

# Network analysis identifies weak and strong links in a metapopulation system

Alejandro F. Rozenfeld<sup>a,1,2</sup>, Sophie Arnaud-Haond<sup>b,c,2</sup>, Emilio Hernández-García<sup>d</sup>, Víctor M. Eguíluz<sup>d</sup>, Ester A. Serrão<sup>b</sup>, and Carlos M. Duarte<sup>a</sup>

<sup>a</sup>Instituto Mediterraneo de Estudios Avanzados (Consejo Superior de Investigaciones Científicas–Universidad de las Islas Baleares), C/Miquel Marqués 21, 07190 Esporles, Mallorca, Spain; <sup>b</sup>Centro de Ciências do Mar do Algarve, Centro Interdisciplinar de Investigação Marinha e Ambiental–Laboratório Associado, Universidade do Algarve, Gambelas, 8005-139 Faro, Portugal; <sup>c</sup>Instituto de Física Interdisciplinar y Sistemas Complejos (Consejo Superior de Investigaciones Científicas–Universidad de las Islas Baleares), Campus Universitat de les Illes Balears, E-07122 Palma de Mallorca, Spain; and <sup>d</sup>Institut Français de Recherche pour l'Exploitation de la Mer, Centre de Brest, BP70, 29280 Plouzané, France

Edited by James H. Brown, University of New Mexico, Albuquerque, NM, and approved October 16, 2008 (received for review June 8, 2008)

The identification of key populations shaping the structure and connectivity of metapopulation systems is a major challenge in population ecology. The use of molecular markers in the theoretical framework of population genetics has allowed great advances in this field, but the prime question of quantifying the role of each population in the system remains unresolved. Furthermore, the use and interpretation of classical methods are still bounded by the need for a priori information and underlying assumptions that are seldom respected in natural systems. Network theory was applied to map the genetic structure in a metapopulation system by using microsatellite data from populations of a threatened seagrass, *Posidonia oceanica*, across its whole geographical range. The network approach, free from a priori assumptions and from the usual underlying hypotheses required for the interpretation of classical analyses, allows both the straightforward characterization of hierarchical population structure and the detection of populations acting as hubs critical for relaying gene flow or sustaining the metapopulation system. This development opens perspectives in ecology and evolution in general, particularly in areas such as conservation biology and epidemiology, where targeting specific populations is crucial.

conservation biology | gene flow | networks | population genetics

Understanding the connectivity between components of a metapopulation system and their role as weak or strong links remains a major challenge of population ecology (1–3). Advances in molecular biology fostered the use of indirect approaches to understand metapopulation structure, based on describing the distribution of gene variants (alleles) in space within the theoretical framework of population genetics (4–7). Yet, the premises of the classical Wright–Fisher model (4, 6), such as “migration-drift” and “mutation-drift” equilibrium (8), “equal population sizes” or symmetrical rate migration among populations, are often violated in real metapopulation systems. Threatened or pathogen species, for example, are precisely studied for their state of demographic disequilibrium due to decline and local extinctions in the first case, or to their complex dynamics of local decline and sudden pandemic burst in the second. Furthermore, the underlying hypotheses of equal population size and symmetrical migration rates hamper the identification of putative population “hubs” centralizing migration pathways or acting as sources in a metapopulation system, which is a central issue in ecology in general, and in conservation biology or epidemiology in particular. Finally, complementary methods of genetic structure analyses, such as hierarchical AMOVA and coalescent methods rely on a priori information (or priors) as to the clustering or demographic state of populations, requiring either subjective assumptions or the availability of reliable demographic, historical or ecological information that are seldom available.

Network theory is emerging as a powerful tool to understand the behavior of complex systems composed of many interacting

units (9–11). This approach has been applied to solve a broad array of problems (12–14). In an ecological context, it has been applied to represent geographical landscape connectivity (15, 16), but only recently it has been adapted to represent genetic relationships among populations or individuals (17, 18). Yet, relevant properties of networks, such as resistance (9) to perturbations (i.e., node paralysis or destruction), the ability to host coherent oscillations (19) or the predominant importance of nodes or clusters of nodes in maintaining the integrity of the system or relaying information through it can be deduced from the network topology and specific characteristics (10, 11). Here, we apply network theory to population genetics data of a threatened species, the Mediterranean clonal seagrass *Posidonia oceanica*. This clonal seagrass endemic to the Mediterranean has the slowest clonal growth rate of all seagrasses (20) and, although major sexual reproduction events are rare in time (21), the relative contribution of sexual versus clonal reproduction may be relatively high in some populations (22). *P. oceanica* clones produce both male and female flowers, and the fruits are buoyant and can drift tens of kilometers before they lose buoyancy and the seed settles in the seafloor initiating a clone (21). *P. oceanica* can also disperse from fragments, which can be transported by currents and can eventually become rooted at distant populations, as shown by recent analysis of patch formation in *P. oceanica* populations (23). Moreover, limited dispersal has been inferred (22) from its high population genetic structure at geographic scales ranging from the whole Mediterranean to regional and even local meadow scales.

We start by analyzing data at the Mediterranean scale, where clustering of the populations distributed in 2 basins connected by a narrow strait and that were almost isolated during the last glaciation, was rather obvious a priori. This particular geographical and historical context facilitated the classical analysis of molecular variance (AMOVA) and allows using this example to validate our network analysis and confirm its potential, without a priori knowledge or assumptions, to characterize population genetic structure and to identify populations that are critical to the dynamics and sustainability of the whole system. We then compared classical and network tools at the regional scale of the

Author contributions: A.F.R., S.A.-H., E.H.-G., V.M.E., E.A.S., and C.M.D. designed research; A.F.R., S.A.-H., and E.H.-G. performed research; A.F.R. and V.M.E. contributed new reagents/analytic tools; A.F.R., S.A.-H., E.A.S., and C.M.D. analyzed data; and A.F.R., S.A.-H., E.H.-G., V.M.E., E.A.S., and C.M.D. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Freely available online through the PNAS open access option.

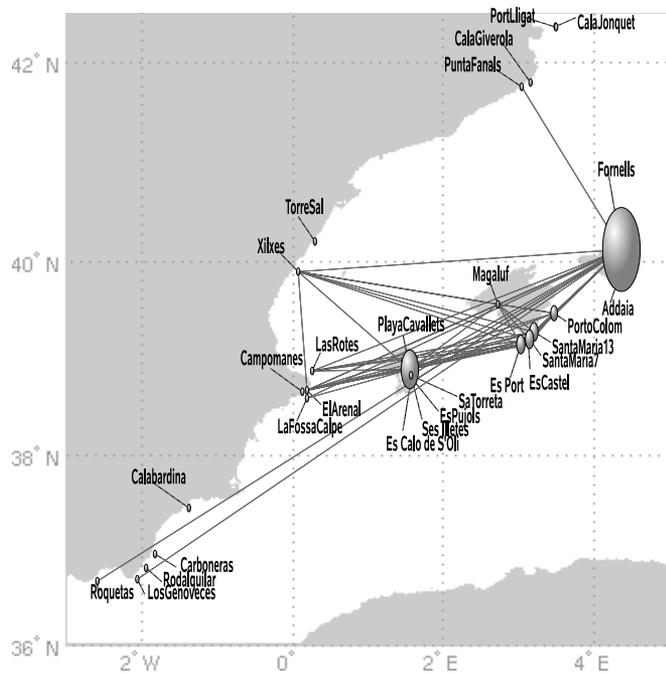
<sup>1</sup>To whom correspondence should be addressed. E-mail: alex@ifisc.uib.es.

<sup>2</sup>A.F.R. and S.A.-H. contributed equally to this work.

This article contains supporting information online at [www.pnas.org/cgi/content/full/0805571105/DCSupplemental](http://www.pnas.org/cgi/content/full/0805571105/DCSupplemental).

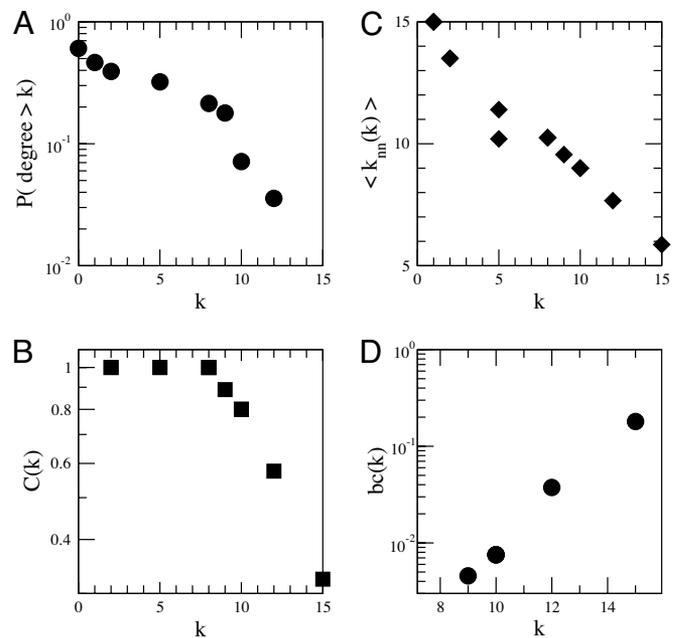
© 2008 by The National Academy of Sciences of the USA





**Fig. 2.** The network constructed for the Spanish meadows with the “geographic threshold” criterion (see Fig. 6). Nodes are shown at the populations’ geographic locations. Node sizes characterize their betweenness centrality (i.e., the proportion of all shortest paths getting through the node).

pair) without clear relationship with geographic location, no clear pattern of allelic richness, and a “comb-like” topology in an unweighted pair group method with arithmetic mean (UPGMA) tree, which forces dichotomous branching of the metapopulation network (Fig. 4A; see *Materials and Methods*). These methods were unable to highlight neither any particularly central position for the populations analyzed nor the clustering of some subgroups that would have suggested preferential roads for gene flow or a dominant role of some populations in the metapopulation system. On the contrary, network analyses of these populations (Fig. 2 and Table S2) revealed a centralized structure with particularly important roles for certain populations. The degree distribution,  $P(k)$ , i.e., the proportion of nodes with  $k$  connections to other nodes, decays rapidly for large  $k$  (Fig. 3A) and the 6 highest values are all observed in samples collected in the Balearic Islands (Fig. 2 and Table S2). The average clustering coefficient of  $\langle C \rangle = 0.4$  was significantly higher than that obtained in the corresponding randomized networks ( $\langle C_0 \rangle = 0.13$  with  $\sigma_0 = 0.05$  after 10,000 realizations), whereas the local clustering decays as a function of the degree  $k$  (Fig. 3B), which indicates that the central core is substructured into a small set of hubs, with high connectivity and low clustering, linking groups of closely connected nodes (i.e., with high clustering). Examination of the relationship between the degree of a node and the average degree of the populations connected to it showed an abundance of links between highly connected and poorly connected nodes (Fig. 3C), a property termed dissortativity, present in many biological networks (25), and confirms again a centralized topology. Observation of Fig. 2 indicates that seagrass populations along the Spanish continental coasts are genetically closer to Balearic populations than to geographically closer populations. The highest values of betweenness centrality (Table S2) are also attained at the Balearic populations, suggesting that the meadows of this region play or have played a central role in relaying gene flow at the scale of the Spanish coasts. Moreover, the betweenness centrality increases exponentially with the

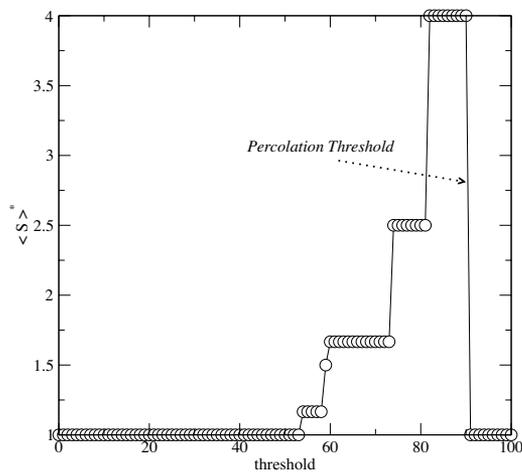


**Fig. 3.** Main topological properties found by analyzing the structure of the network of meadows at the Spanish basin scale (Fig. 2). (A) The complementary cumulative degree distribution  $P(\text{degree} > k)$ . (B) The local clustering  $C(k)$ . (C) The average degree  $\langle k_{nn}(k) \rangle$  in the neighborhood of a meadow with degree  $k$ . (D) The degree-dependent betweenness,  $bc(k)$ , as a function of the connectivity degree  $k$ .

connectivity degree  $k$  (Fig. 3D). The combination of all these findings implies a star-like structure where hubs are connected in cascade and the central core is the set of Balearic populations. A clear but more constrained perspective of this pattern is partly shown by the resulting minimum spanning tree (MST) of populations (Fig. 4B; see *Materials and Methods*), which, when analyzed with the network index of betweenness centrality highlights 3 of the major hubs encountered on the network. Yet, some other populations identified on the network appear as poorly connected on the MST, as a consequence of being a more constrained method, which finds the minimal paths required to maintain connectivity but not all of the important ones. The importance of nonminimal paths was underlined in ref. 15 in the context of geographical connectivity. This emphasizes again the advantage of the network illustration and analysis. The biological implication of these results is a great centrality of the Balearic Islands, acting or having acted as a hub for gene flow through the system.

Populations with high degree  $k$  might either be sources sustaining the system (i.e., spreading propagules), or sinks receiving gene flow from all of the other populations, or both. The extremely low rate of sexual recruitment inferred in populations with low clonal diversity ( $R$ , see *Materials and Methods*) renders those, if highly connected, much more likely to disperse than to receive. The presence in the Balearic Islands of the 2 populations with the lowest observed clonal diversity and the highest connectivity (Es Port,  $R = 0.1$ ;  $k = 10$ ; and Fornells  $R = 0.1$ ;  $k = 15$ ), likely representing populations supplying “genetic material” to neighbor populations, suggests again that the Balearic islands are a key region for the dynamics and connectivity of the metapopulation system at the scale of the Spanish coast. Furthermore, 8 among 16 continental populations show extremely low connectivity ( $k = 0$ ), thereby allowing identification of those least likely to be rescued by other populations if threatened. As in any genetic approach to metapopulation management, the role of currently observed connectivity in





**Fig. 5.** The average cluster size excluding the largest one, as a function of the imposed genetic threshold, at the whole Mediterranean scale. This identifies  $D_p = 91$  as the percolation threshold.

decreasing order, until the network reaches the percolation point (29), beyond which it loses its integrity and fragments into small clusters. This means that gene flow across the whole system is disabled if connections at a distance smaller than this critical one,  $D_p$ , are removed. The precise location of this percolation point is made with the standard methodology adequate for finite systems (29), i.e., by calculating the average size of the clusters excluding the largest one,

$$\langle S \rangle^* = \frac{1}{N} \sum_{s < S_{max}} s^2 n_s, \quad [1]$$

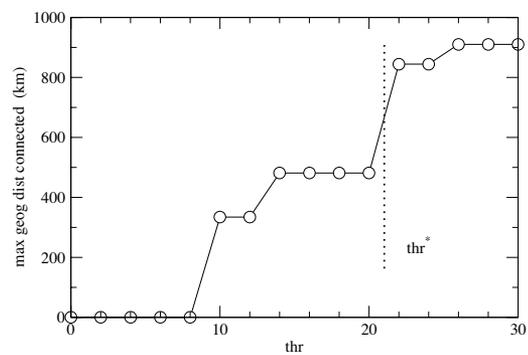
as a function of the last distance value removed,  $thr$ , and identifying the critical distance with the one at which  $\langle S \rangle^*$  has a maximum.  $N$  is the total number of nodes not included in the largest cluster and  $n_s$  is the number of clusters containing  $s$  nodes. Here, we find  $D_p = 91$ , as shown in Fig. 5.

Once the network at percolation point is obtained, we analyzed its topology and characteristics (See Fig. 1 and Table S1) and interpret those biologically. The first column in Table S1 contains also the estimated clonal diversity  $R$  of the different populations, defined as the proportion of different genotypes found with respect to the total number of collected shoots.

At the Spanish coasts scale, no percolation point is found by using the above procedure, meaning that the genetic structure in this area is rather different from the one at the whole-Mediterranean scale. To construct a useful network representation of the meadows' genetic similarity, the following alternative process was applied to determine a relevant distance threshold,  $thr$ , above which links are discarded (see Movie S1 with the network of Spanish basin at a full sequence of thresholds). At a very low threshold ( $thr = 16$ ), only the inner part of a central core, constituted by some meadows from the Balearic Islands, is connected. As the threshold is increased, new meadows (from the central Spanish coast) become connected ( $thr = 20$ ). Beyond that value, more peripheral meadows are connected from the northern and southern Spanish coasts. The geographical extension of the connected cluster (Fig. 6) grows with the distance threshold and an important jump occurs at  $thr = 22$ , when the northern and southern coasts get connected for the first time. Further distance-threshold increase does not contribute to geographical extension. Therefore, we find the value  $thr = 22$  and the resulting network as appropriate for topological characterization, because at this point the network contains a rich mixture of strong and weak links spanning all of the available geographic scales within the Mediterranean Spanish coasts. The determination of this "geographic threshold" is similar in spirit to the determination in ref. 30 of plateaus and discontinuities in network descriptors to identify relevant spatial scales.

**Estimates of Global and Local Properties of the Network.** The degree  $k_i$  of a given node  $i$  is the number of other nodes linked to it (i.e., the number of neighbor nodes). The distribution  $P(k)$  gives the proportion of nodes in the network having degree  $k$ .

We denote by  $E_i$  the number of links existing among the neighbors of node  $i$ . This quantity takes values between 0 and  $E_i^{(max)} = k_i(k_i - 1)/2$ , which is the



**Fig. 6.** The maximal geographic distance connected (at the Spanish coasts scale) as a function of the imposed distance threshold ( $thr$ ). Above  $thr = 22$  the maximal geographic distance covered by connected populations nearly duplicates, and this value—a "geographic threshold"—is chosen to construct the corresponding network.

case of a fully connected neighborhood. The clustering coefficient  $C_i$  of node  $i$  is defined as:

$$C_i = \frac{E_i}{E_i^{(max)}} = \frac{2E_i}{k_i(k_i - 1)}. \quad [2]$$

The clustering coefficient of the whole network ( $\langle C \rangle$ ) is defined as the average of all individual clustering coefficients in the system. The degree dependent clustering  $C(k)$  is obtained after averaging  $C_i$  for nodes with degree  $k$ .

Real networks exhibit correlations among their nodes (25, 31–36) that play an important role in the characterization of the network topology. Those node correlations are, furthermore, essential to understand the dynamical aspects such as spreading of information or their robustness against targeted or random removal of their elements. In social networks, nodes having many connections tend to be connected with other highly connected nodes. This characteristic is usually referred to as assortativity, or assortative mixing. On the other hand, technological and biological networks show rather the property that nodes having high degrees are preferably connected with nodes having low degrees, a property referred to as disassortativity. Assortativity is usually studied by determining the properties of the average degree ( $k_{nn}$ ) of neighbors of a node as a function of its degree  $k$  (25, 35, 37). If this function is increasing, the network is assortative, because it shows that nodes of high-degree connect, on average, to nodes of high degree. Alternatively, if the function is decreasing, as in our present case, the network is disassortative, as nodes of high degree tend to connect to nodes of lower degree. In this last case, the nodes with high degree are therefore central hubs ensuring the connection of the whole system.

The betweenness centrality (38) of node  $i$ ,  $bc(i)$ , counts the fraction of shortest paths between pairs of nodes that pass through node  $i$ . Let  $\sigma_{st}$  denote the number of shortest paths connecting nodes  $s$  and  $t$  and  $\sigma_{st}(i)$  the number of those passing through the node  $i$ . Then,

$$bc(i) = \sum_{s \neq t \neq i} \frac{\sigma_{st}(i)}{\sigma_{st}}. \quad [3]$$

The degree-dependent betweenness,  $bc(k)$ , is the average betweenness value of nodes having degree  $k$ .

**Minimum Spanning Tree.** Given a connected, undirected graph, a spanning tree of that graph is a subgraph without cycles that connects all of the vertices together. A single graph can have many different spanning trees. Provided each edge is labeled with a cost (in our analysis the genetic distance among the connected populations) each spanning tree can be characterized by the sum of the cost of its edges. A minimum spanning tree is then a spanning tree with minimal total cost. A minimum spanning tree is in fact the minimum-cost subgraph connecting all vertices, because subgraphs containing cycles necessarily have more total cost. Fig. 4 shows the minimum spanning tree for the Spanish meadows. The star-like structure centered at Balearic populations is evident, although the restriction of being a tree prevents some of the well-connected populations of the network approach to be identified here.

