events.

In conclusion, SQUIRREL provides an efficient and combinatorially sound approach for reconstructing semi-directed level-1 networks from dense sets of quarnets. The promising consistency scores achieved in our tests underscore SQUIRREL's ability to retain network topology even when faced with noisy data. Together with our  $\delta$ -heuristic, SQUIRREL allows rapid insight into large-scale sequence data. Looking forward, we hope that this approach can complement more time-intensive methods and support the preliminary exploration of network hypotheses.

## Materials and Methods

We start this section by presenting formal definitions surrounding phylogenetic networks and quarnets in the first subsection. The high-level idea of SQUIRREL is described in the second subsection, while its subroutines are formalized in the third and fourth subsection. We end with the description of the  $\delta$ -heuristic in the fifth subsection, and a brief description of the consistency and implementation of SQUIRREL in the sixth and seventh subsection, respectively.



**Fig. 9.** a) A rooted phylogenetic level-1 network  $\mathcal{N}'$  on leaf set  $\mathcal{X} = \{1, \dots, 8\}$ , with the dashed reticulation edges pointing towards its reticulation vertices. b) The triangle-free semi-directed level-1 network  $\mathcal{N}$  which can be obtained from  $\mathcal{N}'$  by suppressing its root and keeping only the dashed reticulation edges directed. c) The blobtree  $\mathcal{T}$  of the semi-directed network  $\mathcal{N}$ , obtained by collapsing all cycles into single vertices. d) Two of the tf-quarnets induced by  $\mathcal{N}$ . When ignoring the leaf labels, these are the two possible tf-quarnet shapes. The top tf-quarnet is a quartet tree and the bottom tf-quarnet is a 4-cycle.

## Phylogenetic Networks and Quarnets

**Phylogenetic networks.** A rooted phylogenetic network on a set of at least four leaves  $\mathcal{X}$  (representing a set of taxa) is a directed acyclic graph with a single root, no parallel edges and no directed cycles such that (i) the root has two children; (ii) each leaf (i.e. a vertex with no children) has one parent and is uniquely labeled by an element from  $\mathcal{X}$ ; (iii) all other vertices either have one parent and two children, or two parents and one child. A vertex of the latter type is a *reticulation (vertex)*, and the two edges directed towards it are *reticulation edges*. See Fig. 9a for an example. *Semi-directed phylogenetic networks*, the type of network this paper is concerned with, can be obtained from a rooted phylogenetic network by suppressing its root and undirecting all edges except for the reticulation edges. For the sake of brevity, we refer to these networks simply as *semi-directed networks*. Since the reticulation edges remain directed, we can still refer to the reticulation vertices and edges of a semi-directed network (see Fig. 9b). We call a semi-directed network *triangle-free* if it does not contain any triangles (3-cycles). Note that a semi-directed network without any reticulations is an (unrooted) phylogenetic tree in the usual sense.

In this paper, we consider semi-directed networks which are *level-1* (again see Fig. 9b), meaning that every reticulation is part of exactly one undirected cycle (ignoring the directions of the reticulation edges). The (possibly nonbinary) phylogenetic tree obtained by collapsing every such cycle into a single vertex is called the *blobtree* (or *tree of blobs*) of the semi-directed network (see Fig. 9c).

Given a semi-directed network  $\mathcal{N}$  on  $\mathcal{X}$ , a partition  $A | B$  of  $\mathcal{X}$  (with  $A$  and  $B$  both nonempty) is a *split* of  $\mathcal{N}$  if there exists an edge of  $\mathcal{N}$  whose removal disconnects the leaves in  $A$  from those in  $B$ . Such a split is *nontrivial* if the corresponding partition is nontrivial, that is, if  $|A|, |B| \geq 2$ . As an example,  $\{1, 2, 3, 4, 8\} | \{5, 6, 7\}$  is a nontrivial split of the network from Fig. 9b. We sometimes omit the set notation for splits with few elements, meaning that we write  $ab | cd$  instead of the split  $\{a, b\} | \{c, d\}$  of the set  $\{a, b, c, d\}$ .

**Quarnets.** A semi-directed network  $q$  on a set of four leaves  $\mathcal{L}(q) = \{a, b, c, d\}$  is called a (*semi-directed*) *quarnet*. Recall that up to relabeling the leaves, there are six different level-1 quarnets (see Fig. 1c). Here, we mostly focus on *tf-quarnets*:

triangle-free level-1 quarnets. For a given leaf set  $\mathcal{X} = \{a, b, c, d\}$  and up to relabeling of the leaves, there are only two such tf-quarnets on  $\mathcal{X}$ : the *quartet tree* and the *4-cycle* (see Fig. 9d). We often denote a quartet tree by its induced split (e.g.  $ab | cd$ ), while we describe a 4-cycle by its circular ordering [e.g.  $(a, b, c, d)$ ] and mention the leaf below the reticulation separately. Note that tf-quarnets either have no nontrivial split at all, or they have exactly one nontrivial split [e.g. for  $\mathcal{X} = \{a, b, c, d\}$  the splits  $ab | cd$ ,  $ac | bd$ , or  $ad | bc$ ].

## SQUIRREL: Main Algorithm

SQUIRREL uses as input a set  $\mathcal{Q}$  of tf-quarnets on some leaf set  $\mathcal{X}$  with  $n = |\mathcal{X}| \geq 4$ . In particular, this set needs to be *dense*, meaning that it contains exactly one tf-quarnet for each subset of four leaves of  $\mathcal{X}$  (see also the Introduction). Such a set can be created from a MSA using QNR-SVM (Barton et al. 2022), the MML algorithm (Martin et al. 2023) or our own  $\delta$ -heuristic (see the fifth subsection of this section). We also allow for a function  $w: \mathcal{Q} \rightarrow [0, 1]$  to give weights to the tf-quarnets, which can e.g. be used to model confidence or bootstrap support. Unweighted tf-quarnets are assumed to have unit weights.

The main idea behind the SQUIRREL algorithm is to first build a sequence of  $n - 3$  phylogenetic trees on the given  $n$  leaves, each one less refined than the other (see Algorithm 2). These trees will function as candidate blobtrees. By expanding all the high-degree nodes in these trees into cycles (and introducing reticulations), we obtain a set of semi-directed candidate networks (see Algorithm 3). Finally, out of these networks, we choose the network  $\mathcal{N}$  with the highest *weighted tf-quarnet consistency score*, defined as

$$C(\mathcal{Q}, \mathcal{N}) = \frac{w(\mathcal{Q} \cap \mathcal{Q}(\mathcal{N}))}{w(\mathcal{Q})}. \quad (3)$$

Here,  $\mathcal{Q}$  is the input set of tf-quarnets and  $\mathcal{Q}(\mathcal{N})$  is the set of tf-quarnets which are induced by the output network  $\mathcal{N}$ . A tf-quarnet  $q$  is *induced* by the network  $\mathcal{N}$  if it is the restriction of  $\mathcal{N}$  to  $\mathcal{L}(q)$ , which is formally defined as the network obtained from  $\mathcal{N}$  by deleting all leaves not in  $\mathcal{L}(q)$  and exhaustively applying the following operations: deleting unlabeled leaves, deleting degree-2 reticulations, suppressing nonreticulate degree-2 vertices, suppressing parallel edges, and

suppressing triangles. For completeness, we mention that an induced quartet can be defined similarly but without suppressing the triangles.

The pseudo-code of SQUIRREL is shown as [Algorithm 1](#). The blobtree construction algorithm ([Algorithm 2](#)) and the cycle expansion algorithm ([Algorithm 3](#)) are explained in detail in the following two subsections. Even though this is not specified in the pseudo-code, SQUIRREL does allow the user to specify an outgroup as input. Then, it makes sure that all candidate networks can be rooted using this outgroup (see also the fourth subsection of this section).

#### Algorithm 1 SQUIRREL

---

**Input:** dense set  $\mathcal{Q}$  of weighted tf-quarnets on  $\mathcal{X} = \{x_1, \dots, x_n\}$   
**Output:** triangle-free semi-directed level-1 network on  $\mathcal{X}$   
1  $(\mathcal{T}_1, \dots, \mathcal{T}_{n-3}) \leftarrow$  candidate blobtrees, using [Algorithm 2](#)  
2  $(\mathcal{N}_1, \dots, \mathcal{N}_{n-3}) \leftarrow$  semi-directed candidate networks obtained from the candidate blobtrees  $\mathcal{T}_i$ , using [Algorithm 3](#)  
3 **return** network  $\mathcal{N}_i$  with highest weighted tf-quarnet consistency score

---

### SQUIRREL: Constructing Candidate Blobtrees

In the following three steps, we describe how SQUIRREL creates the sequence of candidate blobtrees on leaf set  $\mathcal{X}$  from the dense set  $\mathcal{Q}$  of tf-quarnets. The pseudo-code of this procedure is shown as [Algorithm 2](#) at the end of this subsection.

**Step A1.** We first create a phylogenetic tree  $\mathcal{T}^*$  on  $\mathcal{X}$  as described in [Berry and Gascuel \(2000\)](#). Their algorithm takes as input a (possibly nondense) set  $\mathcal{Q}'$  of quartet trees and returns as  $\mathcal{T}^*$  the unique most refined phylogenetic tree on  $\mathcal{X}$  which does not induce a quartet with a different nontrivial split than one of the quartets in  $\mathcal{Q}'$  (see [supplementary material Section A, Supplementary Material](#) online for a more formal definition). By taking  $\mathcal{Q}'$  to be the subset of quartet trees in our set of dense tf-quarnets  $\mathcal{Q}$  (see line 1 of [Algorithm 2](#)), we can employ the algorithm from [Berry and Gascuel \(2000\)](#) to obtain  $\mathcal{T}^*$  (see line 2 of [Algorithm 2](#)). As we show in [supplementary material Lemma A.2, Supplementary Material](#) online, in the case all tf-quarnets are induced by a unique network,  $\mathcal{T}^*$  coincides with the blobtree of that network.

**Step A2.** Since the set  $\mathcal{Q}$  (and thus  $\mathcal{Q}'$ ) is constructed from real data, we expect there to be a fair amount of quartets that contradict each other. Hence, in practice, the tree  $\mathcal{T}^*$  constructed in Step A1 will be highly unresolved. To remedy this problem, we use a method to refine the tree  $\mathcal{T}^*$ , specifically, an adapted version of the QUARTETJOINING algorithm ([Grünwald et al. 2009](#)). QUARTETJOINING takes as input a function  $\omega$  that assigns a nonnegative real number to each possible nontrivial split of four leaves in  $\mathcal{X}$ . Starting with the star-tree with central vertex  $v$  and leaf set  $\mathcal{X}$ , QUARTETJOINING sequentially introduces edges between  $v$  and two of its neighbors (according to some criterion involving the function  $\omega$ ) until the tree is fully resolved.

In our case, we instead start with the tree  $\mathcal{T}^*$  (which might already be partially resolved) and adapt QUARTETJOINING to resolve  $\mathcal{T}^*$  further. This eventually leads to a fully resolved phylogenetic tree  $\mathcal{T}_1$  on  $\mathcal{X}$ , which functions as the first tree in our sequence of candidate blobtrees (see line 3 of [Algorithm 2](#)). In our adaptation, instead of considering all

combinations of neighbors of the central vertex  $v$ , we consider all such combinations of neighbors of any of the internal (i.e. nonleaf) vertices with degree at least 4. We construct the function  $\omega$  used as input to QUARTETJOINING as follows. For any tf-quarnet  $q \in \mathcal{Q}$  with leaf set  $\mathcal{L}(q) = \{a, b, c, d\}$  such that  $q$  is a quartet tree (say with split  $ab|cd$ ), we set  $\omega(ab|cd) = w(q)$  for the input weight function  $w$  mentioned at the beginning of the previous subsection. All other nontrivial splits of four leaves of  $\mathcal{X}$  are assigned an  $\omega$ -value of 0.

**Step A3.** Finally, we explain how we create the full sequence of candidate blobtrees from the phylogenetic tree  $\mathcal{T}_1$ . Given an edge  $uv$  of the tree  $\mathcal{T}_1$  that induces a nontrivial split  $A|B$ , we collect all the quartet trees in  $\mathcal{Q}$  for which their induced splits restrict to quartet splits of  $A|B$  in a set  $\mathcal{Q}'(A|B)$  by first defining  $\mathcal{Q}(A|B) = \{q \in \mathcal{Q} : |A \cap \mathcal{L}(q)| = 2, |B \cap \mathcal{L}(q)| = 2\}$  and then  $\mathcal{Q}'(A|B) = \{q \in \mathcal{Q}(A|B) : q \text{ has split } A \cap \mathcal{L}(q) | B \cap \mathcal{L}(q)\}$ . This allows us to define the *split-support* of  $uv$  as

$$\text{supp}(uv) = \frac{w(\mathcal{Q}'(A|B))}{w(\mathcal{Q}(A|B))}, \quad (4)$$

i.e. as the weighted ratio of the tf-quarnets in  $\mathcal{Q}$  that support the split induced by the edge  $uv$ . For each of the  $n - 3$  edges of the tree  $\mathcal{T}_1$  we then compute this split-support (see line 4 of [Algorithm 2](#)). Afterwards, we sort the edges of  $\mathcal{T}_1$  in increasing order, according to their split-support. To create the trees  $(\mathcal{T}_2, \dots, \mathcal{T}_{n-3})$ , we keep contracting the least supported edge (see line 6 of [Algorithm 2](#)). In other words, the tree  $\mathcal{T}_i$  is obtained from  $\mathcal{T}_1$  by contracting the  $i - 1$  least supported edges. Crucial for our consistency proof in [supplementary material Section A, Supplementary Material](#) online is that  $\mathcal{T}_1$  is a refinement of  $\mathcal{T}^*$ , and therefore one of the trees in the sequence  $(\mathcal{T}_1, \dots, \mathcal{T}_{n-3})$  will be the tree  $\mathcal{T}^*$ .

#### Algorithm 2 Constructing candidate blobtrees

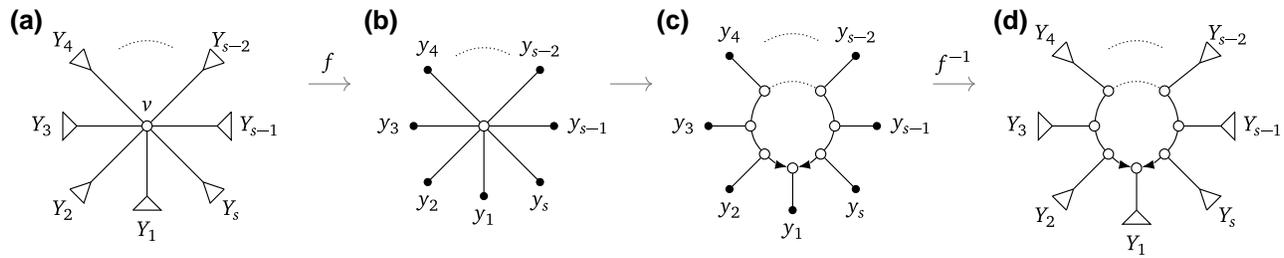
---

**Input:** dense set  $\mathcal{Q}$  of weighted tf-quarnets on  $\mathcal{X} = \{x_1, \dots, x_n\}$   
**Output:** sequence of candidate blobtrees  $(\mathcal{T}_1, \dots, \mathcal{T}_{n-3})$  on  $\mathcal{X}$   
/\* Step A1 \*/  
1  $\mathcal{Q}' \leftarrow$  set of all quartet trees in  $\mathcal{Q}$   
2  $\mathcal{T}^* \leftarrow$  phylogenetic tree on  $\mathcal{X}$  obtained from  $\mathcal{Q}'$ , as described in [Berry and Gascuel \(2000\)](#)  
/\* Step A2 \*/  
3  $\mathcal{T}_1 \leftarrow$  phylogenetic tree on  $\mathcal{X}$  obtained by applying the adapted QUARTETJOINING algorithm to  $\mathcal{T}^*$  and  $\mathcal{Q}$   
/\* Step A3 \*/  
4 **compute** the split-support for every edge in  $\mathcal{T}_1$   
5 **for**  $i \in \{2, \dots, n - 3\}$  **do**  
6      $\mathcal{T}_i$  is constructed from  $\mathcal{T}_{i-1}$  by contracting the least supported edge  
7 **end**  
8 **return**  $(\mathcal{T}_1, \dots, \mathcal{T}_{n-3})$

---

### SQUIRREL: Expanding Cycles in a Tree

Once SQUIRREL has constructed the sequence of candidate blobtrees using [Algorithm 2](#), we transform them into triangle-free semi-directed level-1 networks using the dense set of tf-quarnets  $\mathcal{Q}$ . In this subsection, we describe how we transform a phylogenetic tree  $\mathcal{T}$ —representing one of our candidate blobtrees—into such a network  $\mathcal{N}$ . In particular, we replace every internal vertex of the given tree by a suitable cycle. Since our aim is to build triangle-free networks, we



**Fig. 10.** a) A blobtree on some leaf set  $\mathcal{X}$  with an internal vertex  $v$  inducing the partition  $Y_1 | \dots | Y_s$  of  $\mathcal{X}$ . b) Illustration of the mapping  $f$  which maps every leaf  $x$  of  $\mathcal{X}$  to a leaf in  $\{y_1, \dots, y_s\}$ , depending on which set  $Y_i$  contains  $x$ . c) Illustration of Step B2 and B3 of SQUIRREL, where the single internal vertex is replaced by a cycle. d) Illustration of how the cycle on the leaves  $y_i$  is mapped back to a cycle on the sets  $Y_i$  with the inverse function  $f^{-1}$ .

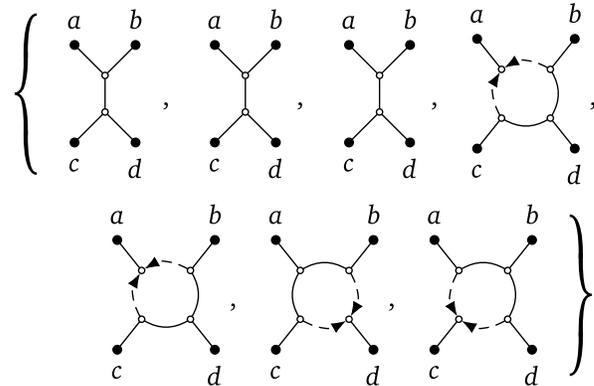
replace vertices incident to  $s \geq 4$  edges by an  $s$ -cycle with a reticulation (see also the illustration in Fig. 10). To this end, we repeat the following three steps for every such internal vertex  $v$  (starting with the ones with the highest degree). The corresponding high-level pseudo-code is shown as Algorithm 3.

**Step B1.** The first step in our approach is to assign a dense set of representative *tf-quarnets*  $\tilde{Q}_v$  to each internal vertex  $v$  of  $\mathcal{T}$  with degree  $s \geq 4$ . In particular, the set  $\tilde{Q}_v$  will be a dense set of *tf-quarnets* on the leaf set  $\mathcal{Y} = \{y_1, \dots, y_s\}$ , where each  $y_i$  represents the set  $Y_i$  which is part of the partition  $Y_1 | \dots | Y_s$  of  $\mathcal{X}$  induced by  $v$  (see Fig. 10a and b). In the next step, these sets will then be used to determine by what cycle to replace  $v$ .

First, let  $f: \mathcal{X} \rightarrow \mathcal{Y}$  be the function that maps every leaf  $x \in \mathcal{X}$ , with  $x$  being in some set  $Y_i$ , to the leaf  $y_i$  (see line 3 of Algorithm 3). To construct the *tf-quarnets* in  $\tilde{Q}_v$  (see line 4 of Algorithm 3), we repeat the following procedure for every subset  $\{y_i, y_j, y_k, y_l\}$  of four leaves in  $\mathcal{Y}$ . Let  $\mathcal{Q}_{\{i,j,k,l\}} = \{q \in \mathcal{Q} : \mathcal{L}(q) = \{x_i, x_j, x_k, x_l\} \text{ with } x_p \in Y_p \text{ for all } p \in \{i, j, k, l\}\}$  be the subset of  $\mathcal{Q}$  containing only *tf-quarnets* with one leaf in each of the four sets  $Y_i, Y_j, Y_k$  and  $Y_l$ . By relabeling the leaves of all *tf-quarnets* in  $\mathcal{Q}_{\{i,j,k,l\}}$  with the function  $f$ , we obtain a multiset of *tf-quarnets* which all have the same leaf set  $\{y_i, y_j, y_k, y_l\}$ . With slight abuse of notation, we denote this multiset by  $f(\mathcal{Q}_{\{i,j,k,l\}})$ . Then, we choose one of the *tf-quarnets* in the multiset  $f(\mathcal{Q}_{\{i,j,k,l\}})$  to assign to  $\tilde{Q}_v$  as the *tf-quarnet* on the four-leaf set  $\{y_i, y_j, y_k, y_l\}$  (see next paragraph). As mentioned before, this is repeated for every subset  $\{y_i, y_j, y_k, y_l\}$  of four leaves in  $\mathcal{Y}$ , resulting in a dense set of *tf-quarnets* on  $\mathcal{Y}$ .

To choose a *tf-quarnet* from the multiset  $f(\mathcal{Q}_{\{i,j,k,l\}})$ , we first choose its *skeleton*: its underlying undirected graph. In particular, for each of the six possible skeletons  $t$  (three quartet trees and three undirected 4-cycles) we let  $w(t)$  be the sum of weights of all *tf-quarnets* in  $f(\mathcal{Q}_{\{i,j,k,l\}})$  with the given skeleton  $t$ . We then choose the skeleton  $t$  with the highest weight (with ties resolved randomly) and assign it a new weight of  $w(t)/w(f(\mathcal{Q}_{\{i,j,k,l\}}))$ . Note that in the unweighted case this simply means that we choose the skeleton that appears most in the multiset. We first choose the skeleton since determining the location of the reticulation in a *quarnet* from data seems especially hard (Martin et al. 2023). If our chosen skeleton is one of the quartet trees, we assign that as our *tf-quarnet* on  $\{y_i, y_j, y_k, y_l\}$ . On the other hand, if one of the undirected 4-cycles appears most, we still need to determine the location of the reticulation. This is done by checking which leaf appears most often below the reticulation in all 4-cycles with the chosen skeleton.

As an example of this voting procedure to choose a *tf-quarnet* from the multiset  $f(\mathcal{Q}_{\{i,j,k,l\}})$ , suppose our multiset  $f(\mathcal{Q}_{\{i,j,k,l\}})$  contains only *tf-quarnets* with weight 1 and is as in Fig. 11.



**Fig. 11.** A multiset of 7 *tf-quarnets* on leaf set  $\{a, b, c, d\}$ .

Then, we choose the 4-cycle with circular ordering  $(a, b, c, d)$  as our skeleton, after which we assign  $a$  to be the leaf below the reticulation. The new *tf-quarnet* is then a 4-cycle with a weight of  $4/7$  because that 4-cycle appears 4 times out of a total of 7 *tf-quarnets*.

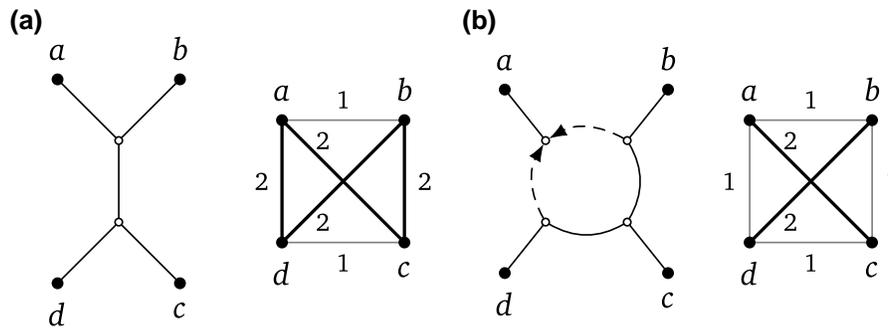
**Step B2.** The next step of our approach is to determine a circular ordering of the leaves in the set  $\mathcal{Y}$  based on the *tf-quarnets* in  $\tilde{Q}_v$ . Note that we repeat this for every internal vertex  $v$  of  $\mathcal{T}$  with degree at least 4. First, we use the set  $\tilde{Q}_v$  to create a distance  $D_{\tilde{Q}_v}$  between every pair of leaves in  $\mathcal{Y}$  (see line 5 of Algorithm 3). Formally, given two leaves  $a$  and  $b$  in  $\mathcal{Y}$ , we define the distance  $D_{\tilde{Q}_v}$  as follows:

$$D_{\tilde{Q}_v}(a, b) = \begin{cases} 0 & \text{if } a = b, \\ \sum_{q \in \tilde{Q}_v : a, b \in \mathcal{L}(q)} \tau_q(a, b) & \text{if } a \neq b. \end{cases} \quad (5)$$

For every *tf-quarnet*  $q \in \tilde{Q}_v$  the exact value of  $\tau$  depends on the weight of  $q$  and the position of the leaves  $a$  and  $b$  within it. In particular, the values are defined on the skeleton of the *tf-quarnets* and hence do not depend on the position of the reticulations. Given two leaves  $a$  and  $b$  of a *tf-quarnet*  $q$  (with weight  $w(q) \in [0, 1]$ ), we define  $\tau_q$  as follows:

$$\tau_q(a, b) = \begin{cases} (3 - w(q))/2 & \text{if } q \text{ is the quartet tree } ab | cd \\ & \text{or if } q \text{ is a 4-cycle with } a, b \\ & \text{as neighbors,} \\ (3 + w(q))/2 & \text{otherwise.} \end{cases} \quad (6)$$

Here, we say that two leaves of a 4-cycle are *neighbors* if they are not on opposite sides of the cycle. The  $\tau_q$ -values reduce to 1



**Fig. 12.** Two tf-quarnets  $q$  with leaf set  $\{a, b, c, d\}$ : a quartet tree a) and 4-cycle b). The values  $\tau_q$  [as defined by equation (6), assuming the quarnets have weight 1] between any two leaves are illustrated by the two complete graphs, where the thin grey edges have length 1 and the thick black length 2.

or 2 for tf-quarnets  $q$  with a weight of 1. Specifically, two leaves on the same side of a split in a quartet tree  $q$  have a  $\tau_q$ -value of 1, otherwise they have a  $\tau_q$ -value of 2. Similarly, two neighboring leaves in a 4-cycle  $q$  have a  $\tau_q$ -value of 1, while two opposite leaves have a  $\tau_q$ -value of 2. See Fig. 12 for an illustration of these values. Note that these pairwise distances between leaves resemble the *quartet distances* used in NANUQ (Allman et al. 2019) and NANUQ<sup>+</sup> (Allman et al. 2024b).

Once the distances  $D_{\tilde{Q}_v}$  are computed, we create a complete graph  $G$  with vertex set  $\mathcal{Y}$ , where the distances between the vertices are given by  $D_{\tilde{Q}_v}$ . By solving the TSP on this graph, we obtain a circular ordering of the elements in  $\mathcal{Y}$  (see line 6 of Algorithm 3). The goal of a TSP instance is to find a shortest *Hamiltonian cycle* (or *TSP-tour*): a cycle that visits each vertex exactly once. The default setting for SQUIRREL is to use the Held–Karp algorithm (Bellman 1962; Held and Karp 1962) for up to and including 13 leaves and to use simulated annealing to heuristically solve instances with more leaves. To obtain true consistency (see supplementary material Section A, Supplementary Material online) this setting can be changed to always solve TSP to optimality, at the cost of a longer running time.

**Step B3.** After solving TSP, SQUIRREL obtains a circular ordering  $\theta$  of the leaves in  $\mathcal{Y}$ . It remains to determine which leaf  $y_i$  needs to be the leaf below the reticulation in the resulting cycle. To ensure SQUIRREL always returns a valid (that is, rootable) semi-directed network, we create a *reticulation ranking*  $\rho$  of the leaves in  $\mathcal{Y}$  instead of picking a single leaf (see line 7 of Algorithm 3). If the set  $\mathcal{Y}$  contains at least five elements, we order them according to how often they appear in a 4-cycle of  $\tilde{Q}_v$  (as defined in Step B1). That is, the first leaf in our ranking  $\rho$  appears most often in a 4-cycle and is our first option to be the leaf below the reticulation. The case where  $|\mathcal{Y}| = 4$  is special, since  $\tilde{Q}_v$  then only contains a single quarnet. If this is a 4-cycle, then the leaf below the reticulation of that 4-cycle is the first leaf in our ranking  $\rho$ . The other three leaves (or in the case that the single tf-quarnet is a quartet tree: all four leaves) are ordered randomly.

Finally, we map every leaf  $y_i$  back to the corresponding leaf set  $Y_i$  of the original tree  $\mathcal{T}$  with the inverse function  $f^{-1}$ . While slightly abusing notation, this results in an ordering  $f^{-1}(\theta)$  of the sets  $Y_i$ . Then, we replace the internal vertex  $v$  in the tree  $\mathcal{T}$  by a cycle that follows this ordering  $f^{-1}(\theta)$  (see line 9 of Algorithm 3 and Fig. 10b and c for an illustration).

We determine the location of the reticulation by looking at the first element  $\rho_1$  of the reticulation ranking  $\rho$ . In particular, we let the leaf set in  $f^{-1}(\rho_1)$  be below the reticulation (again see line 9 of Algorithm 3). This could possibly create a partially constructed network that is *invalid*: one without a valid root location (e.g. if two reticulations are oriented towards each other). Hence, if this is the case we instead pick the leaves in  $f^{-1}(\rho_2)$ . If this is still an invalid option, we keep iterating through the ranking  $\rho$  until we find a valid partial network (see line 10 of Algorithm 3). Note that This procedure ensures that we always return a valid semi-directed network at the end of Algorithm 3. Our implementation of SQUIRREL also allows the user to specify a known outgroup. Then, a (partially constructed) semi-directed network is only valid if it is not only rootable, but if it can also be rooted at the edge incident to the outgroup. Iterating through the reticulation ranking ensures that we always return a valid semi-directed network at the end of Algorithm 3, even in the case of a specified outgroup (see supplementary material Lemma A.6 in Section A, Supplementary Material online for a proof).

### Algorithm 3 Expanding cycles in a tree

---

**Input:** dense set  $\mathcal{Q}$  of weighted tf-quarnets on  $\mathcal{X} = \{x_1, \dots, x_n\}$ , phylogenetic tree  $\mathcal{T}$  on  $\mathcal{X}$

**Output:** triangle-free semi-directed level-1 network on  $\mathcal{X}$

```

1 for internal vertex  $v$  of  $\mathcal{T}$  with degree  $\geq 4$  do // in decreasing order of degree
  /* Step B1 */
2   $Y_1 \dots Y_s \leftarrow$  partition of  $\mathcal{X}$  induced by  $v$ 
3   $f \leftarrow$  function that maps a leaf  $x \in \mathcal{X}$  to a leaf  $y_i$ , depending on the set  $Y_i$  that contains  $x$ 
4   $\tilde{Q}_v \leftarrow$  set of representative quarnets of  $v$  on leaf set  $\mathcal{Y} = \{y_1, \dots, y_s\}$ 
  /* Step B2 */
5  compute the distances  $D_{\tilde{Q}_v}(y_i, y_j)$  for all  $i, j \in \{1, \dots, s\}$ 
6   $\theta \leftarrow$  optimal TSP-tour on  $\{y_1, \dots, y_k\}$  with respect to distances  $D_{\tilde{Q}_v}$ 
  /* Step B3 */
7   $\rho \leftarrow$  reticulation ranking of the leaves in  $\mathcal{Y}$ 
8  for  $j \in \{1, \dots, s\}$  do
9    replace  $v$  in  $\mathcal{T}$  by a cycle  $C$  with ordering  $f^{-1}(\theta)$  and with  $f^{-1}(\rho_j)$  below the reticulation
10   if  $\mathcal{T}$  has a valid root location then
11     break
12   end
13 end
14 end
15 return  $\mathcal{T}$ 

```

---

### $\delta$ -Heuristic: Inferring Quarnets from Sequence Data

As explained before, two model-based methods that use algebraic invariants exist to generate tf-quarnets (Barton et al. 2022; Martin et al. 2023). To allow SQUIRREL to function as a stand-alone tool, we also include a method to infer weighted tf-quarnets from an MSA on a set of taxa  $\mathcal{X}$ : the  $\delta$ -heuristic. Our  $\delta$ -heuristic is based on the concept of  $\delta$ -plots, which function as a measure of treelikeness for sets of four taxa and which were able to pick out recombinants in many simulations (Holland et al. 2002). The algorithm also resembles some aspects of the heuristic to generate trinets from sequences in Oldman et al. (2016). We are now ready to present the steps to create a dense set of weighted tf-quarnets from an MSA on leaf set  $\mathcal{X}$ .

**Step I.** For each pair of taxa  $\{a, b\}$ , we consider the gap-free subalignment of the MSA on  $\{a, b\}$ . That is, we consider only the columns where both taxon  $a$  and  $b$  contain no gaps. Using this subalignment, we assign a distance value  $h_{ab}$  to the pair  $\{a, b\}$ . In particular,  $h_{ab}$  is the *normalized Hamming distance*: the number of columns of the subalignment where taxon  $a$  and  $b$  differ, divided by the total length of the subalignment. Recall that if a tf-quarnet on  $\{a, b, c, d\}$  has a non-trivial split, it has one of the three splits  $ab|cd$ ,  $ac|bd$  or  $ad|bc$ . For each four-taxon subset and for each of these three splits, say  $ab|cd$ , we then let  $h_{ab|cd} = h_{ab} + h_{cd}$ .

The  $\delta$ -value (introduced in Holland et al. 2002) of such a subset  $\{a, b, c, d\}$  of  $\mathcal{X}$  is now defined as follows (assuming we have that  $h_{ab|cd} \geq h_{ac|bd} \geq h_{ad|bc}$ ):

$$\delta_{\{a,b,c,d\}} = \frac{h_{ab|cd} - h_{ac|bd}}{h_{ab|cd} - h_{ad|bc}}, \quad (7)$$

where  $\delta_{\{a,b,c,d\}} = 0$  if  $h_{ab|cd} = h_{ac|bd} = h_{ad|bc}$ . Intuitively, the  $\delta$ -value indicates how much support there is from the subalignment that the tf-quarnet on  $\{a, b, c, d\}$  has a split. That is, if the value of  $\delta_{\{a,b,c,d\}}$  is close to 1, we expect the split  $ab|cd$  to be present.

**Step II.** With the  $\delta$ -values computed for each subset of four taxa, we partition the 4-taxon sets into two subsets  $S_\lambda$  and  $F_\lambda$  for a predefined threshold value  $\lambda \in (0, 1)$ . The set  $S_\lambda$  will contain all 4-taxon subsets for which the  $\delta$ -value is at least  $\lambda$ , while the set  $F_\lambda$  contains those sets with an  $\delta$ -value smaller than  $\lambda$ . We then expect the sets in  $S_\lambda$  to come from a tf-quarnet with a nontrivial split, while those in  $F_\lambda$  are likely to have come from 4-cycle tf-quarnets. Experiments from Holland et al. (2002) show that an average  $\delta$ -value higher than 0.3 is often enough to determine whether recombination was present (or equivalently, whether a tf-quarnet has a nontrivial split). Hence, we settle for a value of  $\lambda = 0.3$ .

**Step III.** Every 4-taxon set  $\{a, b, c, d\}$  in  $S_\lambda$  is assigned a quartet tree. Its split is simply determined by the split  $s \in \{ab|cd, ac|bd, ad|bc\}$  for which  $h_s$  is the highest. On the other hand, the sets in  $F_\lambda$  will be assigned a 4-cycle. Observe that any 4-cycle tf-quarnet with circular ordering  $(a, c, b, d)$  (irrespective of the position of the reticulation) can be turned into the quartet trees with splits  $ac|bd$  or  $ad|bc$  by deleting exactly one reticulation edge, while this is not possible for the quartet tree with split  $ab|cd$ . Assuming that the taxa set  $\{a, b, c, d\}$  is in the set  $F_\lambda$  and that  $h_{ab|cd} \geq h_{ac|bd} \geq h_{ad|bc}$ ,

we therefore assign a 4-cycle with circular ordering  $(a, c, b, d)$  to the taxa set. This aligns with the group-based models (see, e.g. Gross et al. 2021; Barton et al. 2022) which also assume that DNA independently evolves along the trees that can be obtained from a network by deleting reticulation edges.

We also assign a weight  $w(q)$  to each tf-quarnet  $q$ , corresponding to the difference its  $\delta$ -value has from  $\lambda$ . In some sense, this weight signifies the confidence we have in having estimated the correct tf-quarnet. In particular,

$$w(q) = \begin{cases} \frac{|\delta_q - \lambda|}{\lambda} & \text{if } \delta_q \leq \lambda, \\ \frac{|\delta_q - \lambda|}{1 - \lambda} & \text{if } \delta_q > \lambda. \end{cases} \quad (8)$$

**Step IV.** It remains to determine where to place the reticulations in the 4-cycles obtained from the set  $F_\lambda$ . Taking inspiration from Holland et al. (2002) and Oldman et al. (2016), we first compute the value  $\delta(x)$  for each taxon  $x$ , defined as the mean value of all  $\delta$ -values for four-taxon sets containing  $x$ . For each 4-cycle, we then let the leaf  $x$  with the highest  $\delta(x)$ -value be below the reticulation.

### Consistency of SQUIRREL

In [supplementary material Section A, Supplementary Material online](#), we prove that SQUIRREL is combinatorially consistent given, an unweighted dense set of tf-quarnets. We use the word “combinatorially” to emphasize that we do not make any claims regarding statistical consistency. More formally, we prove the following theorem.

**Theorem 1** Let  $\mathcal{N}$  be a triangle-free semi-directed level-1 network and let  $\mathcal{Q}$  be the set of unweighted tf-quarnets induced by  $\mathcal{N}$ , then SQUIRREL applied to  $\mathcal{Q}$  reconstructs  $\mathcal{N}$ .

The first ingredient of the proof is the fact that if a set of tf-quarnets is induced by a network, the tree  $\mathcal{T}^*$  is equal to the blobtree of that network. The other important step of the proof is to show that in this case the distances  $D$  [as defined in equation (6)] form a *Kalmanson metric* (Kalmanson 1975), which have nice properties with respect to the TSP.

### Implementation

A graphical user interface (implemented in Python) of SQUIRREL and the  $\delta$ -heuristic is freely available at <https://github.com/nholtgreffe/squirrel>. The program takes as input a sequence alignment in NEXUS or FASTA format, or a file specifying a dense set of tf-quarnets [e.g. coming from QNR-SVM {Barton et al. 2022} or the MML algorithm {Martin et al. 2023}]. The interface allows the user to specify an optional outgroup, view the different generated candidate networks, and export them in the eNewick file-format (Cardona et al. 2008) (with an arbitrary rooting if no outgroup was specified).

### Supplementary Material

[Supplementary material](#) is available at *Molecular Biology and Evolution* online.

## Acknowledgments

We thank the authors of [Allman et al. \(2024b\)](#) for making us aware of their new approach and the reviewers for their helpful comments and suggestions to improve the paper.

## Funding

This work received funding from grants OCENW.M.21.306 (N.H., L.v.I., and M.J.) and OCENW.KLEIN.125 (L.v.I. and M.J.) of the Dutch Research Council (NWO). Part of this work was done while some of the authors were in residence at the Institute for Computational and Experimental Research in Mathematics (ICERM) in Providence (RI, USA) during the *Theory, Methods, and Applications of Quantitative Phylogenomics* program [supported by grant DMS-1929284 of the National Science Foundation (NSF)].

## Data Availability

The generated networks, Python scripts, sequence alignments and numerical results of the experiments in this paper are available at <https://github.com/nholtgreffe/squirrel>.

## References

- Allman ES, Baños H, Mitchell JD, Rhodes JA. TINNIk: inference of the tree of blobs of a species network under the coalescent model. *Algorithms Mol Biol.* 2024a;19(1):23. <https://doi.org/10.1186/s13015-024-00266-2>.
- Allman ES, Baños H, Rhodes JA. NANUQ: a method for inferring species networks from gene trees under the coalescent model. *Algorithms Mol Biol.* 2019;14(1):1–25. <https://doi.org/10.1186/s13015-019-0159-2>.
- Allman ES, Baños H, Rhodes JA, Wicke K. NANUQ+: a divide-and-conquer approach to network estimation. bioRxiv 621146. <https://doi.org/10.1101/2024.10.30.621146>, 2024b, preprint: not peer reviewed.
- Baños H. Identifying species network features from gene tree quartets under the coalescent model. *Bull Math Biol.* 2019;81(2):494–534. <https://doi.org/10.1007/s11538-018-0485-4>.
- Barton T, Gross E, Long C, Rusinko J. 2022. Statistical learning with phylogenetic network invariants, arXiv, arXiv:2211.11919, preprint: not peer reviewed.
- Bellman R. Dynamic programming treatment of the travelling salesman problem. *J ACM (JACM).* 1962;9(1):61–63. <https://doi.org/10.1145/321105.321111>.
- Berry V, Gascuel O. Inferring evolutionary trees with strong combinatorial evidence. *Theor Comput Sci.* 2000;240(2):271–298. [https://doi.org/10.1016/S0304-3975\(99\)00235-2](https://doi.org/10.1016/S0304-3975(99)00235-2).
- Cardona G, Rosselló F, Valiente G. Extended Newick: it is time for a standard representation of phylogenetic networks. *BMC Bioinformatics.* 2008;9(1):1–8. <https://doi.org/10.1186/1471-2105-9-1>.
- Chifman J, Kubatko L. Quartet inference from SNP data under the coalescent model. *Bioinformatics.* 2014;30(23):3317–3324. <https://doi.org/10.1093/bioinformatics/btu530>.
- Cui R, Schumer M, Kruesi K, Walter R, Andolfatto P, Rosenthal GG. Phylogenomics reveals extensive reticulate evolution in *Xiphophorus* fishes. *Evolution.* 2013;67(8):2166–2179. <https://doi.org/10.1111/evo.12099>.
- Culumber Z, Fisher H, Tobler M, Mateos M, Barber P, Sorenson M, Rosenthal G. Replicated hybrid zones of *Xiphophorus* swordtails along an elevational gradient. *Mol Ecol.* 2011;20(2):342–356. <https://doi.org/10.1111/mec.2010.20.issue-2>.
- Diop A, Torrance EL, Stott CM, Bobay LM. Gene flow and introgression are pervasive forces shaping the evolution of bacterial species. *Genome Biol.* 2022;23(1):239. <https://doi.org/10.1186/s13059-022-02809-5>.
- Du K, Ricci JMB, Lu Y, Garcia-Olazabal M, Walter RB, Warren WC, Dodge TO, Schumer M, Park H, Meyer A, et al. Phylogenomic analyses of all species of swordtail fishes (genus *Xiphophorus*) show that hybridization preceded speciation. *Nat Commun.* 2024;15(1):6609. <https://doi.org/10.1038/s41467-024-50852-6>.
- Ehrendorfer F. Differentiation-hybridization cycles and polyploidy in achillea. *Cold Spring Harb Symp Quant Biol.* 1959;24:141–152. <https://doi.org/10.1101/SQB.1959.024.01.014>.
- Frohn M, Holtgreffe N, van Iersel L, Jones M, Kelk S. Reconstructing semi-directed level-1 networks using few quartets. *J Comput Syst Sci.* 2025;152:103655. <https://doi.org/10.1016/j.jcss.2025.103655>.
- Green RE, Krause J, Briggs AW, Maricic T, Stenzel U, Kircher M, Patterson N, Li H, Zhai W, Fritz MHY, et al. A draft sequence of the Neandertal genome. *Science.* 2010;328(5979):710–722. <https://doi.org/10.1126/science.1188021>.
- Gross E, van Iersel L, Janssen R, Jones M, Long C, Murakami Y. Distinguishing level-1 phylogenetic networks on the basis of data generated by Markov processes. *J Math Biol.* 2021;83(3):1–24. <https://doi.org/10.1007/s00285-021-01653-8>.
- Grünwald S, Moulton V, Spillner A. Consistency of the QNet algorithm for generating planar split networks from weighted quartets. *Discret Appl Math.* 2009;157(10):2325–2334. <https://doi.org/10.1016/j.dam.2008.06.038>.
- Held M, Karp RM. A dynamic programming approach to sequencing problems. *J Soc Ind Appl Math.* 1962;10(1):196–210. <https://doi.org/10.1137/01110015>.
- Holland BR, Huber KT, Dress A, Moulton V.  $\delta$  plots: a tool for analyzing phylogenetic distance data. *Mol Biol Evol.* 2002;19(12):2051–2059. <https://doi.org/10.1093/oxfordjournals.molbev.a004030>.
- Huber KT, van Iersel L, Jones M, Moulton V, Veenema-Nipius L. 2024. When are quartets sufficient to reconstruct semi-directed phylogenetic networks? arXiv, arXiv:2408.12997, preprint: not peer reviewed.
- Huber KT, van Iersel L, Kelk S, Suchecchi R. A practical algorithm for reconstructing level-1 phylogenetic networks. *IEEE/ACM Trans Comput Biol Bioinform.* 2010;8(3):635–649. <https://doi.org/10.1109/TCBB.2010.17>.
- Jiao Y, An M, Zhang N, Zhang H, Zheng C, Chen L, Li H, Zhang Y, Gan Y, Zhao J, et al. Multiple third-generation recombinants formed by CRF55\_01B and CRF07\_BC in newly diagnosed HIV-1 infected patients in Shenzhen city, China. *Virol J.* 2024;21(1):306. <https://doi.org/10.1186/s12985-024-02563-z>.
- Kalmanson K. Edgeconvex circuits and the traveling salesman problem. *Can J Math.* 1975;27(5):1000–1010. <https://doi.org/10.4153/CJM-1975-104-6>.
- Kang JH, Schartl M, Walter RB, Meyer A. Comprehensive phylogenetic analysis of all species of swordtails and platies (Pisces: Genus *Xiphophorus*) uncovers a hybrid origin of a swordtail fish, *Xiphophorus monticolus*, and demonstrates that the sexually selected sword originated in the ancestral lineage of the genus, but was lost again secondarily. *BMC Evol Biol.* 2013;13(1):1–19. <https://doi.org/10.1186/1471-2148-13-25>.
- Kong S, Swofford DL, Kubatko LS. Inference of phylogenetic networks from sequence data using composite likelihood. *Syst Biol.* 2024;74(1):53–69. <https://doi.org/10.1093/sysbio/syae054>.
- Lemey P, Salemi M, Vandamme AM. *The phylogenetic handbook: a practical approach to phylogenetic analysis and hypothesis testing.* Cambridge: Cambridge University Press; 2009.
- Martin S, Moulton V, Leggett RM. Algebraic invariants for inferring 4-leaf semi-directed phylogenetic networks. bioRxiv 557152. <https://doi.org/10.1101/2023.09.11.557152>, 2023, preprint: not peer reviewed.
- Maxfield M, Xu J, Ané C. A dissimilarity measure for semidirected networks. *IEEE Trans Comput Biol Bioinform.* 2025;1–14. <https://doi.org/10.1109/TCBBIO.2025.3534780>.
- Meier JI, Stelkens RB, Joyce DA, Mwaiko S, Phiri N, Schlieven UK, Selz OM, Wagner CE, Katongo C, Seehausen O. The coincidence of

- ecological opportunity with hybridization explains rapid adaptive radiation in Lake Mweru cichlid fishes. *Nat Commun.* 2019;10(1):5391. <https://doi.org/10.1038/s41467-019-13278-z>.
- Meyer A, Salzburger W, Scharl M. Hybrid origin of a swordtail species (Teleostei: *Xiphophorus clemenciae*) driven by sexual selection. *Mol Ecol.* 2006;15(3):721–730. <https://doi.org/10.1111/mec.2006.15.issue-3>.
- Oldman J, Wu T, van Iersel L, Moulton V. TriLoNet: piecing together small networks to reconstruct reticulate evolutionary histories. *Mol Biol Evol.* 2016;33:2151–2162. <https://doi.org/10.1093/molbev/msw068>.
- Patterson N, Richter DJ, Gnerre S, Lander ES, Reich D. Genetic evidence for complex speciation of humans and chimpanzees. *Nature.* 2006;441:1103–1108. <https://doi.org/10.1038/nature04789>.
- Pekar J, Worobey M, Moshiri N, Scheffler K, Wertheim JO. Timing the SARS-CoV-2 index case in Hubei province. *Science.* 2021;372:412–417. <https://doi.org/10.1126/science.abf8003>.
- Perelman P, Johnson WE, Roos C, Seuánez HN, Horvath JE, Moreira MAM, Kessing B, Pontius J, Roelke M, Rumpler Y, et al. A molecular phylogeny of living primates. *PLoS Genet.* 2011;7(3):1–17. <https://doi.org/10.1371/journal.pgen.1001342>.
- Rambaut A, Grass NC. Seq-Gen: an application for the Monte Carlo simulation of DNA sequence evolution along phylogenetic trees. *Bioinformatics.* 1997;13(3):235–238. <https://doi.org/10.1093/bioinformatics/13.3.235>.
- Rhodes JA, Baños H, Xu J, Ané C. Identifying circular orders for blobs in phylogenetic networks. *Adv Appl Math.* 2025;163(1):102804. <https://doi.org/10.1016/j.aam.2024.102804>.
- Rieseberg LH, Raymond O, Rosenthal DM, Lai Z, Livingstone K, Nakazato T, Durphy JL, Schwarzbach AE, Donovan LA, Lexer C. Major ecological transitions in wild sunflowers facilitated by hybridization. *Science.* 2003;301:1211–1216. <https://doi.org/10.1126/science.1086949>.
- Rosenthal GG, de la Rosa Reyna XF, Kazianis S, Stephens MJ, Morizot DC, Ryan MJ, García de León FJ. Dissolution of sexual signal complexes in a hybrid zone between the swordtails *Xiphophorus birchmanni* and *Xiphophorus malinche* (Poeciliidae). *Copeia.* 2003;2003(2):299–307. [https://doi.org/10.1643/0045-8511\(2003\)003\[0299:DOSSCI\]2.0.CO;2](https://doi.org/10.1643/0045-8511(2003)003[0299:DOSSCI]2.0.CO;2).
- Salemi M, Vandamme AM. *The phylogenetic handbook: a practical approach to DNA and protein phylogeny*. Cambridge: Cambridge University Press; 2003.
- Schmidt HA, Strimmer K, Vingron M, von Haeseler A. Tree-puzzle: maximum likelihood phylogenetic analysis using quartets and parallel computing. *Bioinformatics.* 2002;18:502–504. <https://doi.org/10.1093/bioinformatics/18.3.502>.
- Schumer M, Cui R, Boussau B, Walter R, Rosenthal G, Andolfatto P. An evaluation of the hybrid speciation hypothesis for *Xiphophorus clemenciae* based on whole genome sequences. *Evolution.* 2013;67:1155–1168. <https://doi.org/10.1111/evo.12009>.
- Simmons MP, Gatesy J. Coalescence vs. concatenation: sophisticated analyses vs. first principles applied to rooting the angiosperms. *Mol Phylogenet Evol.* 2015;91:98–122. <https://doi.org/10.1016/j.ympev.2015.05.011>.
- Solís-Lemus C, Ané C. Inferring phylogenetic networks with maximum pseudolikelihood under incomplete lineage sorting. *PLoS Genet.* 2016;12(3):e1005896. <https://doi.org/10.1371/journal.pgen.1005896>.
- Solís-Lemus C, Bastide P, Ané C. PhyloNetworks: a package for phylogenetic networks. *Mol Biol Evol.* 2017;34(12):3292–3298. <https://doi.org/10.1093/molbev/msx235>.
- Soucy SM, Huang J, Gogarten JP. Horizontal gene transfer: building the web of life. *Nat Rev Genet.* 2015;16(8):472–482. <https://doi.org/10.1038/nrg3962>.
- Springer MS, Meredith RW, Gatesy J, Emerling CA, Park J, Rabosky DL, Stadler T, Steiner C, Ryder OA, Janecka JE, et al. Macroevolutionary dynamics and historical biogeography of primate diversification inferred from a species supermatrix. *PLoS One.* 2012;7:1–23. <https://doi.org/10.1371/journal.pone.0049521>.
- Steenwyk JL, Li Y, Zhou X, Shen XX, Rokas A. Incongruence in the phylogenomics era. *Nat Rev Genet.* 2023;24(12):834–850. <https://doi.org/10.1038/s41576-023-00620-x>.
- Taylor SA, Larson EL. Insights from genomes into the evolutionary importance and prevalence of hybridization in nature. *Nat Ecol Evol.* 2019;3:170–177. <https://doi.org/10.1038/s41559-018-0777-y>.
- Than C, Ruths D, Nakhleh L. PhyloNet: a software package for analyzing and reconstructing reticulate evolutionary relationships. *BMC Bioinformatics.* 2008;9(1):1–16. <https://doi.org/10.1186/1471-2105-9-322>.
- Tung J, Barreiro LB. The contribution of admixture to primate evolution. *Curr Opin Genet Dev.* 2017;47:61–68. <https://doi.org/10.1016/j.gde.2017.08.010>.
- Vanderpool D, Minh BQ, Lanfear R, Hughes D, Murali S, Harris RA, Raveendran M, Muzny DM, Hibbins MS, Williamson RJ, et al. Primate phylogenomics uncovers multiple rapid radiations and ancient interspecific introgression. *PLoS Biol.* 2020;18(12):1–27. <https://doi.org/10.1371/journal.pbio.3000954>.
- Warnow T, Tabatabaee Y, Evans SN. Advances in estimating level-1 phylogenetic networks from unrooted SNPs. *J Comput Biol.* 2024;32(1):3–27. <https://doi.org/10.1089/cmb.2024.0710>.
- Worobey M, Gemmel M, Teuwen DE, Haselkorn T, Kunstman K, Bunce M, Muyembe JJ, Kabongo JMM, Kalengayi RM, Van Marck E, et al. Direct evidence of extensive diversity of HIV-1 in Kinshasa by 1960. *Nature.* 2008;455(7213):661–664. <https://doi.org/10.1038/nature07390>.
- Wu Z, Solís-Lemus C. Ultrafast learning of four-node hybridization cycles in phylogenetic networks using algebraic invariants. *Bioinform Adv.* 2024;4(1):vbae014. <https://doi.org/10.1093/bioadv/vbae014>.
- Xu J, Ané C. Identifiability of local and global features of phylogenetic networks from average distances. *J Math Biol.* 2023;86(1):12. <https://doi.org/10.1007/s00285-022-01847-8>.
- Yu Y, Nakhleh L. A maximum pseudo-likelihood approach for phylogenetic networks. *BMC Genomics.* 2015;16(1):1–10. <https://doi.org/10.1186/1471-2164-16-1>.
- Zhang C, Mirarab S. Weighting by gene tree uncertainty improves accuracy of quartet-based species trees. *Mol Biol Evol.* 2022;39(12):msac215. <https://doi.org/10.1093/molbev/msac215>.
- Zhang C, Ogilvie HA, Drummond AJ, Stadler T. Bayesian inference of species networks from multilocus sequence data. *Mol Biol Evol.* 2018a;35(2):504–517. <https://doi.org/10.1093/molbev/msx307>.
- Zhang C, Rabiee M, Sayyari E, Mirarab S. ASTRAL-III: polynomial time species tree reconstruction from partially resolved gene trees. *BMC Bioinformatics.* 2018b;19(1):15–30. <https://doi.org/10.1186/s12859-018-2021-9>.
- Zhang W, Dasmahapatra KK, Mallet J, Moreira GR, Kronforst MR. Genome-wide introgression among distantly related *Heliconius* butterfly species. *Genome Biol.* 2016;17(1):1–15. <https://doi.org/10.1186/s13059-015-0866-z>.