

Improving intraspecific allele networks inferred by maximum parsimony

Vincent Branders* and Patrick Mardulyn*

Evolutionary Biology and Ecology, Université Libre de Bruxelles, av. FD Roosevelt 50, 1050 Brussels, Belgium

Summary

1. Allele (or haplotype) networks are often used in phylogeographic studies to display genetic variation within a species or a group of closely related species. A global maximum parsimony approach to infer allele networks, arguably the method of choice to display genetic variation at the intraspecific level, consists in inferring all most parsimonious trees from a DNA sequence alignment and combining the corresponding phylograms into a single graph. However, it has been suggested that, while classic phylogenetic programs generate a single phylogram per most parsimonious tree, deriving all possible phylograms from them would allow identifying additional most parsimonious paths among alleles, thereby improving this network inference method.

2. We test this prediction by analysing both simulated and empirical DNA sequence alignments. For this purpose, a computer program, CPN, was developed to implement the entire procedure, starting with a set of most parsimonious trees and combining all derived phylograms into a network.

3. We show that including all possible most parsimonious phylograms indeed often results in finding additional most parsimonious paths in the network graph, thereby improving the search for a global maximum parsimony solution.

4. We highly recommend the use of this approach in future phylogeographic studies, to ensure that all most parsimonious paths are included in the allele network, instead of an arbitrarily selected subset of those.

Key-words: haplotype network, intraspecific DNA sequence variation, phylogeography, phylogram

Introduction

In phylogeographic studies, allele or haplotype networks are typically used to display genetic variation within a species or a group of related species. While other types of phylogenetic networks exist (e.g. Huson, Rupp & Scornavacca 2011; Morrison 2011), we focus here on those graphs generated by, for example, a median-joining (Bandelt, Forster & Röhl 1999) or statistical parsimony (TCS network; Clement, Posada & Crandall 2000) algorithm, in which nodes represent different allelic sequences and are joined by edges (branches) whose length shows the number of nucleotides that differ between them (e.g. Huson, Rupp & Scornavacca 2011). In this context, an important advantage of a network graph over a phylogenetic tree lies in the use of cycles (loops) to display ambiguous signal in the data (Posada & Crandall 2001; Mardulyn 2012).

Several methods are available to infer an allele network from a set of DNA sequences, but their performances have been shown to differ (Cassens, Mardulyn & Milinkovitch 2005; Woolley, Posada & Crandall 2008; Mardulyn, Cassens & Milinkovitch 2009; Salzburger, Ewing & Von Haeseler 2011). The use of traditional phylogenetic methods, mostly maximum parsimony, appears better suited for this task, as they use an

objective function to explore the space of all possible phylogenetic hypotheses and to identify the best ones. In this sense, they provide, or seek to provide, solutions that are globally optimal. Studies that have compared the performances of network inference methods tend indeed to agree that they perform better in general (Woolley, Posada & Crandall 2008; Salzburger, Ewing & Von Haeseler 2011), or at least equally well (Cassens, Mardulyn & Milinkovitch 2005), than methods that build a network step-by-step following a well-defined procedure, but that do not compare it to other possible networks using a criterion of global optimality. Nonetheless, Mardulyn, Cassens & Milinkovitch (2009) and Mardulyn (2012) have shown examples in which a median-joining algorithm found alternative paths among alleles that were not considered by a classic maximum parsimony phylogenetic analysis. They suggested that a problem may lie with the inference of ancestral sequences for the interior nodes of the most parsimonious (MP) trees. All classic phylogenetic inference programs provide only one MP reconstruction of ancestral sequences per tree (while other equally parsimonious reconstructions are often possible), thereby generating a single phylogram (i.e. a tree that includes branch length information, as opposed to a cladogram, that conveys only topological information), which is sufficient when inferring a species tree, but may become problematic when studying intraspecific variation, because other possible MP phylograms can identify other equally parsi-

*Correspondence author: E-mails: vbranders@gmail.com, pmardulyn@ulb.ac.be

monious paths in the network graph. This suggests a simple improvement to the global parsimony approach, consisting in generating all MP phylograms corresponding to the set of identified MP cladograms (i.e. all MP combinations of ancestral sequences for the set of MP trees inferred from a DNA sequence alignment), before combining them into a network graph.

Here, we explore the added value of taking into account all MP combinations of ancestral sequences to build a network. We developed a program, CPN, that takes as input a set of MP cladograms and the corresponding DNA sequence alignment and outputs the set of corresponding MP phylograms, by taking into account all MP reconstructions of ancestral sequences, identified by an algorithm previously described in Maddison & Maddison (1992; unordered characters, pp. 92–93). Indeed, phylograms instead of cladograms are used to build the final network, as they contain branch length information that is crucial to identify the MP paths of the allele network. Taking all MP combinations of ancestral sequences into account allows the identification of additional MP phylograms which sometimes results in one or more additional MP path(s) in the final network. The same program also combines all identified phylograms in a single network using the algorithm described in Cassens, Mardulyn & Milinkovitch (2005), modified to allow reducing the length of cycles, by merging some of its edges and nodes, where possible (see program manual for a detailed description of this modification). We then simulated DNA sequence data using a classic coalescence model, analysed the resulting data sets under maximum parsimony and produced two network graphs with the resulting MP phylograms, one combining only the original set of phylograms produced by the classic phylogenetic analysis program, and the other by combining the potentially larger set of phylograms, inferred from the identified MP cladograms by taking all possible MP reconstructions of ancestral sequences into account. The comparison of these two networks allowed us to investigate the improvement provided by taking all possible phylograms into account for generating the MP network graph, or, in other words, the extent with which failing to do so produces a globally incomplete MP solution.

Materials and methods

An allele or haplotype network graph can be constructed by combining all MP trees inferred by a classic phylogenetic program from a DNA sequence alignment. Such phylogenetic program provides in general a single phylogram per MP tree, but alternative phylograms can be derived as well, as there are usually other equally parsimonious combinations of ancestral sequences for the interior nodes. We propose to improve the search for a global maximum parsimony solution by deriving all most parsimonious phylograms to produce the final graph. Using a program we developed to produce all most parsimonious phylograms from the original set of most parsimonious cladograms, we evaluate whether, and to what extent, it improves allele network inference.

We generated intraspecific DNA sequence data by simulating gene genealogies under a classic coalescent model using MS (Hudson 2002)

and by simulating the evolution of DNA sequences along the resulting genealogies with Seq-Gen (Rambaut & Grassly 1997). As shown in Table 1, we varied sample size, sequence length and mutation rate parameter among simulations, to span a large set of conditions. These conditions were chosen so as to generate a range of values regarding tree length (i.e. total number of mutations, which varies from ± 20 to ± 300 ; see Table 1) and number of cycles (denoting homoplasious states) similar to what is typically found in empirical phylogeographic data sets. A few empirical data sets were also analysed to explore the usefulness of our suggested approach for real sequences.

Most parsimonious trees were inferred from each DNA sequence alignment using PAUP (Swofford 2003), by implementing a heuristic search (simple addition sequence). The resulting set of MP cladograms, or corresponding MP phylograms produced by the same program, was given as input to our program CPN (<http://ebe.ulb.ac.be/ebe/CPN.html>). CPN was used first to simply combine the set of MP phylograms generated by PAUP into a network, then to derive the complete set of MP phylograms associated with the set of MP trees inferred by PAUP, and to combine them in a second network graph. We conducted 20 independent runs of combining phylograms with CPN, each time randomly reordering the list of phylograms to combine, as the order in which they are added to the graph can influence the final network (Cassens, Mardulyn & Milinkovitch 2005). When different networks are found among runs, CPN identifies the best network as the one with the lowest number of connections (defined as a path between two labelled/union nodes; a union node is defined as unlabelled, that is, unsampled alleles connected to more than two edges). Only data from the best network are shown. For each inferred network, three statistics were computed for comparison: number of connections, number of mutations (total network length) and number of union nodes. Defining network complexity as the number of connections that must be removed from the network to produce one of the initial MP phylogram, we use the total number of connections in the graph as a proxy to compare the complexity of graphs generated with or without inferring extra MP phylograms. CPN has an option allowing reduction of the length of each cycle by merging some of its edges and nodes, wherever possible (i.e. for each potential cycle reduction, the program checks first that each MP phylogram is still included in the network). While this reduction is desirable to produce the final network, it will also increase the total number of connections, which will no longer be proportional to network complexity. For the purpose of our comparison, the three statistics reported here are thus the one measured before cycle reduction.

Results and discussion

Tables 1 and 2 summarize the increase in complexity of the inferred network graph when considering all MP phylograms (i.e. after running the CPN program), instead of only the subset of MP phylograms generated by PAUP, respectively, for simulated and real DNA sequence alignments. While in some cases, mostly involving shorter trees (lower overall number of mutations), the complete set of inferred MP phylograms was identical to the set generated by PAUP; in many others, additional MP phylograms were found. In these cases, the complexity of the network graph increased, as evidenced by its larger number of connections. The fact that it occurs also when analysing real data sets (Table 2; three cases out of seven) indicates that it is worth implementing the extra step of searching for all MP phylograms associated with a sequence alignment,

Table 1. Comparison between networks of simulated DNA sequences (1) built from the set of MP phylograms generated by a classic phylogenetic inference program (PAUP) and (2) built from the set of all MP phylograms derived from the set of MP trees by CPN

Simulation parameters			Inferred phylograms			Network characteristics							
			Length	Count		Connections		Mutations		Union nodes		Time of analysis ¹	
<i>n</i>	<i>L</i>	<i>s</i>		PAUP	CPN	PAUP	CPN	PAUP	CPN	PAUP	CPN	PAUP	CPN
20	200	0.05	23	2	3	12	13	33	24	1	2	<1 s	<1 s
20	200	0.05	28	2	3	14	15	39	29	2	3	<1 s	<1 s
20	200	0.05	30	2	3	16	17	32	31	3	4	<1 s	<1 s
100	200	0.05	31	1	1	17	17	31	31	2	2	<1 s	<1 s
100	200	0.05	37	2	9	21	23	38	39	2	3	<1 s	<1 s
20	200	0.1	41	1	1	16	16	41	41	4	4	<1 s	<1 s
20	200	0.05	46	3	4	21	22	62	63	7	7	6 s	8 s
100	200	0.05	47	1	1	32	32	47	47	6	6	<1 s	<1 s
20	200	0.1	48	1	3	18	21	48	49	6	8	<1 s	1 s
100	200	0.05	53	1	1	28	28	53	53	6	6	<1 s	<1 s
20	200	0.1	55	2	7	21	34	58	70	7	13	5 s	56 s
50	200	0.1	58	5	216	34	63 ²	91	161	5	13	<1 s	3 h
50	200	0.1	66	2	3	32	33	68	67	5	6	<1 s	<1 s
20	200	0.1	69	4	18	21	31 ²	86	92	7	11	1 s	41 s
50	500	0.05	72	1	1	23	23	72	72	6	6	<1 s	12 min
50	200	0.1	72	2	27	36	42	76	75	8	12	2 s	5 min
100	200	0.1	74	2	2	43	43	87	87	8	8	8 s	8 s
50	200	0.1	78	1	9	30	47 ²	78	100	7	14	<1 s	56 s
50	500	0.05	81	1	3	36	39	81	82	6	8	<1 s	1 s
50	500	0.05	85	2	2	35	35	92	92	8	8	2 s	2 s
50	200	0.1	86	2	9	34	38	95	88	6	9	<1 s	59 s
50	500	0.05	87	5	6	42	43	101	97	13	14	4 min	6 min
50	200	0.1	88	1	1	29	29	88	88	4	4	<1 s	<1 s
50	200	0.1	88	1	3	36	44	88	122	9	13	<1 s	18 s
50	500	0.05	92	1	3	35	38	92	93	8	10	<1 s	6 s
50	500	0.05	92	1	18	29	42	92	106	5	11	<1 s	14 min
100	200	0.1	94	2	510	31	228 ²	97	1555	5	36	<1 s	6 days
100	200	0.1	95	1	18	40	51	95	137	6	12	<1 s	58 s
50	200	0.1	96	2	318	32	458 ²	136	5707	7	97	<1 s	6 days
50	500	0.05	98	1	2	39	43	98	106	10	12	<1 s	14 s
50	500	0.05	99	1	1	34	34	99	99	6	6	<1 s	<1 s
50	200	0.1	108	1	3	34	37	108	109	10	12	<1 s	9 s
100	200	0.1	109	4	126	43	63 ²	112	138	7	17	6 s	32 h
50	200	0.1	115	1	360	29	487	115	6861	6	156	<1 s	2 days
20	1000	0.05	126	2	4	26	29	132	133	9	11	20 s	44 s
20	1000	0.05	127	2	21	26	42	139	221	7	15	10 s	14 min
50	500	0.05	149	8	112	45	66 ²	255	180	13	21	6 min	4 days
50	500	0.05	152	2	54	36	66	170	363	7	20	1 s	3 h
20	1000	0.05	161	1	6	30	38	161	167	11	16	<1 s	4 min
20	1000	0.05	187	1	30	26	63	187	976	10	26	<1 s	1 h
20	1000	0.1	211	1	18	30	40	211	219	13	19	<1 s	1 h 39 min
20	1000	0.1	328	3	48	30	74 ²	346	654	13	27	32 s	14 h
20	1000	0.1	331	1	234	28	614	331	22 744	10	225	<1 s	6 days

Data sets ordered by tree length.

Analysis time refers to the process of combining all phylograms into a network; inferring the complete set of all MP phylograms took 10 s at most for the cases investigated here.

n, sample size; *L*, sequence length; *s*, mutation rate parameter. A connection is defined as a path between two labelled/union nodes; a union node is a unlabelled node connected to more than two branches (edges).

¹All analyses performed with the command-line version of CPN, on a computer cluster, each run associated to a 2× AMD Opteron CPU (2.2–2.4 GHz) and to a maximum of 1 GB of RAM.

²The 20 runs of CPN generated at least two different networks; data reported for the network associated with the lowest number of connections.

before combining them into a network graph. For all data sets we analysed, the extra time needed to infer all MP phylograms amounted to maximum 30 s. The time needed to produce the graph (i.e. 20 independent runs of combining all MP phylograms into a network) amounted in general to somewhere

between a few seconds and a few minutes. In the most complex cases, however, this step extended to more than 1 h of analysis per run (up to 7 h 30 per run, for the set of longest MP phylograms in Table 1). Note that when we performed several independent runs of combining a given set of MP phylograms into

Table 2. Comparison between networks of empirical DNA sequences (1) built from the set of MP phylograms generated by a classic phylogenetic inference program (PAUP) and (2) built from the set of all MP phylograms derived from the set of MP trees by CPN

Data set	Inferred phylograms			Network characteristics							
	Length	Count		Connections		Mutations		Union nodes		Time of analysis ¹	
		PAUP	CPN	PAUP	CPN	PAUP	CPN	PAUP	CPN	PAUP	CPN
Pittra <i>et al.</i> (2002)	18	4	4	11	11	19	19	2	2	<1 s	<1 s
Worheide, Hooper & Degnan (2002)	35	3	8	23	27	36	39	4	6	16 s	2 min
Caicedo & Schaal (2004)	39	18	18	33	33	63	63	6	6	2 min 24 s	2 min 30 s
Mardulyn, Mikhailov & Pasteels (2009)	39	16	16	33	33	51	51	2	2	<1 s	<1 s
Olsen (2002)	61	2	2	36	36	67	67	9	9	1 min	1 min
Balakrishnan <i>et al.</i> (2003)	69	9	28	40	48	112	120	14	16	5 min	62 min
Neiman & Lively (2004)	76	60	90	51	55	85	89	6	8	12 min	46 min

Data sets ordered by tree length.

Analysis time refers to the process of combining all phylograms into a network; inferring the complete set of all MP phylograms took 30 s at most for the cases investigated here.

A connection is defined as a path between two labelled/union nodes; a union node is a unlabelled node connected to more than two branches (edges).

¹All analyses performed with the command-line version of CPN, on a computer cluster, each run associated to a 2× AMD Opteron CPU (2.2–2.4 GHz) and to a maximum of 1 GB of RAM.

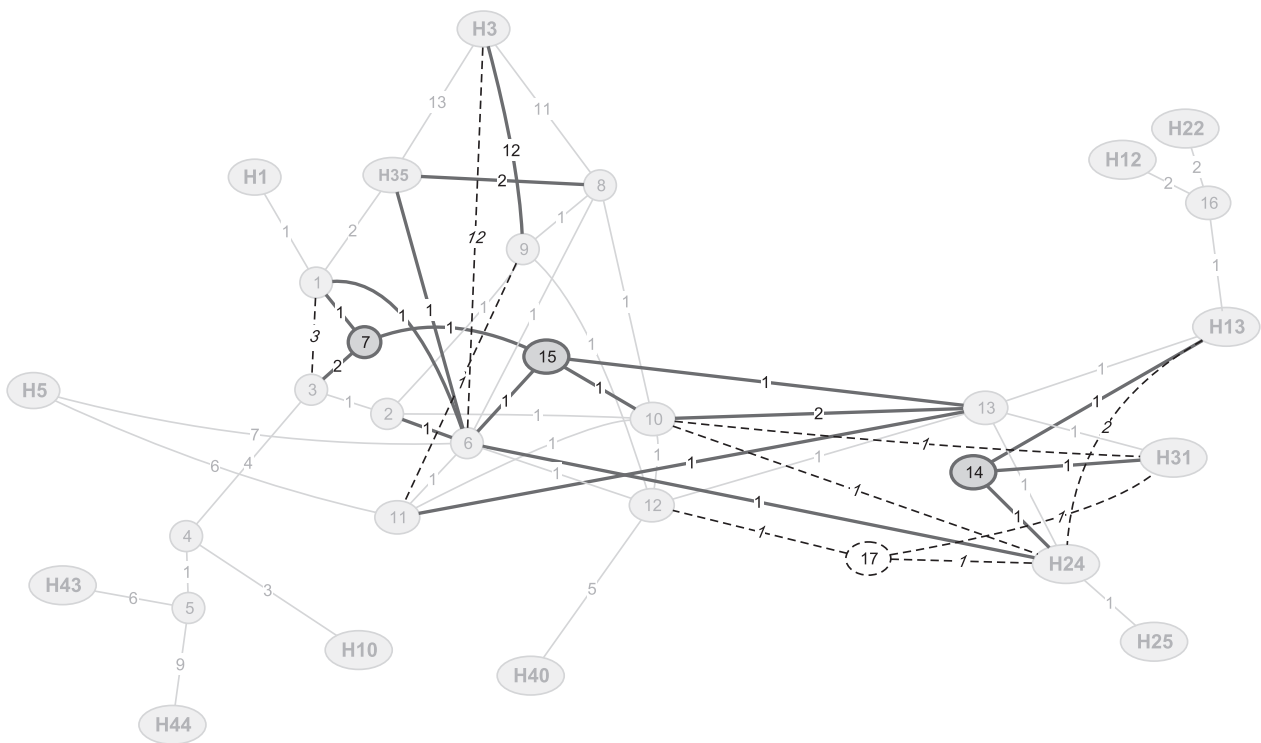


Fig. 1. Comparison between the network built by CPN using all MP phylograms (CPN network) and the network built only from the initial set of MP phylograms produced by PAUP (PAUP network), for the data set of Balakrishnan *et al.* (2003). Alleles present in the data set are labelled with the letter H, followed by a unique number, in contrast to union nodes labelled only with a number. Numbers over branches indicate branch lengths. Connections that are found in both networks are displayed as grey lines and those that are unique to each of the original network are displayed as black (CPN network) or dashed (PAUP network) lines.

a network (randomly reordering the list of phylograms and of connections to add to the network being built between runs), it resulted in more than one network only in a few cases (identified by superscript 2 in Table 1). In these cases, the difference between networks was small.

For a single simulated data set (result not shown), the algorithm could not combine all phylograms in a reasonable amount of time (limit set to 8 h for a single run) given the large number of inferred MP phylograms and the high complexity of the resulting network. However, this concerned one of the

longest MP trees (296 mutations) analysed, and such a high number of mutations is seldom observed for intraspecific data sets of similar sequence length (1000 bp). In this case, the probability of multiple substitutions occurring at the same sites is no longer negligible, and other criteria, such as maximum likelihood, might in fact become more appropriate.

Is the observed increase in complexity of the inferred network desirable? One could be tempted to avoid searching for all MP phylograms, preferring to generate a simpler allele network, which would then be easier to interpret. However, doing so would amount to arbitrarily select only a subset of all MP phylograms, and therefore, only a subset of the MP paths in the network graph. On the contrary, we argue that if relying on the maximum parsimony criterion, it is important to include all MP connections in the final graph, otherwise it represents only a partial MP network.

To illustrate further the increase in network complexity that can accompany the inclusion of all MP phylograms, Fig. 1 compares the network built by CPN from all MP phylograms (CPN network), for one of the empirical data set of Table 2, with the network built from the initial set of MP phylograms generated by a standard MP analysis (PAUP network). While the number of connections is higher for the CPN network (as shown in Table 2), some connections of the PAUP network are not included in the CPN network. This is because the inclu-

sion of additional MP phylograms by CPN, compared to the initial set identified by PAUP, led the program to select alternative connections to produce the best final network (i.e. the network with the overall minimum number of connections). Interestingly, the median-joining network generated by the program Network (Bandelt, Forster & Röhl 1999; available at <http://www.fluxus-engineering.com/sharenet.htm>) with the same data set produces also unique connections compared to the CPN network (Fig. 2). However, those connections are less parsimonious: in fact, every phylogram that can be produced from the median-joining network by deleting alternative connections (thus breaking cycles in the graph) are less parsimonious. It was already noted by Bandelt, Forster & Röhl (1999) and further highlighted by Cassens *et al.* (2003); Cassens, Mardulyn & Milinkovitch (2005), that a median-joining network does not always contain the set of all MP trees, which is not surprising given that the median-joining algorithm does not use a criterion of global optimality to evaluate the network.

Note that the global parsimony approach to infer allele networks relies on the ability of the initial maximum parsimony analysis to infer all MP phylogenetic trees. Because the complexity of this task increases dramatically with the number of sequences in the data set (e.g. Felsenstein 2004), it is recommended to collapse all copies of the same allele/haplotype to a single sequence, thereby reducing the size of the alignment

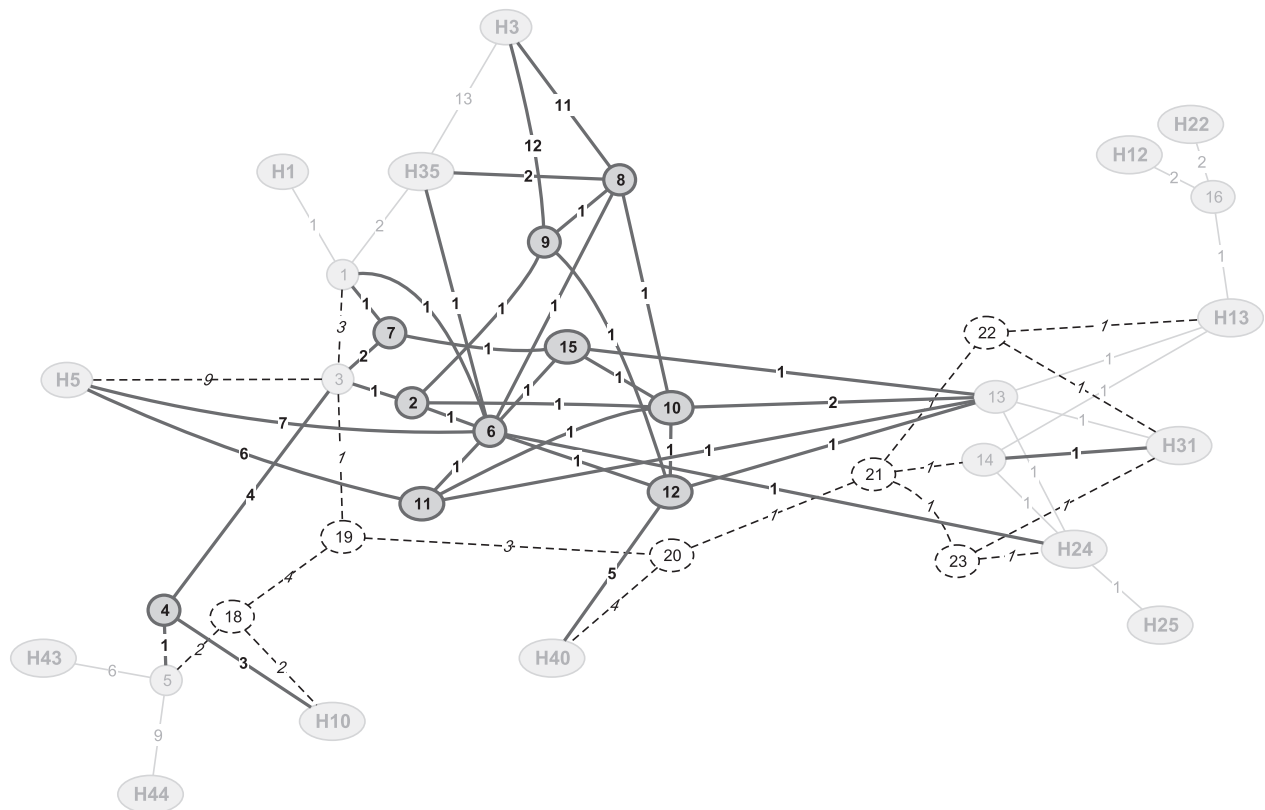


Fig. 2. Comparison between the network built by CPN using all MP phylograms (CPN network) and the median-joining network built with the program network, for the data set of Balakrishnan *et al.* (2003). Alleles present in the data set are labelled with the letter H, followed by a unique number, in contrast to union nodes labelled only with a number. Numbers over branches indicate branch lengths. Connections that are found in both networks are displayed as grey lines, and those that are unique to each of the original network are displayed as black (CPN network) or dashed (median-joining network) lines.

used to infer the network. Nonetheless, given the often large number of alleles analysed when exploring genetic variation at the intraspecies level, the use of heuristic searches to find all MP phylogenetic trees is usually required. This means we can increase the performance of the analysis by using well-designed heuristic strategies, but that we have no guarantee that all MP cladograms, and therefore all MP phylograms, will be identified and included in the final network graph. While this could be seen as a weakness of the global parsimony approach to network inference, we rather see this as an inherent problem associated with analysing large spaces of possible historical hypotheses. Selecting a simpler method that generates a network from a sequence alignment without exploring all possible hypotheses would actually lead even more likely to an incomplete MP network.

In conclusion, when choosing global maximum parsimony to infer a network graph from a set of aligned sequences, it is important to derive all MP phylograms from the set of MP cladograms initially inferred, as extra phylograms can sometimes identify new MP connections that would have otherwise been missed.

Acknowledgements

We are grateful to A. JesusMuñoz-Pajares and David Morrison for helpful comments on an earlier version of the manuscript.

Data accessibility

Empirical data sets analysed (DOI): 10.1046/j.1365-294X.2002.01516.x; 10.1046/j.1365-294X.2002.01570.x; 10.1111/j.1365-294X.2004.02191.x; 10.1111/j.1558-5646.2009.00755.x; 10.1046/j.1365-294X.2002.01493.x; 10.1046/j.1365-294X.2003.01751.x; 10.1111/j.1365-294X.2004.02292.x.

About the program

CPN is written in java, can be run under Windows, Mac OSX or Linux and can be downloaded, along with a user manual, from the following URL: <http://ebe.ulb.ac.be/ebe/CPN.html>.

References

- Balakrishnan, C.N., Monfort, S.L., Gaur, A., Sing, L. & Sorenson, M.D. (2003) Phylogeography and conservation genetics of Eld's deer (*Cervus eldi*). *Molecular Ecology*, **12**, 1–10.
- Bandelt, H.J., Forster, P. & Röhl, A. (1999) Median-joining networks for inferring intraspecific phylogenies. *Molecular Biology and Evolution*, **16**, 37–48.
- Caicedo, A.L. & Schaal, B.A. (2004) Population structure and phylogeography of *Solanum pimpinellifolium* inferred from a nuclear gene. *Molecular Ecology*, **13**, 1871–1882.
- Cassens, I., Mardulyn, P. & Milinkovitch, M.C. (2005) Evaluating intraspecific “network” construction methods using simulated sequence data: do existing algorithms outperform the global maximum parsimony approach? *Systematic Biology*, **54**, 363–372.
- Cassens, I., Van Waerebeek, K., Best, P.B., Crespo, E.A., Reyes, J. C. & Milinkovitch, M.C. (2003) The phylogeography of dusky dolphins (*Lagenorhynchus obscurus*): a critical examination of network methods and rooting procedures. *Molecular Ecology*, **12**, 1781–1792.
- Clement, M., Posada, D. & Crandall, K.A. (2000) TCS: a computer program to estimate gene genealogies. *Molecular Ecology*, **9**, 1657–1659.
- Felsenstein, J. (2004) *Inferring Phylogenies*. Sinauer Associates, Sunderland, Massachusetts, USA.
- Hudson, R.R. (2002) Generating samples under a Wright-Fisher neutral model of genetic variation. *Bioinformatics*, **18**, 337–338.
- Huson, D.H., Rupp, R. & Scornavacca, C. (2011) *Phylogenetic Networks: Concepts, Algorithms and Applications*. Cambridge University Press, Cambridge, UK.
- Maddison, W.P. & Maddison, D.R. (1992) *MacClade Version 3. Analysis of Phylogeny and Character Evolution*. Sinauer Associates, Sunderland, Massachusetts, USA.
- Mardulyn, P. (2012) Trees and/or networks to display intraspecific DNA sequence variation? *Molecular Ecology*, **21**, 3385–3390.
- Mardulyn, P., Cassens, I. & Milinkovitch, M.C. (2009) A comparison of methods for constructing evolutionary networks from intraspecific DNA sequences. *Population Genetics for Animal Conservation*, (eds G. Bertorelle, M.D. Bruford, H.C. Hauffe, A. Rizzoli & C. Vernesi), pp. 102–118. Cambridge University Press, Cambridge, UK.
- Mardulyn, P., Mikhailov, Y.E. & Pasteels, J.M. (2009) Testing phylogeographic hypotheses in a Euro-Siberian cold-adapted leaf beetle with coalescent simulations. *Evolution*, **63**, 2717–2729.
- Morrison, D.A. (2011) *Introduction to Phylogenetic Networks*. RJR Productions, Uppsala, Sweden.
- Neiman, M. & Lively, C.M. (2004) Pleistocene glaciation is implicated in the phylogeographical structure of *Potamopyrgus antipodarum*, a New Zealand snail. *Molecular Ecology*, **13**, 3085–3098.
- Olsen, K.M. (2002) Population history of *Manihot esculenta* (Euphorbiaceae) inferred from nuclear DNA sequences. *Molecular Ecology*, **11**, 901–911.
- Pitra, C., Hansen, A.J., Lieckfeldt, D. & Arctander, P. (2002) An exceptional case of historical outbreeding in African sable antelope populations. *Molecular Ecology*, **11**, 1197–1208.
- Posada, D. & Crandall, K.A. (2001) Intraspecific gene genealogies: trees grafting into networks. *Trends in Ecology and Evolution*, **16**, 37–45.
- Rambaut, A. & Grassly, N.C. (1997) Seq-Gen: an application for the Monte Carlo simulation of DNA sequence evolution along phylogenetic trees. *Computer Applications in the Biosciences*, **13**, 235–238.
- Salzburger, W., Ewing, G.B. & Von Haeseler, A. (2011) The performance of phylogenetic algorithms in estimating haplotype genealogies with migration. *Molecular Ecology*, **20**, 1952–1963.
- Swofford, D.L. (2003) *PAUP*, Phylogenetic Analysis Using Parsimony (*and Other Methods)*, v. 4b10. Sinauer Associates, Sunderland, Massachusetts, USA.
- Woolley, S.M., Posada, D. & Crandall, K.A. (2008) A comparison of phylogenetic network methods using computer simulation. *PLoS One*, **3**, e1913.
- Worheide, G., Hooper, J.N.A. & Degnan, B.M. (2002) Phylogeography of western Pacific *Leucetta 'chagosensis'* (Porifera: Calcarea) from ribosomal DNA sequences: implications for population history and conservation of the Great Barrier Reef World Heritage Area (Australia). *Molecular Ecology*, **11**, 1753–1768.

Received 5 May 2015; accepted 21 July 2015

Handling Editor: Oscar Gaggiotti