

# Complexity in human transportation networks: a comparative analysis of worldwide air transportation and global cargo-ship movements

O. Woolley-Meza<sup>1</sup>, C. Thiemann<sup>1,2</sup>, D. Grady<sup>1</sup>, J.J. Lee<sup>1</sup>, H. Seebens<sup>3</sup>, B. Blasius<sup>3</sup>, and D. Brockmann<sup>4,a</sup>

<sup>1</sup> Department of Engineering Sciences and Applied Mathematics, Northwestern University, Evanston, Illinois, USA

<sup>2</sup> Max-Planck-Institut für Dynamik und Selbstorganisation, Göttingen, Germany

<sup>3</sup> ICBM, University of Oldenburg, 26111 Oldenburg, Germany

<sup>4</sup> Department of Engineering Sciences and Applied Mathematics & Northwestern Institute on Complex Systems, Northwestern University, Evanston, Illinois, USA

Received 17 March 2011 / Received in final form 28 September 2011

Published online 8 December 2011 – © EDP Sciences, Società Italiana di Fisica, Springer-Verlag 2011

**Abstract.** We present a comparative network-theoretic analysis of the two largest global transportation networks: the worldwide air-transportation network (WAN) and the global cargo-ship network (GCSN). We show that both networks exhibit surprising statistical similarities despite significant differences in topology and connectivity. Both networks exhibit a discontinuity in node and link betweenness distributions which implies that these networks naturally segregate into two different classes of nodes and links. We introduce a technique based on effective distances, shortest paths and shortest path trees for strongly weighted symmetric networks and show that in a shortest path tree representation the most significant features of both networks can be readily seen. We show that effective shortest path distance, unlike conventional geographic distance measures, strongly correlates with node centrality measures. Using the new technique we show that network resilience can be investigated more precisely than with contemporary techniques that are based on percolation theory. We extract a functional relationship between node characteristics and resilience to network disruption. Finally we discuss the results, their implications and conclude that dynamic processes that evolve on both networks are expected to share universal dynamic characteristics.

## 1 Introduction

Large-scale human transportation networks are essential for global travel, international trade, the facilitation of international partnerships and relations, and the advancement of science and commerce. The worldwide air transportation network supports the traffic of over three billion passengers travelling between more than 4000 airports on more than 50 million flights in a year [1]. The global cargo-ship network accounts for up to 90% of the international exchange of goods; approximately 60 000 cargo-ships are connecting more than 5000 ports worldwide with about a million ship movements every year [2,3]. These two networks constitute the operational backbone of our globalized economy and society.

Although they are immensely important for facilitating exchange between geographically distant regions, the ever-increasing amount of traffic over such complex, densely-connected transportation networks introduces serious problems. Rising energy costs, pollution, and global warming are obvious concerns, but globalized traffic also plays a key role in the worldwide dissemination of infectious diseases and invasive species [4–19]. The first decade

of the 21st century has witnessed the emergence and worldwide spread of two major global epidemics: the severe acute respiratory syndrome (SARS) in 2003 [20–23], and the recent H1N1 pandemic of 2009 [24,25]. Both diseases rapidly spread across the globe in a matter of weeks to months, a process linked directly to long-distance traffic routes over which infected individuals dispersed infectious agents. In combination with increasing worldwide population size, which is expected to pass the 7 billion threshold within the next decade, and the concentration of the majority of the world's population in mega-cities and urban areas [26], the impact of global pandemic events is expected to become one of the most challenging problems of the 21st century. The spread of invasive species into new habitats and ecosystems presents a similar and equally-significant problem [17,27–29]. The largest vector of marine bioinvasion is assumed to be global shipping [30]. Human-mediated bioinvasion has become one of the key factors in the global biodiversity crisis [31,32] and may affect the stability of ecosystems, survival of species, and human health [11,27]. The introduction of invasive species to foreign ecosystems has generated annual costs of over \$120 billion in the United States alone [33].

<sup>a</sup> e-mail: brockmann@northwestern.edu

In addition to introducing environmental problems, the complex transportation web itself is subject to external disruptions. For instance, the unexpected eruption of the Icelandic volcano Eyjafjallajökull in 2010 and subsequent closure of major European airports led to a major disruption in global traffic and significant economic stress over a period of several days. Other influences, such as meteorological events like hurricanes, the recent rise in acts of piracy in Somalia, or the financial crisis of 2007 also make flexibility in global cargo traffic necessary and underline the vulnerability of international trade and transportation systems. It is therefore of fundamental importance to understand the resilience of these networks in response to regional and large-scale failure of parts of the network, and to identify “sensitive” regions of the network. This point becomes even more important in light of malicious terrorist activities.

A deep understanding of the structure of human transportation networks will lead to new insights into the geographical spread of diseases and invasive species, allow the development of new computational models for their time courses, and eventually allow us to predict their impact on our environment and society. Computational techniques for investigating the resilience of these networks in the face of partial failure will play a fundamental part in achieving this understanding; complex network theory [34] already provides a powerful theoretical tool in this respect. But, although both the worldwide air transportation network and the global cargo-ship network have already been subjected to a number of network-theoretic analyses [6,9,21,35,36], it is still unclear whether the observed properties of these networks are unique to a specific context, or are universal and generic. A lack of comparative studies in this direction, and indeed a lack of data, has led to a scarcity of universal theories of the structure of transportation networks.

Here we address this issue using a comparative approach. We analyse and compare the structure of the worldwide air-transportation and global cargo-ship networks (WAN and GCSN in the following) and show that a surprising number of properties are shared by both networks despite their different use, economic context, scale, and connectivity structures. The analysis suggests that the same fundamental principles guide the growth of both networks. Most importantly, it suggests that dynamic processes that evolve on these networks will exhibit similar dynamic features, an important insight since it implies that processes as different as emergent human infectious diseases and human-mediated bioinvasion can be investigated along the same line of research.

The dynamics of processes on networks are guided not only by the topology of the network, but also by the interaction strengths between pairs of nodes. One of the characteristic features of transportation networks is a strongly heterogeneous distribution of interaction strength. Among other things, this implies that the shortest topological path between two nodes may not be the path of strongest interaction, and we account for such effects by using the idea of effective shortest paths. These are analogous to

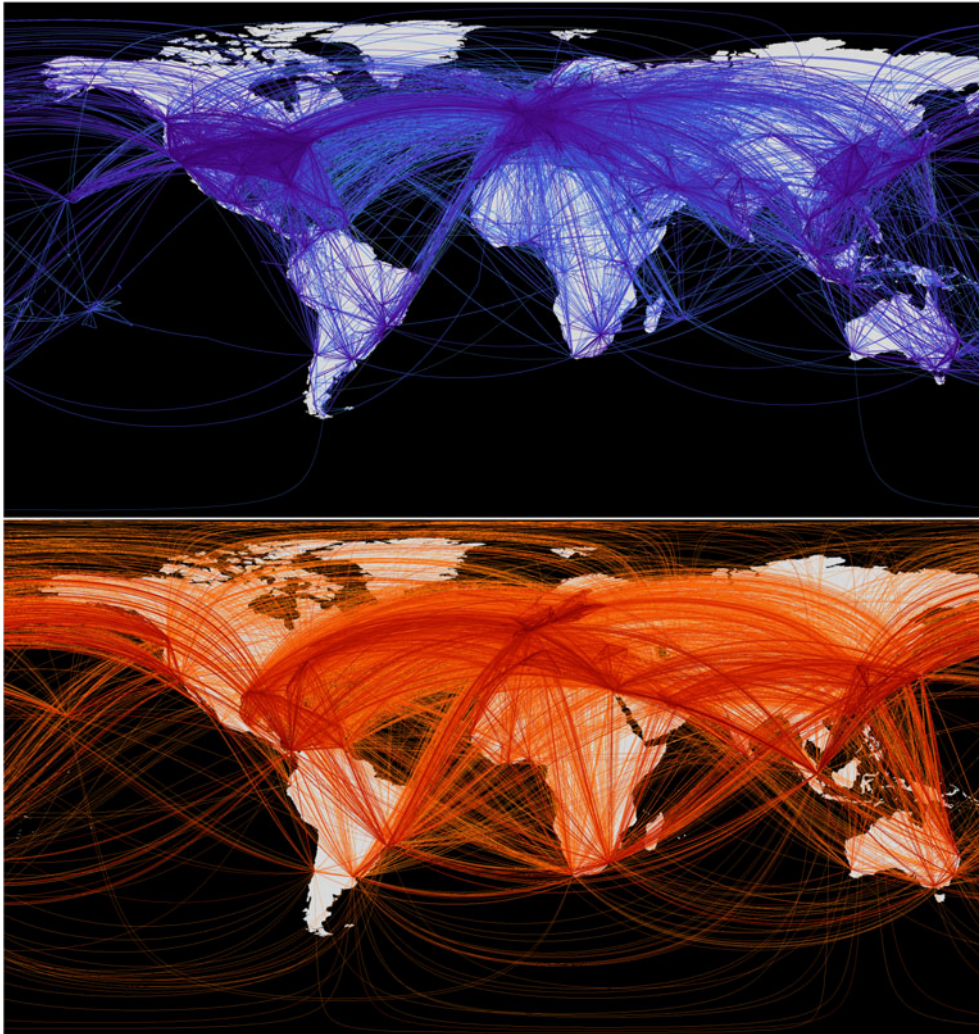
the well-known topological shortest paths, except that the length of an edge is taken to be the reciprocal of the weight of that edge, and the effective shortest path is then the path that minimizes the total effective distance. This approach also reveals surprising similarities between the two networks.

The paper is structured as follows: In Sections 2 and 3 we introduce the WAN and GCSN and discuss their statistical properties and similarities. In Section 4 we apply a recently developed technique based on shortest path trees to compute structural properties of these networks. In Section 5 we use the measures introduced to investigate and compare resilience of these networks in response to targeted attacks and random failure. We discuss the implications of our results in Section 6.

## 2 The worldwide air transportation and cargo-ship networks

The WAN and GCSN are infrastructure systems on which we travel and transport commodities on a worldwide scale. Complex network theory [34,37,38] provides the most plausible quantitative description of these systems: pairs of nodes  $i$  and  $j$  are connected by links with non-negative weights  $w_{ij} > 0$  if transport occurs directly between these nodes, and  $w_{ij} = 0$  if they are not directly connected. It is generally possible in transportation networks to begin at any node and locate a path to any other node, but the  $w_{ij}$  measure only direct connections and quantify the magnitude of traffic between pairs of nodes. In the WAN nodes represent airports and the weight matrix could be defined as the total number of passengers per unit time, the number of passenger planes or the number of scheduled flights. In the GCSN nodes represent ports, and  $w_{ij}$  could quantify the number of cargo-ships or the net tonnage of cargo per unit time. For a comparative analysis we choose  $w_{ij}$  to represent the number of carrier vehicles (passenger planes or cargo-ships) that travel from node  $j$  to  $i$  per unit time. For the WAN,  $w_{ij}$  is the average number of scheduled commercial flights per year between airports in the three-year period 2004–2006, as reported by OAG Worldwide Ltd. [1]. The GCSN was established from data of worldwide ship movements provided by IHS Fairplay [3]. This information was used to reconstruct the journey of 15,415 ships travelling around the globe during 2007 [9]. The GCSN is restricted to vessels bigger than 10 000 gross tonnages which accounts for 86% of the world fleet. For both WAN and GCSN, the resulting weight matrix  $W$  is virtually symmetric, i.e.  $w_{ij} \approx w_{ji}$  for all pairs  $ij$ . In order to guarantee strict symmetry we symmetrized the matrix according to  $w_{ij} \rightarrow (w_{ij} + w_{ji})/2$ . This differs from the definition that was used in [9], where the number of directed links was reported.

Both networks are depicted in Figure 1. Despite their global coverage and structural similarity, these networks exhibit distinct features. The WAN comprises approximately five times as many nodes but almost the same number of links, yielding a less-densely-connected network



**Fig. 1.** (Color online) Global transportation and mobility: the top panel illustrates the worldwide air-transportation network (WAN) consisting of 4069 airports connected by 25 453 links that represent the number of passenger planes travelling between nodes per unit time. The bottom panel depicts the global cargo-ship network (GCSN) that connects 951 international ports along 25 819 routes. In both panels, the saturation of the lines indicates the total flux along a route with darker lines carrying more flux.

**Table 1.** Network characteristics of WAN and GCSN. Number of nodes  $N$ , number of links  $L$ , network connectivity  $\sigma \approx 2L/N^2$ , network diameters  $\phi$  and  $d_T$ , network clustering coefficient  $c$  and network length scale  $\langle r \rangle$  (in units km), total traffic  $C$  (vehicles/yr.) reflect global properties of the network. The table also lists the mean link weights  $\langle w \rangle$ , traffic per node (flux)  $\langle F \rangle$  (both in units of vehicles/yr.) and node degree  $\langle k \rangle$  and corresponding coefficients of variation.

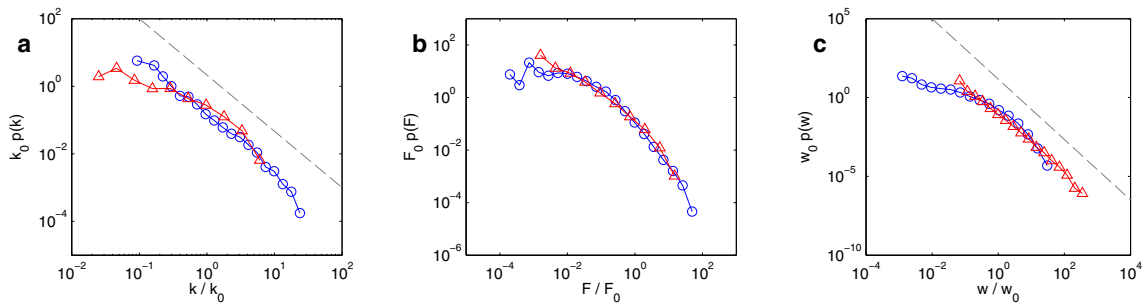
	$N$	$L$	$\sigma$	$\phi$	$d_T$	$c$	$\langle r \rangle$	$C$	$\langle w \rangle$	$\langle F \rangle$	$\langle k \rangle$	$CV(w)$	$CV(F)$	$CV(k)$
WAN	4069	25 453	$3.07 \times 10^{-3}$	29.11	4.16	0.55	1109	$5.68 \times 10^7$	1116.49	$1.39 \times 10^4$	12.51	2.03	3.60	2.15
GCSN	951	25 819	$5.72 \times 10^{-2}$	13.57	2.34	0.57	1857	$9.87 \times 10^5$	19.11	$1.04 \times 10^3$	54.30	7.27	2.34	1.22

( $\sigma_{\text{WAN}} = 3.07 \times 10^{-3}$  as compared to  $\sigma_{\text{GCSN}} = 5.72 \times 10^{-2}$ , see Tab. 1). As the the GCSN is restricted to large vessels, the number of seaports in the world may be an order of magnitude higher.

Node flux and degree are key characteristics, defined according to

$$F_i = \sum_j w_{ij} \quad \text{and} \quad k_i = \sum_j a_{ij}, \quad (1)$$

where  $a_{ij}$  are elements of the adjacency matrix, that is  $a_{ij} = 1$  if nodes  $i$  and  $j$  are connected and  $a_{ij} = 0$  otherwise. On average a node in the WAN dispatches  $\langle F \rangle_{\text{WAN}} = 1.39 \times 10^4$  vehicles per year. The average number of cargo-ships leaving a port in the GCSN is  $\langle F \rangle_{\text{GCSN}} = 1.03 \times 10^3$ . Higher connectivity of the GCSN is also reflected in the mean degree,  $\langle k \rangle_{\text{WAN}} = 12.51$  and  $\langle k \rangle_{\text{GCSN}} = 54.30$ .



**Fig. 2.** (Color online) Statistical properties of WAN (blue) and GCSN (red): panels a, b, c depict the probability density functions of node degree  $k$ , node flux  $F$  and link weight  $w$ . Dashed lines are power laws with exponents 1.5 (panel a) and 2 (panel c). The abscissal scaling factors  $k_0, F_0, w_0$  are given by the mean in each distribution. Despite differences in  $p(k)$  both networks' degrees are broadly distributed, the WAN more closely resembling a power law than the GCSN. The flux distributions  $p(F)$  are almost identical with deviations in the small flux regime. The weight distributions range over 3.8 and 3.5 orders of magnitude and for intermediate values exhibit approximately the same decay.

The typical traffic per link is given by the mean link weight  $\langle w \rangle$ , and although in the WAN it exceeds the GCSN by two orders of magnitude, the variability reflected in the coefficient of variation is significantly higher in the GCSN. The clustering coefficient  $c$  indicates the abundance of triangular motifs in the network, and in spite of the GCSN's higher connectivity the clustering coefficient is nearly identical in both networks. This indicates that both networks can be considered sparse and saturation effects are not significant.

A typical length scale of the network can be defined by

$$\langle r \rangle = \sum_{j,i} p_{ji} r_{ji} \quad (2)$$

where  $p_{ji} = w_{ji}/F_j$  and  $r_{ji}$  the geographical distance between nodes  $j$  and  $i$ . The quantity  $p_{ji}$  is the relative fraction of traffic from  $i$  to  $j$  with respect to the entire traffic through node  $j$ . Thus  $\langle r \rangle$  represents the mean distance traveled by a carrier in the network. According to this definition the typical length scale of the GCSN is over one and a half times the typical length scale of the WAN (see Tab. 1). Related to the geographic distance are topological distance measures defined by the connectivity of the networks. The diameter of a network can be defined as the average shortest path length that connects a pair of nodes,  $d_T$ . For WAN and GCSN  $d_T = 4.16$  and  $2.34$ , respectively (in a fully connected network  $d_T = 1$ ).

### 3 Universal statistics in large-scale transportation networks

A key feature that many large-scale technological networks share is their strong structural heterogeneity in terms of link and node statistics and centrality measures. These networks typically contain a small fraction of hubs characterized by strong connectivity and high centrality scores complemented by a large number of smaller nodes that connect to the hubs. This structural property is captured in several models with scale-free distributions of centrality measures [36,39].

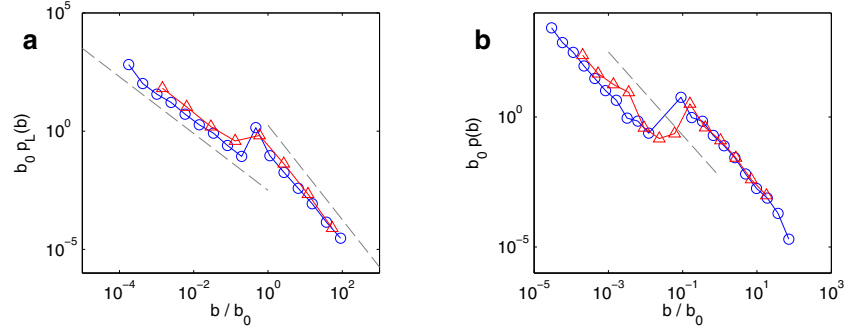
Figure 2 shows that  $w$ ,  $k$ , and  $F$  have qualitatively the same distribution in both networks, ranging over many orders magnitude and differing only by a scaling factor. Their surprisingly similar shape, in particular in  $p(F)$ , supports the claim that these networks have evolved according to similar fundamental processes. It has been pointed out [35,38,40,41] that degree, flux, and weight approximately follow power laws in many networks. However, in spite of the similarity of these distributions and the suggestive scaling behavior, we do not find a reasonable power-law fit over more than one decade, suggesting that simple models for generating scale-free statistics are not sufficient to describe these empirical networks.

#### 3.1 Weighted betweenness centrality of links and nodes

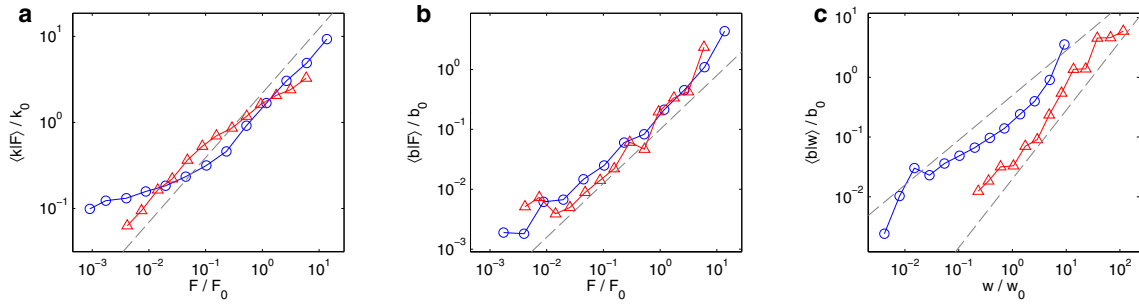
Another commonly investigated measure for link and node centrality is betweenness centrality. The betweenness  $b$  of a link (or a node) is the fraction of shortest paths in the entire network of which the link (or node) is part of. Betweenness requires the definition of length of a path which in turn requires the definition of length of a link. In weighted networks a plausible choice for the effective length of a link connecting nodes  $i$  and  $j$  is given by the proximity  $\lambda_{ij}$  defined by

$$\lambda_{ij} = \frac{\langle w \rangle}{w_{ij}}. \quad (3)$$

This definition accounts for the notion that strongly connected nodes are effectively more proximate than nodes that are weakly coupled. The numerator  $\langle w \rangle$  sets the typical distance scale  $\lambda_0 = 1/\langle w \rangle$  and  $\lambda_{ij}$  is defined relative to it. Based on this effective proximity one can define the length of a path  $P(i_0, \dots, i_k)$  that starts at node  $i_0$  and terminates at node  $i_k$  connecting a sequence of intermediate nodes  $i_n$ ,  $n = 1, \dots, k-1$  along direct connections of weights  $w_{i_n i_{n+1}}$  by summing of the proximities of each leg in the path. This integrated distance  $l(i_0, \dots, i_k)$  is given



**Fig. 3.** (Color online) Structure of betweenness centrality of WAN (blue) and GCSN (red): panels a and b depict the distributions  $p(b)$  of betweenness of links and nodes, respectively. Link betweenness exhibits two scaling regimes with exponents  $\gamma_2 \approx 1.2$  ( $b \ll b_0$ ) and  $\gamma_2 \approx 2.0$  ( $b \gg b_0$ ). Node betweenness also exhibits two distinct regimes but they are characterized by the same scaling exponent  $\gamma_1 \approx 1.6$ . The abscissal scaling factors  $b_0$  are given by the mean in each distribution.



**Fig. 4.** (Color online) Correlation structure of centrality measures. (a) The conditional mean degree  $\langle k|F \rangle$  as a function of flux  $F$ . The dashed line represents sub-linear scaling with exponent 0.75. Neither networks exhibit a clear scaling, the WAN approaches the scaling line for large arguments whereas the GCSN for small arguments. The overall behavior can be roughly described by an algebraic scaling relation indicated by the dashed line. (b) Conditional mean node betweenness  $\langle b|F \rangle$  as a function of  $F$  with dashed line with scaling exponent of 0.9. Despite deviation for small and large arguments the data for both networks are reasonably well described by the indicated scaling for intermediate arguments across approx. 2 orders of magnitude. (c) Conditional weight betweenness  $\langle b|w \rangle$  as a function of  $w$ . The WAN exhibits sub-linear scaling contrary to the GCSN for which super-linear scaling is observed (roughly following dashed lines with scaling exponent 0.75 and 1.15 respectively)

by the sum

$$l(i_0, \dots, i_k) = \sum_{n=0}^{k-1} \lambda_{i_n i_{n+1}} = \sum_{n=0}^{k-1} \frac{\langle w \rangle}{w_{i_n i_{n+1}}}. \quad (4)$$

For a given pair of nodes  $i_0$  and  $i_k$ , many paths exist that connect these nodes along intermediate nodes  $i_n = 1, \dots, k-1$ . Using the definition of length of a path above, the shortest path between two nodes is defined as one with minimal  $l$ :

$$d(i_0, i_k) = \min_{i_n=1, \dots, k-1} l(i_0, \dots, i_k). \quad (5)$$

We define the effective distance  $d_{ij}$  between nodes  $i$  and  $j$  as this effective length of the shortest path connecting them, i.e.  $d_{ij} = d(i, j)$ , and denote the unique path associated with it by  $P_s(i, j)$ . Based on this definition we define the diameter  $\phi$  of the network as the mean shortest path length over the ensemble of all pairs of nodes. According to this definition the WAN's diameter is slightly more than twice the diameter of the GCSN (see Tab. 1). The reasons for this will be discussed in more detail in Section 4.

We computed betweenness centrality  $b$  for both links and nodes based on the set of all shortest paths  $P(i, j)$ . Figure 3 depicts the distributions  $p(b)$  for both networks. Unlike the centrality measures of degree and flux for nodes and weights for links, the distribution of betweenness exhibits a well-pronounced discontinuity in both networks. Note that this discontinuity is absent from  $p(b)$  if the network are assumed to have uniform weights as for instance in reference [42]. This indicates that in the WAN and GCSN links and nodes segregate into two distinct functional groups. In fact the point  $b_c$  at which the discontinuity occurs can be employed to separate links and nodes that belong to the operational backbone of the network [43]. Both networks exhibit similar scaling behavior in the two betweenness regimes.

### 3.2 Correlations in centrality measures

Degree, flux and betweenness typically exhibit positive correlations and scaling relationships with one another. For instance, recently-investigated mobility networks [9,35,41] exhibit a sub-linear scaling relation  $k \sim F^\eta$  with exponent  $\eta \approx 0.58$  and  $\eta \approx 0.7$ . Figure 4 compiles

scaling relationships we observe in the WAN and GCSN. To extract the scaling relationship we computed the mean of one centrality measure  $x$  conditioned on a second centrality measure  $y$ , that is,

$$\langle x|y \rangle = \frac{\int dx x p(x, y)}{\int dx p(x, y)} \quad (6)$$

where  $p(x, y)$  is the combined distribution of both. Our analysis shows that both networks exhibit a sub-linear correlation of degree with flux

$$\langle k|F \rangle \sim F^\eta \quad (7)$$

with approximately identical exponent  $\eta = 0.75$  for both networks and across four orders of magnitude of  $F$ . This is consistent with previous findings and the intuitive notion that node connectivity increases with traffic. A sub-linear scaling of degree with flux implies that the typical weight  $\langle w|F \rangle$  of links connected to nodes of size  $F$  scales according to

$$\langle w|F \rangle \sim F / \langle k|F \rangle \sim F^{1-\eta}. \quad (8)$$

Since  $\eta < 1$  this implies that high flux nodes typically connect to other nodes with stronger links, as expected for transportation networks. The fact that  $\eta$  is almost identical in both networks is additional evidence that similar universal mechanisms are responsible for shaping the topological structure of both the WAN and GCSN. Similarly, node betweenness scales as

$$\langle b|F \rangle \sim F^\zeta \quad (9)$$

with an exponent  $\zeta \approx 1$  in both networks. A linear relationship between node flux and betweenness can be explained by the heuristic argument that typical betweenness values of a node increase linearly with its degree  $k$ . Likewise, since shortest paths are computed based on link weights, it is reasonable to assume that node betweenness scales linearly with the typical link weight of a node and thus

$$\langle b|F \rangle \sim \langle k|F \rangle \langle w|F \rangle \sim F^\eta F^{1-\eta} = F \quad (10)$$

and hence one expects  $\zeta \approx 1$  as observed. Conditional mean of link betweenness as a function of link weight  $\langle b|w \rangle$  exhibit approximate scaling. Figure 4c suggests sub-linear scaling for the WAN as opposed to super-linear scaling for the GCSN. This difference in scaling in both networks is the first marked difference that we observe in the statistics of centrality measures. Possible explanations are the differences in overall connectivity  $\sigma$  in both networks (see Tab. 1) and that there is a statistically-significant difference in the way that weights are distributed among the nodes.

## 4 Network shortest path trees

Global properties of strongly heterogeneous, multi-scale networks, such as connectivity, average clustering coefficient, and diameter, as well as statistical distributions of

centrality measures, provide important insight and may serve as quantitative classifiers for networks. However, they cannot resolve properties and structures on a local scale. On the other hand, local measures such as a node's individual degree or mean link weight of its connections provide local information only and cannot capture global properties. Transportation networks exhibit important structure on intermediate scales, so it is vital to understand structural properties that are neither local nor global in these networks. One way to approach this is to analyse and investigate the structure of the entire network from the perspective of a chosen node. Clearly, geographic distance is an important parameter in this context as operation costs typically scale with geographic distance. However, in complex multi-scale transportation networks such as the WAN and the GCSN, geographic distance is rarely a good indicator of the effective distance of connected nodes. High-flux hubs in each network are typically connected by strong traffic bonds even across very large geographic distances while smaller-flux nodes can be connected by weak links although they may be geographically close. A spatial representation as depicted in Figure 1 is therefore a misleading way to convey effective distances in these networks.

An alternative representation can be obtained based on the notion of proximity defined by equation (3) and effective shortest paths, equation (4). Based on this notion we compute the shortest paths of a chosen root node  $i$  to all other nodes  $j$ . The collection of links contributing to these paths form a shortest path tree  $\mathbf{T}_i$  rooted at  $i$ . Spatial representations of such trees are depicted in Figure 5 for each network and two different root nodes. The radial distance in these figures represents the effective, shortest path distance  $d_{ij}$ . The lines represent the connections of  $\mathbf{T}_i$ . Note that, although the trees differ in both networks and for different root nodes, high-centrality nodes tend to exhibit the smallest effective (shortest path) distance to the root node. Note also that the geometry of the networks exhibits significant structural differences in both networks: in the WAN the spatial distribution in the new representation is less regular and the scatter in effective distance is larger than in the GCSN where nodes reside in a well defined annular region.

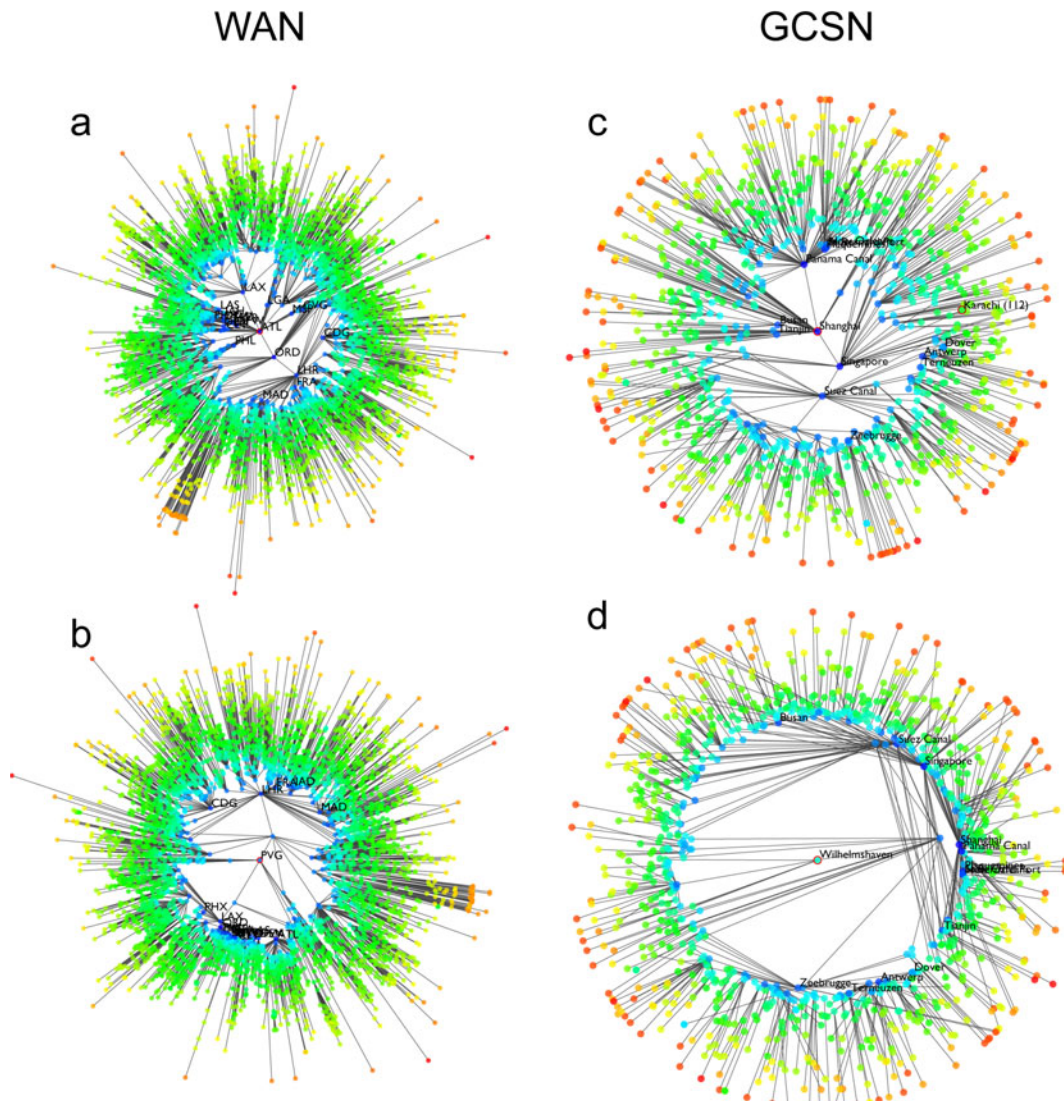
In order to understand these qualitative differences and similarities we investigate the distribution of the shortest path distances conditioned on the type of root node (Fig. 6). Conditioned on the flux of the root node, we compute the distribution of shortest path distance, that is,  $p(d|F)$ . Based on this distribution we determine the expected distance of the network from a node with specified flux as

$$\mu_d(F) = \langle d|F \rangle \quad (11)$$

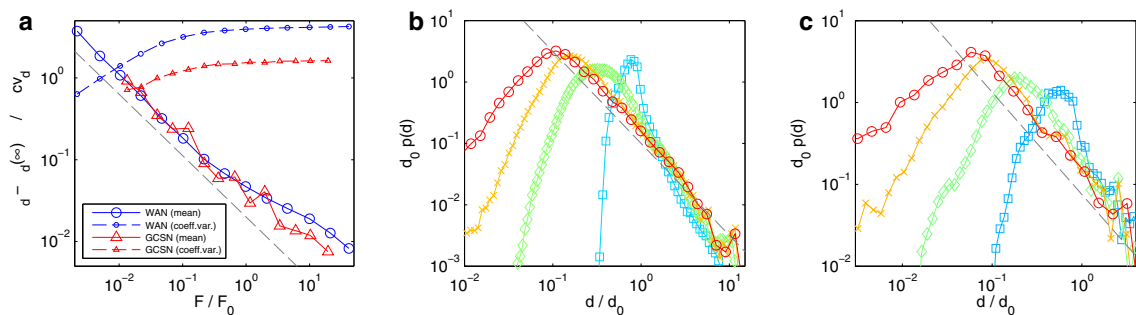
as well as the conditional coefficient of variation:

$$cv_d(F) = \frac{\sqrt{\langle d^2|F \rangle - \langle d|F \rangle^2}}{\langle d|F \rangle}.$$

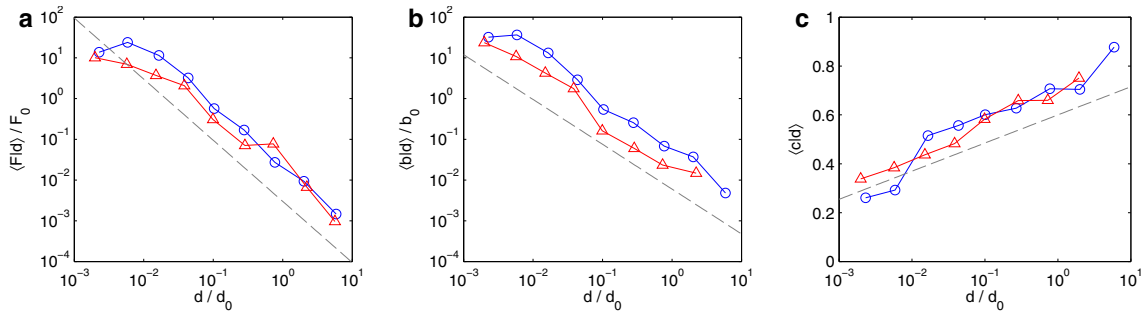
The quantity  $\mu_d(F)$  measures the typical distance from a root node with flux  $F$  to the rest of the network. The



**Fig. 5.** (Color online) Shortest path tree structures and effective distances in WAN and GCSN. Left: the panels depict the shortest path trees of airports ATL (Atlanta) and PVG (Shanghai). The radial distance of the remaining airports with respect to these root nodes represent the logarithm of the shortest path distance to the reference node. Right: the panels depict the GCSN shortest path tree for ports Wilhelmshaven, Germany and Shanghai. Note that the overall structure of both representations is different, yet both networks share the feature of circular arrangement according to node flux, encoded by color (blue represents large flux nodes and orange small node flux). Note that irrespective of the chosen root node, the closest nodes in terms of effective distance are always high flux nodes and small flux nodes are always peripheral in this representation.



**Fig. 6.** (Color online) Shortest path distance statistics. (a) Conditional mean shortest path (Eq. (11)) as a function of node flux as well as the conditional coefficient of variation in shortest paths. (b) The conditional distributions of shortest paths for four subtypes of root nodes (WAN) ranked according to flux, where red markers denote high flux and light blue denote low flux. (c) Same as in (b) for the GCSN. Dashed grey lines are scale free with exponents 0.75 (a), 1.5 (b) or 1.25 (c).



**Fig. 7.** (Color online) Correlation analysis of effective distance  $d$  and centrality measures for WAN (blue) and GCSN (red). (a) Node flux. (b) Weighted node betweenness. (c) The local clustering coefficient  $c$ . The dashed lines illustrate scaling exponents  $\omega_F = 1.5$  (a) and  $\omega_b = 1.1$  (b).

coefficient of variation measures the statistical variability in  $d$ . Figure 6 depicts both quantities for the WAN and GCSN. Note that  $\mu_d(F)$  is remarkably similar for both networks and can be described by

$$\mu_d(F) - \mu_d(\infty) \sim (F/F_0)^{-\tau} \quad (12)$$

with  $\tau \approx 0.75$ . Note that this relation indicates the existence of a lower limit to the typical effective distance for increasing node flux  $\mu_d(\infty) > 0$  which implies that even extremely large hubs exhibit a least distance to the rest of the network. Equation (12) implies that mean effective distance decreases in a systematic way with node centrality and according to the same relation in both networks. However, the coefficient of variation increases monotonically with  $F$ , which implies that the variability in effective distance increases with the centrality of the root node. This can also be observed in Figures 6b and 6c, which depicts the entire distribution  $p(d|F)$  for four categories of root nodes of different centrality. For most central nodes  $p(d|F)$  increases steeply for small values of  $d$  and exhibits an algebraic decay for large distance. As  $F$  decreases,  $p(d|F)$  attains a sharper peak as small distances disappear from the distribution. This qualitative behavior is observed in both networks. The asymptotic behavior for large effective distances is approximately

$$p(d|F) \sim d^{-\theta}$$

with  $\theta \approx 1.5$  for the WAN and  $\theta \approx 1.25$  for the GCSN.

A characteristic property of the network representation in Figure 5 is that regardless of the properties of the root node, the rest of the nodes tend to sort in concentric circles (effective distances) according to centrality measures. A key question is then how effective distance correlates with centrality measures. If there is a strong correlation between effective distance and node centrality measures, this implies that centrality measures dominate the placement of a node in a network.

In order to determine the relationship between effective distance and centrality measures, we select three groups of nodes, the top 2.5% ranked according to degree, flux and betweenness, and combine them into a single high-centrality subset of nodes  $\Omega$ , which accounts for 5% of the entire network. The remaining 95% of the nodes

are denoted by  $\bar{\Omega}$ . Based on this subset we determine the distribution  $p(x, d|\Omega)$ , the probability of finding a node in  $\bar{\Omega}$  with centrality measure  $x$  (degree, flux, betweenness) and effective distance  $d$  to the root nodes in  $\Omega$ . From this we compute the conditional mean

$$\langle x|d \rangle = \int x p(x|d, \Omega) = \int x p(x, d|\Omega)/p(d|\Omega). \quad (13)$$

Figure 7 depicts  $\langle F|d \rangle$  and  $\langle b|d \rangle$  for both networks. Despite their difference, WAN and GCSN exhibit almost identical scaling relations

$$\langle F|d, \Omega \rangle \sim d^{-\omega_F} \quad \text{and} \quad \langle b|d, \Omega \rangle \sim d^{-\omega_b} \quad (14)$$

with  $\omega_F \approx 1.5$  and  $\omega_b \approx 1.1$ , consistent with the intuitive notion that centrality decreases with increasing effective distance from central root nodes. Figure 7 also shows that the local clustering coefficient as a function of  $d$  approximately scales according to

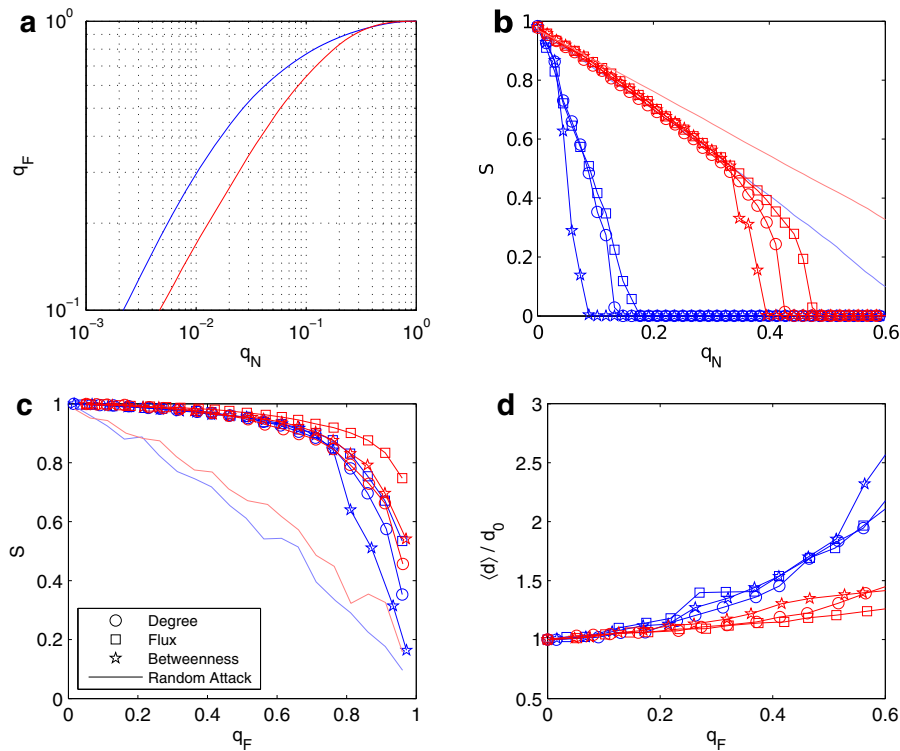
$$\langle c|d, \Omega \rangle \sim \log(d/d_0). \quad (15)$$

The logarithmic increase of the clustering coefficient implies that in their peripheral regions the WAN and GCSN become less tree-like. Although nodes in the two networks are subject to different geographic constraints and clustering may therefore have different geographic structure, they nonetheless organize themselves in a surprisingly similar fashion when effective distance is considered. Nodes from  $\bar{\Omega}$  that are connected to the root nodes in  $\Omega$  do not exhibit large fractions of connections among one another, which indicates that high-centrality root nodes function as “feed-in” hubs to low centrality nodes.

## 5 Network resilience

A key question in the context of large-scale technological and infrastructural networks concerns their response to local failure and resilience to partial accidental breakdown or anticipated attacks. Both the WAN and GCSN are subject to unpredictable, recurrent, and extreme weather conditions that lead to repetitive and regionally-localized failure that must be compensated for by re-routing traffic or re-planning schedules.





**Fig. 8.** (Color online) Resilience properties of WAN (blue) and GCSN (red) in response to random and selected node removal: (a) Fraction  $q_F$  of total traffic that is carried by a fraction  $q_N$  of top-flux nodes. In both networks a few most central nodes carry substantial amount of traffic in the network. (b) The relative size of the largest connected component of the networks as a function of the fraction  $q_N$  of removed nodes. Nodes are removed according to the rank in terms of degree  $k$  (circles), flux  $F$  (squares), betweenness  $b$  (stars), or randomly (no symbols). (c) The largest component as a function of fraction of removed traffic. (d) The response to selected node removal as reflected in network inflation. The panel depicts the diameter of the network as defined by the mean shortest path between all pairs of nodes.

Random failures and targeted attacks are typically investigated using the framework of percolation theory [44,45]. A random (failure) or selected (attack) fraction  $q$  of nodes is removed from the network and structural responses of the network are investigated as a function of  $q$ . Important insight was gained in studies that investigated random or selected node removal in random networks [44–47]. One of the most important findings of these studies was that scale-free networks with power-law degree distributions respond strikingly differently in scenarios that reflect random failures as opposed to selected removal of central nodes. For instance, scale-free networks are relatively immune to random removal of nodes and extremely sensitive to targeted removal of high centrality nodes. Since centrality measures such as degree, betweenness, and flux typically correlate in these networks, this effectively amounts to removal of nodes that function as hubs. One of the essential questions in this context addresses the critical fraction  $q$  of removed nodes that are required to disintegrate the global connectivity of the network. This critical value is the percolation threshold  $q_c$ : for  $q < q_c$  the size of the giant component (the largest subset of nodes that are connected by paths) is typically the size of the entire network. Beyond the percolation threshold ( $q > q_c$ ) the networks falls apart into a family of disconnected, fragmented sub-networks.

The resilience properties of the WAN and GCSN to sequential node removal are depicted in Figure 8. For each centrality measure (degree, betweenness, and flux), we remove fractions  $q$  of nodes randomly and also in order of descending rank with respect to  $k$ ,  $b$ , and  $F$ , respectively. We compare two different removal protocols. Since both networks are strongly inhomogeneous, removing a fraction of nodes is not equivalent to removing a fraction of traffic (see Fig. 8a). For example 1% of the most connected nodes account for 29.7% of the entire traffic in the WAN and 17.6% in the GCSN, and removal of 10% of nodes with highest flux is equivalent to reducing the total traffic in the WAN by 76.9% and in the GCSN by 64.5%. Because of this pronounced nonlinear relationship, we compare resilience of the network as a function of the fraction of removed nodes  $q_N$  as well as the fraction of removed traffic  $q_F$ .

Figure 8b depicts the relative size of the giant component  $S$  as a function of  $q_N$ . Both networks are resilient to random failures (we find that the giant component decreases linearly with the fraction of nodes removed, i.e.  $S \approx 1 - q_N$ ), although the WAN is taking some excess damage from random failures. Furthermore, we observe a percolation threshold for the targeted attacks in both networks. The WAN exhibits a percolation threshold at  $q_N^c = 13.8\%$ ,  $17.2\%$  and  $9.4\%$ , for node removal according

to degree, flux and betweenness. The thresholds are significantly larger for the GCSN at  $q_N^c = 44.3\%$ ,  $49.0\%$  and  $39.5\%$ . In each network the threshold depends only weakly on the choice of centrality measure because of the strong correlation among different centrality measures. Note however that both networks are most susceptible to removal according to betweenness rank, followed by degree and node flux. The overall higher threshold in the GCSN is caused by the greater connectivity  $\sigma$  and mean degree  $\langle k \rangle$  of the network (see Tab. 1).

Figure 8c depicts  $S$  as a function of  $q_F$ . The random failures appear to be more effective here because they remove more nodes for a given fraction of removed traffic than the targeted attacks, since not only high-centrality nodes are selected. However, due to the strong nonlinear relationship between  $q_N$  and  $q_F$  it is evident that both networks are strongly resilient to targeted attacks. Even substantial traffic reduction has virtually no impact on the relative size of the giant component, for instance when 50% of the entire traffic is reduced in both networks, the giant component is still larger than 90% of the original network and no percolation threshold is observed in the range up to 80% of traffic reduction. These percolation thresholds are rarely reached in real networks. Another approach that has been applied [46,48] is based on the topological diameter of the network and its response to network disruption. Typically when high-centrality nodes are removed from the network, the diameter of the network increases as the shortest paths connecting two arbitrary nodes lengthen due to the increasing lack of hubs that can serve as connecting junctions. Figure 8d shows that this inflation of the network in response to node removal is observed in both networks. This effect is relatively independent of the choice of centrality measure used in the removal protocol. Furthermore, the GCSN is more robust to node reduction, which we believe to be a consequence of the high connectivity of the network.

Naive percolation analysis and network inflation have only limited applicability in real world scenarios. Many refinements to percolation have been introduced [38,49] which are more sensitive to the heterogeneous distribution of traffic in the network but like inflation, they only address changes in the global structure of the network.

### 5.1 Resilience and shortest paths

The concept of shortest path trees can be used to study network perturbations in a more refined framework and well below the percolation threshold, on a node by node basis [50]. In response to removal of a fraction  $q_N$  of most central nodes or the equivalent fraction of traffic  $q_F$  in the entire network, we can compute the impact by investigating the change of shortest path trees  $\mathbf{T}_i$  for each root node  $i$ , that is, we can quantify the impact of the network disruption from the perspective of every node. To this end we define a node's impact factor as

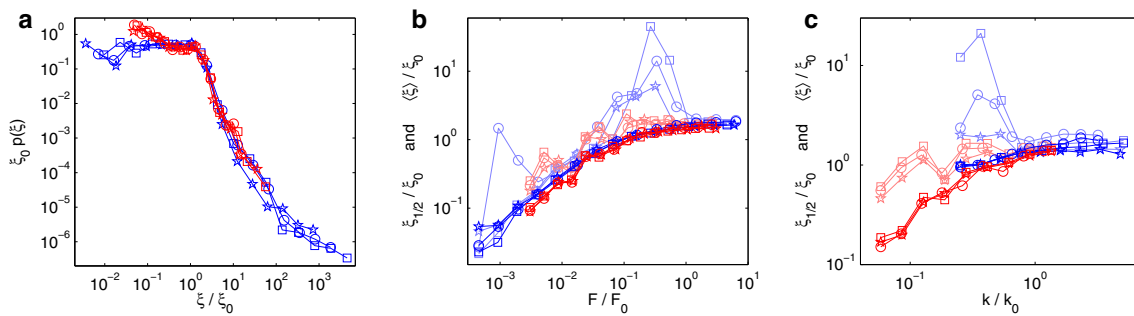
$$\xi_i = \frac{\Delta \bar{d}_i}{\bar{d}_i} \quad (16)$$

where  $\bar{d}_i$  is the median shortest path distance from reference node  $i$  to all other nodes  $j$ , and  $\Delta \bar{d}_i$  the change of this median in response to the network disruption. This impact factor is different for every node and the distribution  $p(\xi)$  gives insight into the variability of how individual nodes are affected by the network disruption [50]. Figure 9a illustrates  $p(\xi)$  for scenarios in which the entire traffic was reduced by 30% through the removal of high-centrality nodes. The distribution  $p(\xi)$  is independent of the measure of centrality and also identical in both networks. Below a typical impact of  $\xi_0$  the distribution of impact factors  $p(\xi)$  is uniform and for  $\xi > \xi_0$  it decreases slowly, ranging over many orders of magnitude. A question that immediately arises is what nodes in the network experience the largest impact. Figures 9b, 9c depict the mean  $\langle \xi \rangle$  and median  $\xi_{1/2}$  conditioned on the flux  $F$  and degree  $k$ . Both the WAN and GCSN exhibit the same dependence, with increasing centrality, the median impact increases monotonically and reaches the typical asymptotic value  $\xi_0$ . However, the mean  $\langle \xi \rangle$  as a function of  $F$  exhibits strong fluctuations, most markedly for intermediate ranges of  $F$ . Low-flux nodes experience a small impact because in the unperturbed network they exhibit large  $\bar{d}$  and changes under the perturbation will typically be small compared to this original  $\bar{d}$ . In contrast, nodes of intermediate centrality have some strong connections to hubs; for them  $\bar{d}$  is initially of intermediate size, and upon removal of the hubs, they experience a relatively larger increase in median effective distance due to the loss of these central connections. A similar effect is seen in the behavior of  $\langle \xi \rangle$  as a function of degree  $k$ .

## 6 Discussion

The comparative analysis of the worldwide air-transportation network and the global cargo-ship network presented here is a first step towards a better understanding of the organizational structure, the evolution and management of large scale infrastructural networks in general. The statistical analysis of node and link centrality measures and their correlations revealed a surprising degree of similarity of both networks despite their different purpose, scale and connectivity. We believe that this is evidence for common underlying principles that govern the growth and evolution of infrastructural networks. This is also supported by the variety of simple, approximate algebraic scaling relations that we extracted from both networks.

Our analysis revealed an unusual discontinuity in the distribution of both link and node betweenness. This suggests that strongly heterogeneous transportation and mobility networks exhibit a natural functional separation of links and nodes into two distinct groups. Interestingly, this discontinuity is localized at the same relative betweenness value and has approximately the same magnitude in both networks. We conclude that this natural separation into different classes of nodes and links might well be a universal feature of these transportation networks as well and



**Fig. 9.** (Color online) The distribution of impact in response to the removal of central nodes of the system as determined by degree (circles), flux (squares), and betweenness (stars) in the WAN (blue) and GCSN (red). (a)  $p(\xi)$  is the probability distribution of impact factors  $\xi$  computed for each node in response to the various attacks, which ranges over many orders of magnitude. (b) Dependence of impact factor as a function of node flux  $F$ . The  $y$ -axis measures the normalized median impact ( $\xi_{1/2}/\xi_0$ , solid lines) and normalized mean impact ( $\langle \xi \rangle / \xi_0$ , faint lines). (c) The dependence of impact factor on node degree.

could be a starting point for further investigations along these lines.

The analysis of network resilience showed that because of their dense connectivity, naive percolation analysis or network diameter inflation do not yield much information. The percolation threshold for both networks lies well beyond many significant real-world network perturbations. The alternative approach based on effective distance, shortest paths, and shortest path trees allows a better, more intuitive representation of networks and resilience analysis, taking into account the fact that nodes that are connected by strong traffic are effectively closer than nodes that are connected by weak links and investigating network perturbations from the viewpoint of chosen reference nodes. Furthermore, the shortest path tree representation revealed an interesting correlation of effective shortest path distance and node centrality measures such as flux, degree, and betweenness and an interesting symmetry in both networks: on average, any node in the network is closest to the subset of nodes with high centrality. This has fundamental implications for spreading phenomena on these types of networks. Where global disease dynamics, for example, are characterized by highly complex spatio-temporal patterns when visualized in conventional geographical coordinates, we expect these patterns become simpler and thus better understood when shortest path tree representations are employed. Since the shortest path tree representations are structurally similar in both networks one might expect a strong dynamic similarity of otherwise unrelated spreading phenomena that occur in these networks, for example the global spread of emergent human infectious diseases on the worldwide air-transportation network and human mediated bio-invasion processes on the global cargo-ship network. We conclude that our results can serve as a starting point for both the development of theories for the evolution of large scale transportation networks and dynamical processes that evolve on them.

The authors wish to acknowledge support from the Volkswagen Foundation.

## References

- OAG Worldwide Ltd. (2007), <http://www.oag.com/>
- UNCTAD, *Review of Maritime Transport 2008*, United Nations Conference on Trade and Development, 2008
- IHS Fairplay, *The source for maritime information and insights* (2008), [www.ihsfairplay.com](http://www.ihsfairplay.com)
- V. Colizza, M. Barthélemy, A. Barrat, A. Vespignani, C.R. Biol. **330**, 364 (2007)
- V. Colizza, A. Barrat, M. Barthélemy, A.J. Valleron, A. Vespignani, PLoS Med. **4**, 95 (2007)
- V. Colizza, A. Barrat, M. Barthélemy, A. Vespignani, Proc. Natl. Acad. Sci. USA **103**, 2015 (2006)
- V. Colizza, A. Barrat, M. Barthélemy, A. Vespignani, B. Math. Biol. **68**, 1893 (2006)
- D. Balcan, V. Colizza, B. Gonçalves, H. Hu, J. Ramasco, A. Vespignani, Proc. Natl. Acad. Sci. USA **106**, 21484 (2009)
- P. Kaluza, A. Koelzsch, M.T. Gastner, B. Blasius, J. R. Soc. Interface **7**, 1093 (2010)
- D. Tompkins, A. White, M. Boots, Ecol. Lett. **6**, 189 (2003)
- G. Ruiz, T. Rawlings, F. Dobbs, L. Drake, T. Mullady, A. Huq, R. Colwell, Nature **408**, 49 (2000)
- K.L.S. Drury, J.M. Drake, D.M. Lodge, G. Dwyer, Ecol. Model. **206**, 63 (2007)
- N.M. Ferguson, D.A.T. Cummings, C. Fraser, J.C. Cajka, P.C. Cooley, D.S. Burke, Nature **442**, 448 (2006)
- M.E. Halloran et al., Proc. Natl. Acad. Sci. USA **105**, 4639 (2008)
- T.D. Hollingsworth, N.M. Ferguson, R.M. Anderson, Nat. Med. **12**, 497 (2006)
- T.D. Hollingsworth, N.M. Ferguson, R.M. Anderson, Emerg. Infect. Dis. **13**, 1288 (2007)
- L.A. Meyerson, H.A. Mooney, Front. Ecol. Environ. **5**, 199 (2007)
- P.E. Hulme, J. Appl. Ecol. **46**, 10 (2009)
- J.M. Levine, C.M. D'Antonio, Conserv. Biol. **17**, 322 (2003)
- L. Hufnagel, D. Brockmann, T. Geisel, Proc. Natl. Acad. Sci. USA **101**, 15124 (2004)
- D. Brockmann, Eur. Phys. J. Special Top. **157**, 173 (2008)
- V. Colizza, A. Barrat, M. Barthélemy, A. Vespignani, BMC Med. **5**, 34 (2007)

23. B.S. Cooper, R.J. Pitman, W.J. Edmunds, N.J. Gay, *PLoS Med.* **3**, 845 (2006)
24. D. Balcan et al., *BMC Med.* **7**, 45 (2009)
25. C. Fraser et al., *Science* **324**, 1557 (2009)
26. United Nations – Department of Economic and Social Affairs, *World population to 2300* (2004), <http://www.un.org/esa/population/unpop.htm>
27. R. Mack, D. Simberloff, W. Lonsdale, H. Evans, M. Clout, F. Bazzaz, *Ecol. Appl.* **10**, 689 (2000)
28. C.S. Kolar, D.M. Lodge, *Science* **298**, 1233 (2002)
29. D. Simberloff, I.M. Parker, P.N. Windle, *Front. Ecol. Environ.* **3**, 12 (2005)
30. G.M. Ruiz, J.T. Carlton, E.D. Grosholz, A.H. Hines, *Am. Zool.* **37**, 621 (1997)
31. O.E. Sala et al., *Science* **287**, 1770 (2000)
32. J.L. Molnar, R.L. Gamboa, C. Revenga, M.D. Spalding, *Front. Ecol. Environ.* **6**, 485 (2008)
33. D. Pimentel, R. Zuniga, D. Morrison, *Ecol. Econ.* **52**, 273 (2005)
34. M.E.J. Newman, *SIAM Rev.* **45**, 167 (2003)
35. A. Barrat, M. Barthelemy, A. Vespignani, *J. Stat. Mech. Theory Exp.* P05003 (2005)
36. A. Vespignani, *Science* **325**, 425 (2009)
37. A. Barrat, M. Barthelemy, A. Vespignani, *Phys. Rev. E* **70**, 066149 (2004)
38. L. Dall’Asta, A. Barrat, M. Barthelemy, A. Vespignani, *J. Stat. Mech. Theory Exp.* P04006 (2006)
39. A. Barabasi, R. Albert, *Science* **286**, 509 (1999)
40. D. Brockmann, L. Hufnagel, T. Geisel, *Nature* **439**, 462 (2006)
41. D. Brockmann, F. Theis, *IEEE Pervas. Comput.* **7**, 28 (2008)
42. R. Guimera, S. Mossa, A. Turttschi, L. Amaral, *Proc. Natl. Acad. Sci. USA* **102**, 7794 (2005)
43. D. Grady, C. Thiemann, D. Brockmann, in preparation (2011)
44. R. Cohen, K. Erez, D. ben Avraham, S. Havlin, *Phys. Rev. Lett.* **85**, 4626 (2000)
45. R. Cohen, S. Havlin, D. ben Avraham, *Phys. Rev. Lett.* **91**, 247901 (2003)
46. R. Albert, H. Jeong, A. Barabasi, *Nature* **406**, 378 (2000)
47. Y. Chen, G. Paul, S. Havlin, F. Liljeros, H.E. Stanley, *Phys. Rev. Lett.* **101**, 058701 (2008)
48. Z. Wu, L.A. Braunstein, V. Colizza, R. Cohen, S. Havlin, H.E. Stanley, *Phys. Rev. E* **74**, 056104 (2006)
49. E. Lopez, R. Parshani, R. Cohen, S. Carmi, S. Havlin, *Phys. Rev. Lett.* **99**, 188701 (2007)
50. O. Woolley Meza, C. Thiemann, D. Grady, D. Brockmann, in preparation (2011)