

Public transport networks: empirical analysis and modeling

C. von Ferber^{1,2}, T. Holovatch^{1,3}, Yu. Holovatch^{4,5,a}, and V. Palchykov⁴

¹ Applied Mathematics Research Centre, Coventry University, Coventry CV1 5FB, UK

² Physikalisches Institut, Universität Freiburg, 79104 Freiburg, Germany

³ Laboratoire de Physique des Matériaux, Université Henri Poincaré, Nancy 1, 54506 Vandœuvre les Nancy Cedex, France

⁴ Institute for Condensed Matter Physics, National Academy of Sciences of Ukraine, 79011 Lviv, Ukraine

⁵ Institut für Theoretische Physik, Johannes Kepler Universität Linz, 4040 Linz, Austria

Received 1st October 2008 / Received in final form 17 December 2008

Published online 14 March 2009 – © EDP Sciences, Società Italiana di Fisica, Springer-Verlag 2009

Abstract. Public transport networks of fourteen cities of so far unexplored network size are analyzed in standardized graph representations: the simple graph of the network map, the bipartite graph of routes and stations, and both one mode projections of the latter. Special attention is paid to the inter-relations and spatial embedding of transport routes. This systematic approach reveals rich behavior beyond that of the ubiquitous scale-free complex network. We find strong evidence for structures in PTNs that are counter-intuitive and need to be explained, among these a pronounced diversity in the expression of typical network characteristics within the present sample of cities, a surprising geometrical behavior with respect to the two-dimensional geographical embedding and an unexpected attraction between transport routes. A simple model based on these observations reproduces many of the identified PTN properties by growing networks of attractive self-avoiding walks.

PACS. 02.50.-r Probability theory, stochastic processes, and statistics – 07.05.Rm Data presentation and visualization: algorithms and implementation – 89.75.Hc Networks and genealogical trees

1 Introduction

The recent general interest in networks of man-made and natural systems has lead to the advancement of a complex network science through careful analysis of various network systems using empirical, simulational, and theoretical tools [1–5]. In this work we strive to identify the distinguishing properties of public transport networks (PTN) of 14 large cities when interpreted as complex network graphs. These networks may be expected to share general features of other transportation networks [3] like the airport [6–13], railway [14], or power grid networks [6,15,16]. These features include evolutionary growth, optimization, and usually an embedding in two dimensional (2D) space.

The evolution of a city's PTN is closely related to the growth of the city and therefore influenced by numerous factors of geographical, historical, and social origin. However, there is ample evidence that PTNs of different cities share common statistical properties that possibly arise due to their functional purposes [17–29]. Some of these properties have been analyzed in former studies. Here, our objective is to systematically analyze PTNs in all standardized graph representations: the simple graph of the network map, the bipartite graph of routes and stations, and both one mode projections of the latter and furthermore, to

identify inter-relations and spatial embedding properties of transport routes which are unique to PTNs. Finally, based on the empirical observations, we embark to formulate a model with simple growth rules for that generate PTNs with network characteristics matching empirical results.

Previous studies have often analyzed specific sub-networks of PTNs [17–20,22–24,26]. Examples are the Boston [17–20] and Vienna [20] subway networks and the bus networks of cities in Poland [22] and China [24,26]. However, as far as the bus-, subway- or tram-subnetworks are not closed systems the inclusion of additional subnetworks has significant impact on the overall network properties as has been shown for the subway and bus networks of Boston [18,19].

All PTNs analyzed within our study are either operated by a single operator or by a small number of operators with a coordinated schedule, as e.g. expressed by a central website from which our data was obtained. Rather than artificially dividing these centrally organized networks into subnetworks of different means of transport like bus and metro or in a 'urban' and an 'sub-urban' part we treat each full PTN as an entity.

Our choice for the selection of fourteen major cities (see Tab. 1) [30,31] was motivated by the idea to collect network samples from cities of different geographical, cultural, and economical background. Apart from the

^a e-mail: hol@icmp.lviv.ua

Table 1. Cities analyzed in this study. N : number of PTN stations; R : number of PTN routes; S : mean route length (mean number of stations per route). Types of transport taken into account: Bus, Electric trolleybus, Ferry, Subway, Tram, Urban train.

City	N	R	S	Type
Berlin	2992	211	29.4	BSTU
Dallas	5366	117	59.9	B
Düsseldorf	1494	124	28.5	BST
Hamburg	8084	708	25.5	BFSTU
Hong Kong	2024	321	39.6	B
Istanbul	4043	414	31.7	BST
London	10937	922	34.2	BST
Los Angeles	44629	1881	52.9	B
Moscow	3569	679	22.2	BEST
Paris	3728	251	38.2	BS
Rome	3961	681	26.8	BT
São Paulo	7215	997	58.3	B
Sydney	1978	596	16.3	B
Taipei	5311	389	70.5	B

systematic analysis explained above this choice also extends to PTNs of much larger size as compared to previous work [21,22] which considered PTNs of typically hundreds of stations.

This paper is organized as follows. The next Section 2 sets up and defines the different representations in which the PTN will be analyzed, Sections 3–4.2 explore the network properties in these representations. We separately analyze in Section 3 local characteristics, such as node degrees and clustering coefficients, and in Section 4 global characteristics, such as path length distributions and betweenness centralities. Paragraphs 4.3 and 4.4 are devoted to characteristics that are unique to PTNs and networks with similar construction principles. Section 4.3 analyses the phenomenon of sequences of routes proceeding in parallel along a sequence of stations, a feature we call ‘harness’ effect. Section 4.4 analyzes the network embedding in geographical space. Our findings for the statistics of real-world PTNs are supported by simulations of an evolutionary model of PTNs as displayed in Section 5. Conclusions and an outlook are given in Section 6. Some of our results have been preliminarily announced in reference [25]. Supplementary material is available to the interested reader in reference [32].

2 PT network topology

A straightforward representation of a PT map in the form of a graph represents every station by a node while the edges correspond to the links that exist between stations due to the PT routes servicing them (see e.g. Figs. 1, 2a). Let us first introduce a simple graph to represent this situation, see Figure 2b. In the following we will refer to this graph as the \mathbb{L} -space graph [22] or simply as \mathbb{L} -space. This graph represents each station by a node, a link between nodes indicates that there is at least one route that services the two corresponding stations consecutively. No multiple links are allowed. In the analysis of

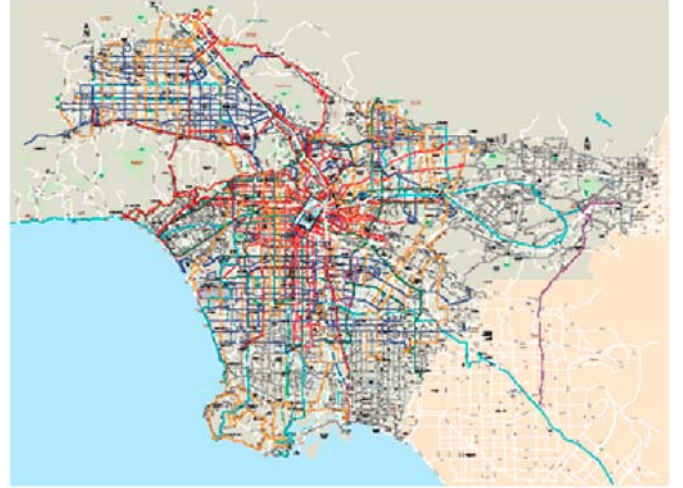


Fig. 1. (Color online) One of the networks we analyze in this study. The Los Angeles PTN consists of $R = 1881$ routes and $N = 44629$ stations, some of them are shown in this map.

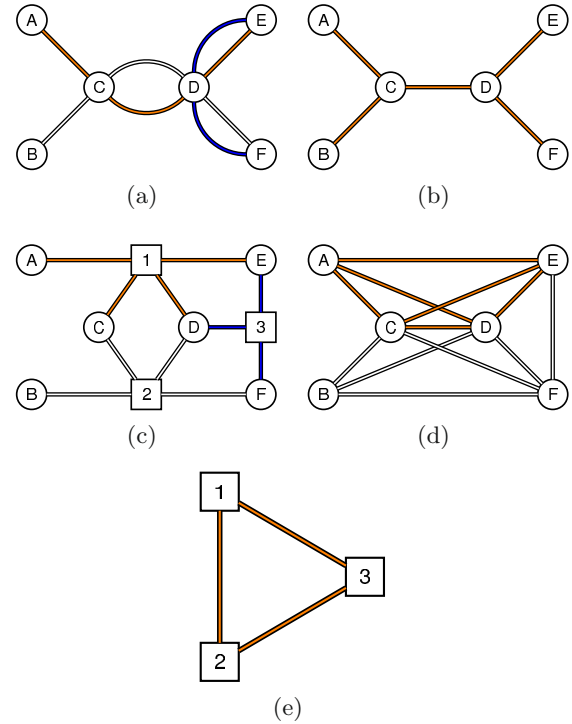


Fig. 2. (Color online) (a) a simple public transport map. Stations A–F are serviced by routes No 1 (shaded orange), No 2 (white), and No 3 (dark blue). (b) \mathbb{L} -space graph. (c) \mathbb{B} -space bipartite graph. Route nodes are shown as squares. (d) \mathbb{P} -space graph, the complete sub-graph corresponding to route No 1 is highlighted (shaded orange). (e) \mathbb{C} -space graph of routes.

PTNs, this \mathbb{L} -space representation has been used in references [18,21–23,26].

A somewhat different concept is that of a bipartite graph which has proven useful in the analysis of cooperation networks [3,33]. In this representation which we call \mathbb{B} -space both routes and stations are represented by nodes [24,25,27]. Each route node is linked to all station

Table 2. PTN characteristics in different spaces (subscripts refer to \mathbb{L} -, \mathbb{P} -, and \mathbb{C} -spaces, correspondingly). k : node degree; $\kappa = \langle z \rangle / \langle k \rangle$ where z is the number of next nearest neighbors; ℓ^{\max} , $\langle \ell \rangle$: maximal and mean shortest path length (10); \mathcal{C}^b : betweenness centrality (14); c : relation of the mean clustering coefficient to that of the classical random graph of equal size (7). Averaging has been performed with respect to corresponding network, only the mean shortest path $\langle \ell \rangle$ is calculated with respect to the largest connected component.

City	$\langle k_{\mathbb{L}} \rangle$	$\kappa_{\mathbb{L}}$	$\ell_{\mathbb{L}}^{\max}$	$\langle \ell_{\mathbb{L}} \rangle$	$\langle \mathcal{C}_{\mathbb{L}}^b \rangle$	$c_{\mathbb{L}}$	$\langle k_{\mathbb{P}} \rangle$	$\kappa_{\mathbb{P}}$	$\ell_{\mathbb{P}}^{\max}$	$\langle \ell_{\mathbb{P}} \rangle$	$\langle \mathcal{C}_{\mathbb{P}}^b \rangle$	$c_{\mathbb{P}}$	$\langle k_{\mathbb{C}} \rangle$	$\kappa_{\mathbb{C}}$	$\ell_{\mathbb{C}}^{\max}$	$\langle \ell_{\mathbb{C}} \rangle$	$\langle \mathcal{C}_{\mathbb{C}}^b \rangle$	$c_{\mathbb{C}}$
Berlin	2.58	1.96	68	18.5	2.6×10^4	52.8	56.61	11.47	5	2.9	2.9×10^3	41.9	27.56	4.43	5	2.2	1.2×10^2	4.75
Dallas	2.18	1.28	156	52.0	1.4×10^5	55.0	100.58	11.23	8	3.2	5.9×10^3	48.6	11.09	3.45	7	2.7	9.2×10^1	5.34
Düsseldorf	2.57	1.96	48	12.5	8.6×10^3	24.4	59.01	10.56	5	2.6	1.2×10^3	19.7	32.18	2.47	4	1.8	4.9×10^1	2.23
Hamburg	2.65	1.85	156	39.7	1.4×10^5	254.7	50.38	7.96	11	4.7	1.4×10^4	132.2	17.51	4.49	10	4.0	9.9×10^2	28.3
Hong Kong	3.59	3.24	60	11.0	1.0×10^4	60.3	125.67	10.20	4	2.2	1.3×10^3	11.7	98.98	2.12	3	1.7	1.2×10^2	2.14
Istanbul	2.30	1.54	131	29.7	5.7×10^4	41.0	76.88	10.59	6	3.1	4.2×10^3	41.5	52.81	3.86	5	2.3	2.6×10^2	5.00
London	2.60	1.87	107	26.5	1.4×10^5	320.6	90.60	16.97	6	3.3	1.2×10^4	90.0	49.91	6.80	6	2.6	7.4×10^2	11.1
Los Angeles	2.37	1.59	210	37.1	7.9×10^5	645.3	97.99	17.21	11	4.4	7.4×10^4	399.6	40.11	8.42	10	3.6	2.3×10^3	22.1
Moscow	3.32	6.25	27	7.0	1.1×10^4	127.4	65.47	26.48	5	2.5	2.7×10^3	38.0	109.37	4.57	4	1.9	3.2×10^2	3.59
Paris	3.73	5.32	28	6.4	1.0×10^4	78.5	50.92	24.06	5	2.7	3.1×10^3	59.6	39.95	4.67	4	1.9	1.1×10^2	2.72
Rome	2.95	2.02	87	26.4	5.0×10^4	163.4	69.05	11.34	6	3.1	4.2×10^3	41.4	59.40	4.86	5	2.5	5.1×10^2	7.04
São Paulo	3.21	4.17	33	10.3	3.4×10^4	268.0	137.46	19.61	5	2.7	6.0×10^3	38.2	151.72	4.25	4	2.0	5.2×10^2	4.27
Sydney	3.33	2.54	34	12.3	7.3×10^3	82.9	42.88	7.79	7	3.0	1.3×10^3	33.6	65.02	2.92	6	2.4	3.5×10^2	6.30
Taipei	3.12	2.42	74	20.9	5.3×10^4	186.2	236.65	12.96	6	2.4	3.6×10^3	15.4	93.33	2.95	5	1.8	1.6×10^2	2.44

nodes that it services. No direct links between nodes of the same type occur (see Fig. 2c). Obviously, in \mathbb{B} -space the neighbors of a given route node are all stations that it services while the neighbors of a given station node are all routes that service it.

There are two one-mode projections of the bipartite graph of \mathbb{B} -space. The projection to the set of station nodes is the so-called \mathbb{P} -space graph, Figure 2d. The complementary projection to route nodes leads to the \mathbb{C} -space graph, Figure 2e, of route nodes where any two route nodes are neighbors if they share a common station.

The \mathbb{P} -space graph representation [14,22] has proven particularly useful in the analysis of PTNs [14,20,22,25,26]. The nodes of this graph are stations and they are linked if they are serviced by at least one common route. In this way the neighbors of a \mathbb{P} -space node are all stations that can be reached without changing means of transport and each route gives rise to a complete \mathbb{P} -subgraph, see Figure 2d.

It is worthwhile to note the real world significance of these seemingly abstract ‘spaces’. To give an example, the average length of a shortest path $\langle \ell_{\mathbb{L}} \rangle$ in an \mathbb{L} -space graph gives the average number of stops one has to pass to travel between any two stations. When represented in \mathbb{P} -space, the mean shortest path $\langle \ell_{\mathbb{P}} \rangle$ counts the average number of changes one has to do to travel between two stations while the corresponding mean \mathbb{C} -space path length $\langle \ell_{\mathbb{C}} \rangle$ counts the average number of changes needed to pass between any two routes. As another example let us note the node degree k : for the \mathbb{L} -space graph the node degree of a station is the number of other stations within one stop distance; in the bipartite \mathbb{B} -space graph the degree of a station is the number of routes servicing it, while the degree of a route is the number of its stations; in the \mathbb{P} -space graph the degree $k_{\mathbb{P}}$ of a station is the number of stations reachable without changing the route; whereas in the \mathbb{C} -space graph

the degree $k_{\mathbb{C}}$ of a route is the number of other routes one can transfer to.

Table 2 lists some of the PTN characteristics we have obtained for the cities under consideration using publicly available data from the web pages of local transport organizations [30,31]. To limit the data presented, this and further tables are restricted to the basic results discussed within this article. The interested reader may find supplementary material in [32].

3 Local network properties

Let us first examine the properties of the PTNs determined by the immediate neighborhood of the nodes as measured by its size, its interconnectedness and the correlations within this neighborhood.

3.1 Neighborhood size (node degree)

The size of the neighborhood of a node as given by its degree often indicates its importance e.g. as a hub within the network. In large networks created by randomly connecting nodes, hubs are rare while in real networks they are often found with much higher probability. Formally this is measured by the behavior of the tail of the node degree distribution. Denoting by $p(k)$ the normalized node degree distribution, the mean node degree k is given by the average

$$\langle k \rangle = \sum_{k=1}^{k^{\max}} p(k)k = \frac{2M}{N}. \quad (1)$$

Here, M is the number of links and N the number of edges of the graph while k^{\max} stands for the maximal

node degree. For the finite size Erdős-Rényi [34,35] random graph the node degree distribution $p(k)$ is binomial, which for fixed $\langle k \rangle$ in the infinite case becomes a Poisson distribution.

The higher organization of real world networks usually leads to slower decaying distributions. Typical classes of networks have either exponential or power law tails. Exponentially decaying distributions for large degrees k are characterized by

$$p(k) \sim \exp(-k/\hat{k}), \quad (2)$$

where the scale \hat{k} is of the order of the mean node degree. Scale-free degree distributions that decay according to a power law have a tail of the form

$$p(k) \sim 1/k^\gamma. \quad (3)$$

The exponent γ further classifies the network [36]. If $\gamma < 2$ the distribution has no finite average $\langle k \rangle$ in the infinite network limit. If $\gamma < 3$ there is no finite second moment and the network has no percolation threshold with respect to a dilution of its nodes. Its connected component remains robust against random failure of any number of its nodes. When $\gamma > 4$, however, its percolation and other properties are expected to be similar to those of exponentially decaying networks.

Both exponential and power law decay of the degree distribution can be modeled by assuming a non equilibrium growth process of the network by which in consecutive time steps nodes and links are added to the existing network [4]. If the added nodes are arbitrarily linked to any of the existing nodes an exponential tail results, however, if the probability to connect to a given existing node is a linear function of its degree one can show that the resulting degree distribution develops a power law tail. The latter mechanism to explain the abundant occurrence of power laws is also referred to as preferential attachment or ‘rich get richer’ [37–39]. As far as PTNs obviously are evolving networks, their evolution may be expected to follow similar mechanisms. However, scale-free networks have also been shown to arise when minimizing both the effort for communication and the cost for maintaining connections [40,41]. Moreover, this kind of optimization was shown to lead to small world properties [42] and to explain the appearance of power laws in a general context [43]. Therefore, scale-free behavior in PTNs could also be related to obvious objectives to optimize their operation.

Figures 3 and 4 show the node degree distributions for PTNs of several cities in \mathbb{L} -, \mathbb{P} -, and \mathbb{C} -spaces. Note, that the monotonously decreasing curves displayed for the \mathbb{P} - and \mathbb{C} -spaces are cumulative distributions defined as:

$$P(k) = \sum_{q=k}^{k_{\max}} p(q). \quad (4)$$

The data for \mathbb{L} - and \mathbb{P} -spaces in Figures 3a, 3b is shown in log-linear plots together with fits to an exponential decay (2). The latter distributions are nicely described by an

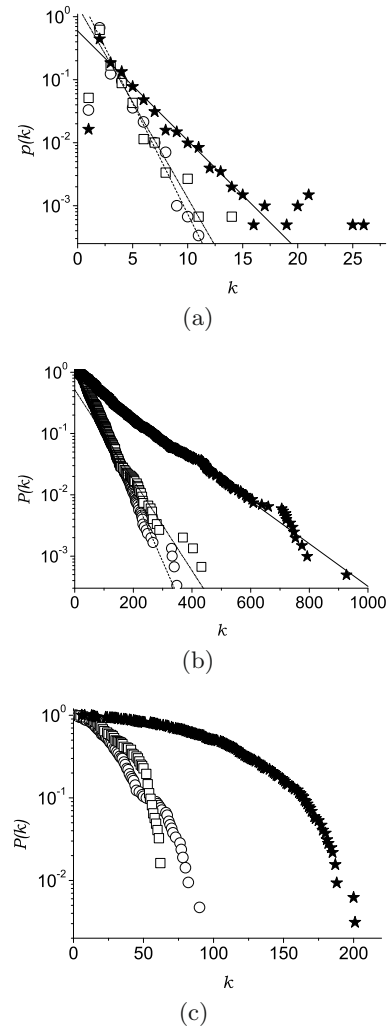


Fig. 3. (a) Node degree distributions of PTN of several cities in \mathbb{L} -space. (b) Cumulative node degree distribution in \mathbb{P} -space. (c) Cumulative node degree distribution in \mathbb{C} -space. Berlin (circles, $\hat{k}_{\mathbb{L}} = 1.24$, $\hat{k}_{\mathbb{P}} = 39.7$), Düsseldorf (squares, $\hat{k}_{\mathbb{L}} = 1.43$, $\hat{k}_{\mathbb{P}} = 58.8$), Hong Kong (stars, $\hat{k}_{\mathbb{L}} = 2.50$, $\hat{k}_{\mathbb{P}} = 125.1$).

exponential decay. As far as the \mathbb{L} -space data is concerned, we find evidence for an exponential decay for about half of the cities analyzed, while the other part rather follow a power law decay (3), see Table 3.

Figures 4a, 4b show the corresponding plots for three other cities on a log-log scale. Here, these plots are shown together with fits to a power law (3). Numerical values of the fit parameters \hat{k} and γ for different cities are given in Table 3. Here, values in parentheses indicate a less reliable fit. In the case when none of the equations (2), (3) lead to reliable data, both fit parameters are given in parentheses in the table. The typical range of data points which could be fitted was of the order of 90% or more both for \mathbb{L} - and \mathbb{P} -spaces. The value of the fit parameters was considered to be reliable if the absolute value of the Pearson correlation coefficient exceeded $R_{\mathbb{L}} = 0.984$ and $R_{\mathbb{P}} = 0.990$ in \mathbb{L} and \mathbb{P} -spaces, correspondingly. Exceptions from this rule are the \mathbb{L} -space fits for the PTNs

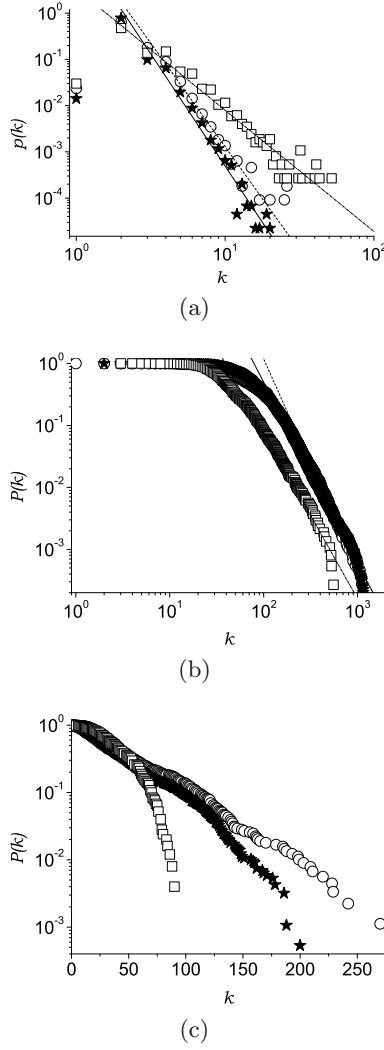


Fig. 4. (a) Node degree distributions of the PTNs of several cities in \mathbb{L} -space. (b) Cumulative node degree distributions in \mathbb{P} -space. (c) Cumulative node degree distribution in \mathbb{C} -space. London (circles, $\gamma_{\mathbb{L}} = 4.48$, $\gamma_{\mathbb{P}} = 3.89$), Los Angeles (stars, $\gamma_{\mathbb{L}} = 4.85$, $\gamma_{\mathbb{P}} = 3.92$), Paris (squares, $\gamma_{\mathbb{L}} = 2.62$, $\gamma_{\mathbb{P}} = 3.70$).

of Paris, Rome ($R_{\mathbb{L}} \simeq 0.97$, but 97% of data points are covered), and London (with 72% of data points covered and $R_{\mathbb{L}} \simeq 0.985$). For \mathbb{P} -space, exceptions are the PTNs of Paris ($R_{\mathbb{P}} = 0.993$) and São Paulo ($R_{\mathbb{P}} = 0.999$), where the fit covered only $\sim 60\%$ of data points. Note, that for \mathbb{L} -space the fit was done for the plain node degree distribution $p(k)$, whereas for \mathbb{P} -space the parameters $\gamma_{\mathbb{P}}$ or $\hat{k}_{\mathbb{P}}$ were determined by fitting the cumulative distribution (4).

While the node degree distribution of almost half of the cities in the \mathbb{L} -space representation display a power law decay (3), this is in general not the case for the \mathbb{P} -space. However, the data for the PTNs of Hamburg, London, Los Angeles, and Paris (see Fig. 4b) give first evidence of power law behavior of $P(k)$ even in the \mathbb{P} -space representation. Previous results concerning node-degree distributions of PTNs in \mathbb{L} - and \mathbb{P} -spaces [22,26] seemed to indicate that in general the degree distribution may be power-law like in \mathbb{L} -space but never in \mathbb{P} -

Table 3. Parameters of the PTN node degree distributions fit to an exponential (2) and power law (3) behavior. Bracketed values indicate less reliable fits. Subscripts refer to \mathbb{L} - and \mathbb{P} -spaces [31].

City	$\gamma_{\mathbb{L}}$	$\hat{k}_{\mathbb{L}}$	$\gamma_{\mathbb{P}}$	$\hat{k}_{\mathbb{P}}$
Berlin	(4.30)	1.24	(5.85)	39.7
Dallas	5.49	(0.78)	(4.67)	64.2
Düsseldorf	3.76	(1.43)	(4.62)	(58.8)
Hamburg	(4.74)	1.46	4.38	(60.7)
Hong Kong	(2.99)	2.50	(4.40)	125.1
Istanbul	4.04	(1.13)	(2.70)	86.7
London	4.48	(1.44)	3.89	(143.3)
Los Angeles	4.85	(1.52)	3.92	(201.0)
Moscow	(3.22)	(2.15)	(2.91)	50.0
Paris	2.62	(3.30)	3.70	(100.0)
Rome	(3.95)	1.71	(5.02)	54.8
São Paolo	2.72	(4.20)	(4.06)	225.0
Sydney	(4.03)	1.88	(5.66)	38.7
Taipei	(3.74)	1.75	(5.16)	201.0

space. This was interpreted [22] as being due to strongly correlated connections between stations in \mathbb{L} -space and nearly randomly linked routes, as also expressed by a low clustering coefficient in \mathbb{C} -space, see below. Our present study, which includes a much less homogeneous selection of cities (Ref. [22] was exclusively based on Polish cities) shows that almost any combination of different distributions in \mathbb{L} - and \mathbb{P} -spaces may occur. We note that even within the small sub-group formed by Hamburg, Los Angeles, London and Paris there is no alignment to ‘typical behavior’.

In \mathbb{C} -space the decay of the node degree distribution is exponential or faster, as one can see from the plots in Figures 3c and 4c. From the cities presented there, only the PTNs of Berlin, London, and Los Angeles are governed by an exponential decay.

For most cities that show a power law degree distribution in \mathbb{L} -space the corresponding exponent $\gamma_{\mathbb{L}}$ is $\gamma_{\mathbb{L}} \sim 4$. Also the exponents found for the PTNs of Polish cities of similar size N lie in this region: $\gamma_{\mathbb{L}} = 3.77$ for Kraków (with number of stations $N = 940$), $\gamma_{\mathbb{L}} = 3.9$ for Łódź ($N = 1023$), $\gamma_{\mathbb{L}} = 3.44$ for Warsaw ($N = 1530$) [22]. According to the general classification of scale-free networks [2] this indicates that in many respect these networks are expected to behave similar to those with exponential node degree distribution. Prominent exceptions to this rule are the PTNs of Paris ($\gamma_{\mathbb{L}} = 2.62$) and São-Paulo ($\gamma_{\mathbb{L}} = 2.72$). Note, that values of $\gamma_{\mathbb{L}}$ in the range $2.5 \div 3.0$ were recently reported for the bus networks of three cities in China: Beijing ($N = 3938$), Shanghai ($N = 2063$), and Nanjing ($N = 1150$) [26].

A conclusion from our survey of the various degree distributions is that they appear much more diverse than expected and that with respect to these there is no simple division of the PTNs at hand into two or even three classes.

3.2 Clustering

While the node degree counts the neighbors of a node, the connectivity within its neighborhood may be quantified in terms of the so called clustering coefficient. The latter is defined as

$$C_i = \frac{2y_i}{k_i(k_i - 1)} \text{ for } k_i \geq 2, \quad (5)$$

where y_i is the number of links between the k_i nearest neighbors of the node i . $C_i \equiv 0$ for $k_i = 0, 1$. The clustering coefficient of a node may also be defined as the probability of any two of its randomly chosen neighbors to be connected. For the mean value of the clustering coefficient of an Erdős-Rényi random graph one finds

$$\langle C \rangle_{\text{ER}} = \frac{\langle k \rangle}{N} = \frac{2M}{N^2}. \quad (6)$$

In Table 2 we give the values of the mean clustering coefficient in \mathbb{L} -, \mathbb{P} -, and \mathbb{C} -spaces. The highest absolute values of the clustering coefficient are found in \mathbb{P} -space, where their range is given by $\langle C_{\mathbb{P}} \rangle = 0.7 \div 0.9$ (c.f. with $\langle C_{\mathbb{L}} \rangle = 0.02 \div 0.1$). This is not surprising since in \mathbb{P} -space each route gives rise to a fully connected (complete) sub-graph between all of its stations. In order to make numbers comparable we normalize the mean clustering coefficient by that of a random graph (6) of the same size:

$$c = N^2 \langle C \rangle / (2M). \quad (7)$$

In \mathbb{L} - and \mathbb{P} -representations we find the mean clustering coefficient to be larger by orders of magnitude relative to the random graph. This difference is less pronounced in \mathbb{C} -space indicating a lower degree of organization in these graphs. Most prominently, we find the values to vary strongly within the sample of the 14 cities.

In \mathbb{P} -space the clustering coefficient of a node is strongly correlated with the node degree. All stations i belonging to the complete subgraph of a single route have $C_i = 1$, while C_i generally decreases if i belongs to more than one route. Averaging the \mathbb{P} -space clustering coefficient over all nodes with given degree k we confirm that it decays as a function of k following a power law

$$\langle C_{\mathbb{P}}(k) \rangle \sim k^{-\beta}. \quad (8)$$

Within a simple model of networks with star-like topology this exponent is found to be $\beta = 1$ [22]. In transport networks, this behavior has been observed before for the Indian railway network [14] as well as for Polish PTNs [22]. In our case, the values of the exponent β for the networks studied range from 0.65 (São Paulo) to 0.96 (Los Angeles) again showing significant diversity within our sample.

These obvious differences in the locally observable structure may be assumed to reflect a strong diversity within the concepts according to which various PTNs are structured. Comparing the division between weak and strongly clustered PTNs we find no alignment with the different classes of degree distributions adding to the idea of an individual profile of each city's PTN with respect to the various network characteristics.

Table 4. Nearest neighbor and next nearest neighbor assortativities $r^{(1)}$ and $r^{(2)}$ in different spaces for the whole PTN.

City	$r_{\mathbb{L}}^{(1)}$	$r_{\mathbb{L}}^{(2)}$	$r_{\mathbb{P}}^{(1)}$	$r_{\mathbb{P}}^{(2)}$	$r_{\mathbb{C}}^{(1)}$	$r_{\mathbb{C}}^{(2)}$
Berlin	0.158	0.616	0.065	0.441	0.086	0.318
Dallas	0.150	0.712	0.154	0.728	0.290	0.550
Düsseldorf	0.083	0.650	0.041	0.494	0.244	0.180
Hamburg	0.297	0.697	0.087	0.551	0.246	0.605
Hong Kong	0.205	0.632	-0.067	0.238	0.131	0.087
Istanbul	0.176	0.726	-0.124	0.378	0.282	0.505
London	0.221	0.589	0.090	0.470	0.395	0.620
Los Angeles	0.240	0.728	0.124	0.500	0.465	0.753
Moscow	0.002	0.312	-0.041	0.296	0.208	0.011
Paris	0.064	0.344	-0.010	0.258	0.060	-0.008
Rome	0.237	0.719	0.044	0.525	0.384	0.619
São Paulo	-0.018	0.437	-0.047	0.266	0.211	0.418
Sydney	0.154	0.642	0.077	0.608	0.458	0.424
Taipei	0.270	0.721	0.009	0.328	0.100	0.041

3.3 Generalized assortativities

To describe correlations between the properties of neighboring nodes in a network the notion of assortativity was introduced measuring the correlation between the node degrees of neighboring nodes in terms of the mean Pearson correlation coefficient [44,45]. Here, we propose to generalize this concept to also measure correlations between the values of other node characteristics (other observables). For any link i let X_i and Y_i be the values of the observable at the two nodes connected by this link. Then the correlation coefficient is given by:

$$r = \frac{M^{-1} \sum_i X_i Y_i - [M^{-1} \sum_i \frac{1}{2}(X_i + Y_i)]^2}{M^{-1} \sum_i \frac{1}{2}(X_i^2 + Y_i^2) - [M^{-1} \sum_i \frac{1}{2}(X_i + Y_i)]^2} \quad (9)$$

where summation is performed with respect to the M links of the network. Taking X_i and Y_i to be the node degrees equation (9) is equivalent to the usual formula for the assortativity of a network [44]. Here, we will call this special case the degree assortativity $r^{(1)}$. In separate work we have investigated generalized assortativities for a number of other network characteristics [32]. Here, besides the assortativity $r^{(1)}$, we discuss the behavior of the generalized assortativity $r^{(2)}$ for the number z of next nearest neighbors. The numerical values of the assortativities $r^{(1)}$ and $r^{(2)}$ of all PTNs are listed in Table 4 for the \mathbb{L} -, \mathbb{P} - and \mathbb{C} -spaces. With respect to the values of the standard node degree assortativity $r_{\mathbb{L}}^{(1)}$ in \mathbb{L} -space, we find two groups of cities. The first is characterized by values $r_{\mathbb{L}}^{(1)} = 0.1 \div 0.3$. Although these values are still small they signal a finite preference for assortative mixing. That is, links tend to connect nodes of similar degree. In the second group of cities these values are very small $r_{\mathbb{L}}^{(1)} = -0.02 \div 0.08$ showing no preference in linkage between nodes with respect to node degrees. PTNs of both large and medium sizes are present in each of the groups. This indicates the absence of correlations between network size and degree assortativity $r_{\mathbb{L}}^{(1)}$ in \mathbb{L} -space. Measuring the same quantity in the \mathbb{P} - and \mathbb{C} -spaces, we observe different behavior. In \mathbb{P} -space

almost all cities are characterized by very small (positive or negative) values of $r_P^{(1)}$ with the exception of the PTNs of Istanbul ($r_P^{(1)} = -0.12$) and Los Angeles ($r_P^{(1)} = 0.12$). On the contrary, in \mathbb{C} -space PTNs demonstrate clear assortative mixing with $r_C^{(1)} = 0.1 \div 0.5$. An exception is the PTN of Paris with $r_C^{(1)} = 0.06$.

As we have seen above, the PTNs demonstrate assortative ($r^{(1)} > 0$) or neutral ($r^{(1)} \sim 0$) mixing with respect to the node degree (first nearest neighbors number) k . Defining an assortativity $r^{(2)}$ with respect to the number z of second next nearest neighbors we explore the correlation of a wider environment of adjacent nodes. Due to the fact that in this case the two connected nodes share at least part of this environment (the first nearest neighbors of a node form part of the second nearest neighbors of the adjacent node) one may expect the assortativity $r^{(2)}$ to be non-negative. The results for $r^{(2)}$ shown in Table 4 appear to confirm this assumption. In all the spaces considered, we find that all PTNs that belong to the group of neutral mixing with respect to k also belong to the same group with respect to the second nearest neighbors. For those PTNs that display significant nearest neighbors assortativity $r^{(1)}$ we find that the second nearest neighbor assortativity $r^{(2)}$ is in general even stronger in line with the above reasoning.

From the above observations on assortativity within our sample of PTNs we note further evidence for diversity ranging from indefinite to clearly pronounced assortativities $r_L^{(1)}$ and $r_C^{(1)}$ which appear uncorrelated with other properties of the network such as the size or the specific behavior of e.g. the degree distribution.

4 Global characteristics

4.1 Shortest paths

Let $\ell_{i,j}$ be the length of a shortest path between sites i and j in a given graph. Note, that $\ell_{i,j}$ is well-defined only if the nodes i and j belong to the same connected component of the graph. In the following we will restrict considerations to the largest (so-called giant) connected component, GCC. Denoting the path length distribution within the GCC as $\Pi(\ell)$, the mean shortest path is

$$\langle \ell \rangle = \sum_{\ell=1}^{\ell^{\max}} \Pi(\ell) \ell, \quad (10)$$

where ℓ^{\max} is the maximal shortest path length found within the GCC. In general, the shortest path length distributions obtained in \mathbb{L} -, \mathbb{P} -, and \mathbb{C} -spaces that we have analyzed [32] are nicely described by an asymmetric unimodal distribution [22]:

$$\Pi(\ell) = A \ell \exp(-B \ell^2 + C \ell), \quad (11)$$

where A, B , and C are parameters. However, additional structures may lead to deviations from this behavior as

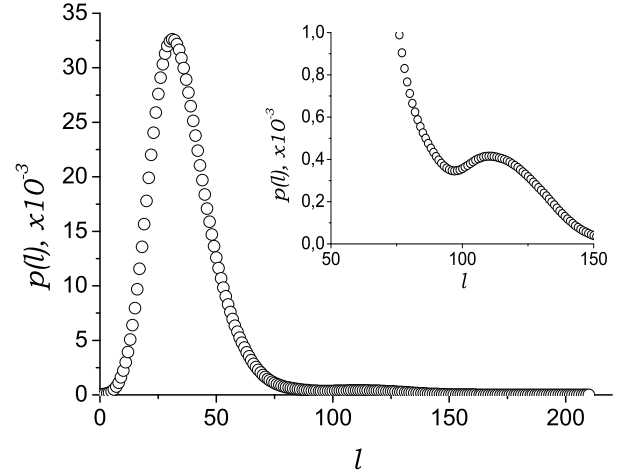


Fig. 5. Shortest path length distribution in \mathbb{L} -space, $P_L(\ell)$, for the PTN of Los Angeles.

can be seen from Figure 5, which shows the mean shortest path length distribution in \mathbb{L} -space $P_L(\ell)$ for Los Angeles. One observes a second local maximum on the right shoulder of the distribution. Qualitatively this behavior may be explained by assuming that the PTN consists of more than one community. For the simple case of one large community and a second smaller one at some distance this situation will result in short intra-community paths which will give rise to a global maximum and a set of longer paths that connect the larger to the smaller community resulting in additional local maxima. Such a situation definitely appears to be present in the case of the Los Angeles PTN, see Figure 1.

Of particular interest is the mean shortest path length between nodes of given degrees k and q , $\ell(k, q)$. As has been shown in [46], this relation can be approximated by

$$\ell(k, q) = A - B \log(kq). \quad (12)$$

For random networks the coefficients A and B can be calculated exactly [47]. A rather good agreement with equation (12) was found for the majority of the \mathbb{L} -space graphs of Polish PTNs analyzed in [22]. Within our study which includes PTNs of much larger size, we do not observe a similar alignment for all cities. The suggested logarithmic dependence (12) does occur also for the \mathbb{L} -space graphs of larger cities, however, with a much more pronounced scatter of data for large values of the product kq . In Figure 6 we plot the mean path $\ell_L(k, q)$ for the \mathbb{L} -space graphs of the PTNs of Berlin, Hong Kong, Rome, and Taipei, where the relation (12) is observed with better accuracy. Note, however, that due to the scatter of data a logarithmic dependence frequently is indistinguishable from a power law with a small exponent.

The dependency of the average path length on the degrees of both end nodes of the path may be reduced to a dependency on the degree of a single end node. We define $\ell(k)$, the mean shortest path between any node of degree k and other nodes of the network. For the majority of the analyzed cities the dependence of the mean path $\ell_L(k)$ on

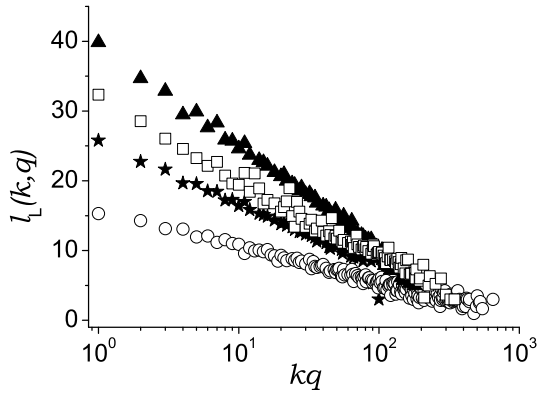


Fig. 6. Mean \mathbb{L} -space paths $\ell_{\mathbb{L}}(k, q)$ as function of kq for the PTNs of Berlin (stars), Hong Kong (circles), Rome (triangles), and Taipei (squares).

the node degree k in \mathbb{L} -space can be approximated by a power law

$$\ell_{\mathbb{L}}(k) \sim k^{-\alpha_{\mathbb{L}}}. \quad (13)$$

We find that the value of the exponent varies in the range $\alpha_{\mathbb{L}} = 0.17 \div 0.27$. It is instructive to compare this result with results obtained in reference [48] for the same characteristics calculated for correlated growing networks. For deterministic scale-free networks $\ell(k)$ was found to be characterized by a logarithmic law with power-law corrections, whereas for stochastic scale-free networks $\ell(k)$ was shown to follow logarithmic behaviour. Furthermore, networks with an exponential node-degree distribution displayed a linear law $\ell(k) \sim a - bk$. Obviously, the small values of the exponent $\alpha_{\mathbb{L}}$ found for the PTNs in our study do not exclude a logarithmic law, however the linear dependence can be ruled out. Note, that within our sample of PTNs one finds both scale-free and exponential node degree distributions. However, an essential difference between the construction principles of PTNs and of the graphs of reference [48] is that the latter are so-called ‘citation graphs’ (where new connections do not emerge between already existing nodes), whereas there is no such restriction for PTNs.

In \mathbb{P} -space, the shortest path length ℓ_{ij} gives the minimal number of routes required to be used in order to reach site j starting from the site i . The higher the node degree, the easier it is to access other routes in the network. Therefore, also in \mathbb{P} -space one expects a decrease of $\ell_{\mathbb{P}}(k)$ when k increases. Apart from an expected decrease we find a tendency to a power-law decay with small powers, sometimes almost indistinguishable from a logarithmic behavior. The value of the exponent $\alpha_{\mathbb{P}}$ varies in the interval $\alpha_{\mathbb{P}} = 0.09$ (for Sydney) to $\alpha_{\mathbb{P}} = 0.17$ (for Dallas) and is centered around $\alpha_{\mathbb{P}} = 0.12 \div 0.13$. The mean path $\ell_{\mathbb{P}}(k, q)$ is found to decrease as a function of kq also in \mathbb{P} -space, but with much more pronounced scattering than in \mathbb{L} -space. An analysis of further characteristics related to shortest path lengths ℓ_{ij} can be found in [32].

Concluding we note that the mean lengths of the shortest paths as function of the end node degrees show no special structure within the sample of PTNs studied. In

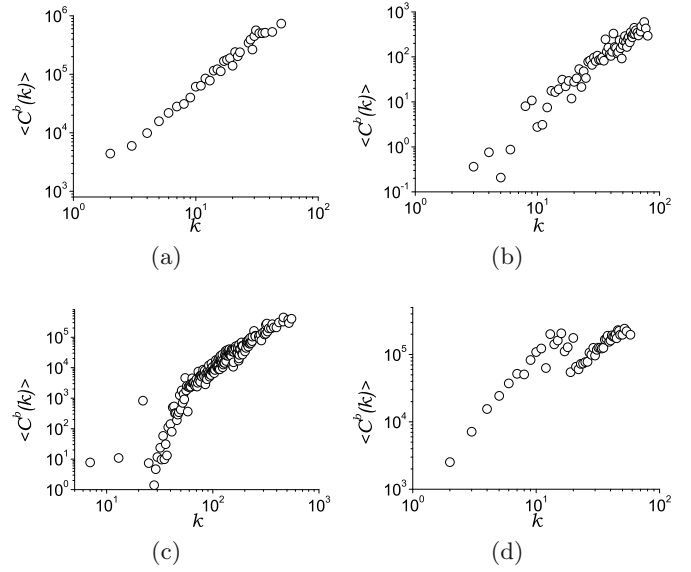


Fig. 7. Mean betweenness centrality $\langle C^b(k) \rangle$ - degree k correlations for the PTN of Paris in (a) \mathbb{L} -, (b) \mathbb{C} -, (c) \mathbb{P} -, and (d) \mathbb{B} -spaces.

general the observed behavior does not significantly deviate from the logarithmic behavior that is expected for random graphs.

4.2 Betweenness centrality

To measure the importance of a given node with respect to different properties of a graph a number of so-called centrality measures have been introduced [49–53]. Referring the interested reader to reference [32] for a more extensive survey on centrality measures of PTNs, we here discuss data related to the betweenness centrality which measures the importance of a node with respect to the connectivity between other nodes of the network. The betweenness centrality $C^b(i)$ of a node i is calculated as

$$C^b(i) = \sum_{j \neq i \neq k} \frac{\sigma_{jk}(i)}{\sigma_{jk}}, \quad (14)$$

where σ_{jk} is the number of shortest paths between nodes j and k and $\sigma_{jk}(i)$ is the number of these paths that go via node i . Numerical values of the mean betweenness centrality (14) are given in Table 1 for the \mathbb{L} -, \mathbb{P} - and \mathbb{C} -space graphs.

The betweenness centrality (14) of a given node measures the share of the shortest paths between nodes that this node mediates. It is obvious that a node with a high degree has a higher probability to be part of any path connecting other nodes. This relation between C^b and the node degree may be quantified by plotting the mean betweenness centrality $\langle C^b(k) \rangle$ averaged among nodes with degree k as function of k . In Figures 7 we present corresponding results for the PTN of Paris in \mathbb{L} -, \mathbb{C} -, \mathbb{P} -, and

B-spaces. Especially well expressed is the betweenness-degree correlation in L-space (Fig. 7a) and with somewhat less precision in C-space (Fig. 7b). In both cases there is a clear tendency to a power law $\langle C^b(k) \rangle \sim k^\eta$ with an exponent $\eta = 2 \div 3$.

In the plots for both B- and P-spaces we observe the occurrence of two regimes which correspond to small and large degrees k . This separation however has a different origin in each of these cases. In the B-space representation, the network consists of nodes of two types, route nodes and station nodes. Typically, station nodes are connected only to a low number of routes while there is a minimal number of stations per route. One may thus identify the low degree behavior as describing the betweenness of station nodes, while the high degree behavior corresponds to that of route nodes. In the overlap region of the two regimes one may observe that when having the same degree station nodes have a higher betweenness than route nodes.

In the P-space representation on the other hand, the occurrence of two regimes is a feature of this representation. Stations that are part of only a single route and thus within the P-graph belong only to the complete subgraph corresponding to this route (recall Fig. 2d) are not part of any shortest P-space path between other nodes and have a betweenness centrality of $C^b = 0$. The decreasing contribution of these stations to the average $\langle C^b(k) \rangle$ leads a steep slope in the low degree regime. For degrees higher than the maximal route length these stations no longer contribute and the slope rather describes the correlation between the degree and finite mean betweenness values. Instead of a steep slope in the low degree regime reference [22] observes a saturation; this may be due to an exclusion of the zero-betweenness nodes from the average. Very similar betweenness – degree relations as shown in Figure 7 are found for most of the other cities in our sample with slightly varying quality of expression. We emphasize however, that this uniformity of the correlation between the degrees of the nodes and their respective betweenness is strictly speaking valid only for the average value $\langle C^b(k) \rangle$. When analyzing the importance of individual nodes e.g. with respect to the vulnerability of the network against failure or attack the betweenness centrality turns out to be a much more sensitive measure than the node degree [29].

4.3 Harness

Besides the local and global properties of networks described above which can be defined in any type of network, there are some characteristics that are unique for PTNs and networks with similar construction principles. A particularly striking example is the fact that as far as the routes share the same grid of streets and tracks often a number of routes will proceed in parallel along shorter or longer sequences of stations. Similar phenomena are observed in networks built with space consuming links such as cables, pipes, neurons, etc. In the present case this behavior may be easily worked out on the basis of sequences

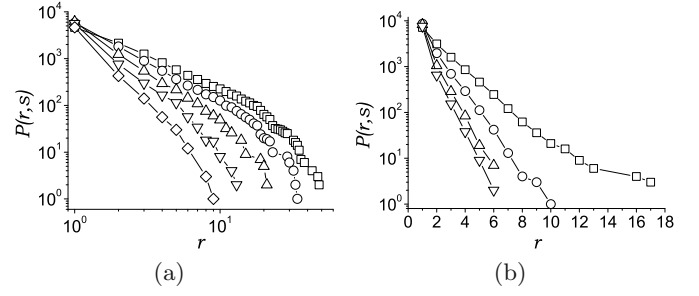


Fig. 8. Cumulative harness distributions. (a) Istanbul PTN ($s = 2(\square), 6(\circ), 11(\triangle), 16(\nabla), 21(\diamond)$). (b) Moscow PTN $s = 3(\square), 6(\circ), 9(\triangle), 11(\nabla)$.

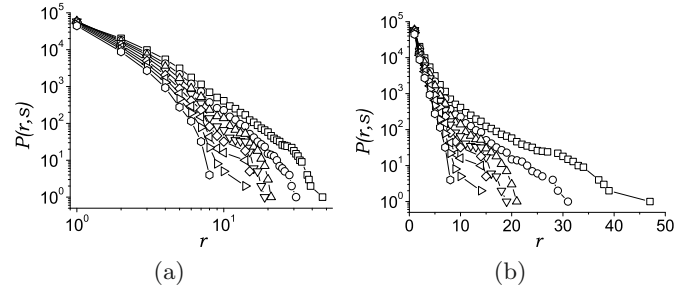


Fig. 9. Cumulative harness distributions for Los Angeles PTN. From above: $s = 2(\square), 4(\circ), 6(\triangle), 9(\nabla), 13(\diamond), 17(\triangleleft), 21(\triangleright), 26(\circ)$. (a) log-log scale; (b) log-linear scale.

of stations serviced by each route. To quantify this behavior we use the recently introduced notion of network harness [25]. It is described by the harness distribution $P(r, s)$: the number of sequences of s consecutive stations that are serviced by r parallel routes. Similar to the node-degree distributions, we observe that the harness distribution for some cities (Hong Kong, Istanbul, Paris, Rome, São Paulo, Sydney) may be described by a power law:

$$P(r, s) \sim r^{-\gamma_s}, \quad \text{for fixed } s, \quad (15)$$

whereas the PTNs of other cities (Berlin, Dallas, Düsseldorf, London, Moscow) are better described by an exponential decay:

$$P(r, s) \sim \exp(-r/\hat{r}_s), \quad \text{for fixed } s. \quad (16)$$

As examples we show the harness distributions for Istanbul (Fig. 8a) and for Moscow (Fig. 8b). Sometimes (we observe this for Los Angeles and Taipei), there is a crossover from a power law to an exponential regime for larger s . We show this crossover for the PTN of Los Angeles in Figure 9, where it is particularly obvious.

As one can observe in Figures 8, 9 the harness distribution $P(r, s)$ for fixed s decays faster for longer sequences s . For PTNs for which the harness distribution follows a power law (15) the corresponding exponents γ_s are found in the range of $\gamma_s = 2 \div 4$. For those distributions with an exponential decay the scale \hat{r}_s (16) varies in the range $\hat{r}_s = 1.5 \div 4$. The power laws observed for the behavior of $P(r, s)$ indicate a certain level of organization and

planning which may be driven by the need to minimize the costs of infrastructure and secondly by the fact that points of interest tend to be clustered in certain locations of a city. Note, that this effect may be seen as a result of the strong interdependence of the evolutions of both the city and its PTN.

As noted above, the notion of harness may be useful also for the description of other networks with similar properties. On the one hand, the harness distribution is closely related to distributions of flow and load on the network. On the other hand, in the situation of space-consuming links (such as tracks, cables, neurons, or pipes) the information about the harness behavior may be important with respect to the spatial optimization of networks.

From our observations we conclude that there is strong evidence for a significant harness effect within the organization of PTN networks according to which network routes are often found aligned following the same geographical path along segments of varying length and ‘thickness’. The details of the harness distribution which quantifies this behavior however differ considerably adding to the diversity of behavior found within our PTN sample for many of the properties measured.

It should be emphasized that with respect to network optimization the harness property may at first seem completely counter-intuitive: why should a route that is e.g. added to the network follow the path of previous, already existing routes, instead of exploring yet unserved nearby areas? We may name at least two possible reasons for the empirically confirmed harness behavior: the first is the minimization of the cost for infrastructure which is most evident for means of transport that need tracks but relevant also with respect to maintaining e.g. bus stops. Other, more operation related reasons are those of interconnectivity minimizing the effort needed to change from one route to the other and of system redundancy, ensuring a higher transport frequency on important segments of the routes.

Related unexpected behavior of the routes concerning their geographical embedding is observed and discussed in the following section.

4.4 Geographical embedding

So far, we have discussed the properties of PTNs without reference to their geographical embedding. The fact that this subject has so far been left aside also by previous studies of PTNs with respect to their complex network behavior, is due mainly to the lack of easily accessible data on the locations of stations and routes. Note, however, a study on the fractal dimension of railway networks, reference [54]. For the present work we have been able to obtain such data for stations of the Berlin PTN as well as for those of the metro subnetwork of Paris. For the Berlin network the positions of the stations were extracted in an automated way from interactive maps provided on the web-pages of the operator [55] which (invisibly) contain the geographical coordinates of the stations. For the

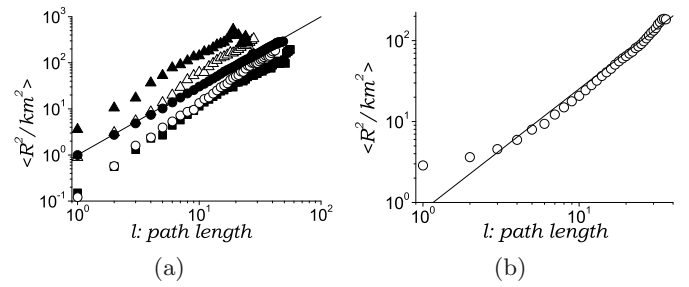


Fig. 10. Distance – path length – relation $\langle \mathcal{R}^2(\ell) \rangle$ in comparison with that of a two dimensional self-avoiding walk (solid line, $\sim \ell^{3/2}$) for (a) different means of transport within the Berlin PTN (bus (■), tram (○), u-bahn (△), s-bahn (▲), simulated city (●), and (b) for the Paris metro network.

Metro network of Paris these coordinates were retrieved by hand using a free web based map service [56].

The question we pose here is, what is the distance \mathcal{R} between initial and final stations of a passenger’s journey traveling for ℓ stops on a single route? For routes optimizing the time of passenger travel a naive consideration might lead to the expectation of distance growing linearly with path length ℓ at least on larger scales. Surprisingly, the empirical data show quite a different behavior (see Fig. 10). For all means of transport analyzed within the Berlin PTN as well as for the metro network the dependence of the mean square distance $\langle \mathcal{R}^2(\ell) \rangle$ on ℓ is well described by a power law

$$\langle \mathcal{R}^2(\ell) \rangle \sim \ell^{2\nu} \quad (17)$$

with an exponent ν that is significantly smaller than one. For most transport routes this exponent appears to be near to $\nu = 3/4$, which is the well known self-avoiding walk (or Flory-) exponent in two dimensions [57] corresponding to a fractal dimension of $D \sim 1.33$. For the different Berlin subnetworks we find exponents ranging from $\nu = 0.82$ for the bus routes to $\nu = 0.9$ and 0.96 for the subway and tram routes. The s-bahn data is distorted due to a ring structure within this sub-network. The Paris metro data supports an exponent of $\nu = 0.82$ when excluding the short distance contributions. For comparison, the fractal dimensions D of some regional railway networks (not individual routes) reported in reference [54] are of the order $D \sim 1.5 \div 1.8$.

Self-avoiding walks, apart from observing the constraint of non-self-intersection evolve randomly. The fact that PT routes at least within the present sample appear to display the same scaling symmetry is quite unexpected. In particular, this behavior seems to be at odds with the requirement of minimizing passengers traveling time between origin to destination. The latter argument, however, ignores the time passengers spend walking to the initial and from the final stations. Including these, one understands the need for the routes to cover larger areas by meandering through neighborhoods. Given the requirements for a PTN to cover a metropolitan area with a limited number of routes while simultaneously offering fast transport across the city one may speculate that routes scaling like SAWs may present an optimal solution. Further research is obviously needed to support this claim.

5 Modeling PTNs

5.1 Motivation and description of the model

Having at hand the above described wealth of empirical data and analysis with respect to typical scenarios found in a variety of real-world PTNs we feel in the position to propose a model that may capture the characteristic features of these networks. In view of the diversity found in our sample, it would be in vein to try to construct a model that quantitatively reproduces the data of a given city. The aim of the present model is to show that a few simple rules and a low number of parameters suffice to generate PTNs that display profiles which with respect to most observables are within the range of those found in real world PTNs. Nonetheless it should be capable of discriminating between some of the various scenarios observed.

Essential basic properties of PTNs that we intend to implement or reproduce within our model are the following: (a) the model is to be based on routes and stations and allow for \mathbb{L} - \mathbb{P} - \mathbb{C} - and \mathbb{B} -space representations; (b) the model should be embedded in two dimensions and reproduce the SAW scaling behavior of the routes; (c) the model should be able to generate realistic degree distributions; (d) the model must generate realistic harness distributions.

If we were only to reproduce the degree distribution of the network, standard models such as random networks [4,58] or preferential attachment type models [6,39,59–62] would suffice. The evolution of such networks however is based on the attachment of nodes. For the description of PTNs the concept of routes as finite sequences of stations is essential [5,23,25,28] and allows for the representation with respect to the spaces defined above. Moreover, taking a route as the essential element of PTN growth allows to account for the bipartite structure of this network [20,24,27,33]. Therefore, the growth dynamics in terms of routes will be a central ingredient of our model. Another obvious requirement is the embedding of this model in two-dimensional space. To simplify matters we will restrict the model to a two-dimensional grid, in particular to a square lattice. Both the observations of power law degree distributions as well as the occurrence of the corresponding harness distributions described above indicate a preference of routes to service common stations (i.e. an attraction between routes).

Let us describe our model in more detail. As noticed above, a route will be modeled as a sequence of stations that are adjacent nodes on a two-dimensional square lattice. Following the observation of SAW scaling symmetry for the geographical embedding we choose each PTN route to be a self-avoiding walk. To incorporate all the above features the model is set up as follows. A model PTN consists of R routes each with S stations constructed on a possibly periodic $X \times X$ square lattice. The dynamics of the route generation adheres to the following rules:

1. Construct the first route as a SAW of S lattice sites.
2. Construct the $R - 1$ subsequent routes as SAWs with the following preferential attachment rules:

- (a) choose a terminal station at \mathbf{x}_0 with probability

$$p \sim k_{\mathbf{x}_0} + a/X^2; \quad (18)$$

- (b) choose any subsequent station \mathbf{x} of the route with probability

$$p \sim k_{\mathbf{x}} + b. \quad (19)$$

In (18), (19) $k_{\mathbf{x}}$ is the number of times the lattice site \mathbf{x} has been visited before (the number of routes that pass through \mathbf{x}). Note, that to ensure the SAW property any route that intersects itself is discarded and its construction is restarted with step 2a).

5.2 Global topology of model PTN

Let us first investigate the global topology of this model as function of its parameters. We first fix both the number of routes R and the number of stations S per route as well as the size of the lattice X . This leaves us with essentially two parameters a and b , equations (18), (19). Dependencies on R and S will be studied below.

For the real-world PTNs as studied in the previous sections, almost all stations belong to a single component, GCC, with the possible exception of a very small number of routes. Within the network however we often observe the harness effect of several routes proceeding in parallel for a sequence of stations. Let us first investigate from a global point of view which parameters a and b reproduce realistic maps of PTNs. In Figure 11 we show simulated PTNs on lattices 300×300 for $R = 1024$, $S = 64$ and different values of the parameters a and b . Each route is represented by a continuous line tracing the path along its sequence of stations. For representation purposes, parallel routes are shown slightly shifted. Thus, the line thickness and intensity of colors indicate the density of the routes.

The parameter a quantifies the possibility to start a new route outside the existing network. For vanishing $a = 0$ the resulting network always consists of a single connected component, while for finite values of a a few or many disconnected components may occur. The results for $a = 0$ and varying b parameters are independent of the lattice size X provided X is sufficiently large to accommodate the network without boundary effects. Parameter b governs the evolution of each single subsequent route. If $a = 0$ and $b = 0$ the only allowed sites according to equations (18), (19) are those of the first SAW route as far as the choice is restricted to sites x with a finite number k_x of previous visits. The shape variation of the simulated PTNs as b is increased for fixed $a = 0$ is shown in the first row of Figure 11. For small values of $b = 0 \div 0.1$ almost all routes of the simulated PTN follow the same path with only a few deviations. Shifting b to $b = 0.2$ the area covered by the routes increases while the majority of the routes are concentrated on a small number of paths. Further shifting b to $b = 0.5$ and beyond we find a wider distributed coverage with the central part of the network remaining the most densely covered area. This is due to the non-equilibrium growth process described by equations (18), (19).

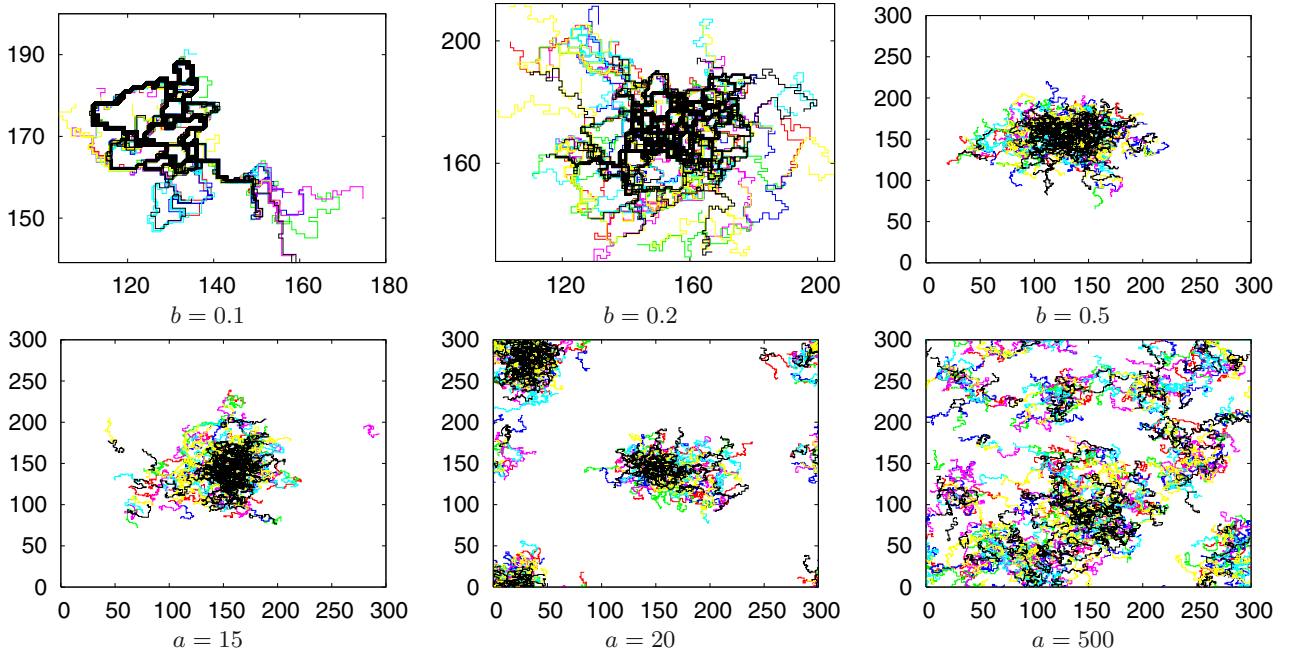


Fig. 11. (Color online) PTN maps of different simulated cities of size 300×300 with $R = 1024$ routes of $S = 64$ stations each (color online). First row: $a = 0$, $b = 0.1 \div 0.5$. Second row: $b = 0.5$, $a = 15 \div 500$. With an increase of b routes cover more and more area. Increase of a leads to clusterisation of the network.

Table 5. Characteristics of the simulated PTN with $X = 300$, $a = 0$ for different parameters R , S , and b . The rest of notations as in Table 2.

R	S	b	$\langle k_L \rangle$	κ_L	ℓ_L^{\max}	$\langle \ell_L \rangle$	$\langle C_L^b \rangle$	$\langle k_P \rangle$	κ_P	ℓ_P^{\max}	$\langle \ell_P \rangle$	$\langle C_P^b \rangle$	c_P	$\langle k_C \rangle$	κ_C	ℓ_C^{\max}	$\langle \ell_C \rangle$	$\langle C_C^b \rangle$	c_C
256	16	0.5	2.92	1.66	61	20.8	4.7×10^3	44.15	3.18	7	3.0	4.7×10^2	7.98	86.39	1.36	6	1.9	1.2×10^2	2.22
256	16	5.0	2.99	1.74	80	21.7	7.5×10^3	42.95	3.76	9	3.4	8.8×10^2	11.7	59.96	1.99	8	2.2	1.5×10^2	2.79
256	32	0.5	2.76	1.60	127	38.1	3.0×10^4	84.45	4.32	8	3.3	1.9×10^3	13.6	60.51	1.75	7	2.2	1.6×10^2	2.90
256	32	5.0	2.90	1.72	177	43.1	5.3×10^4	74.24	5.22	10	4.0	3.8×10^3	23.7	33.06	2.69	9	2.8	2.3×10^2	4.55
512	16	0.5	2.95	1.68	73	22.5	6.7×10^3	50.07	3.39	7	3.1	6.5×10^2	9.14	169.7	1.44	6	1.9	2.3×10^2	2.25
512	16	5.0	3.12	1.78	80	23.3	1.0×10^4	51.56	3.79	10	3.5	1.2×10^3	12.3	115.3	2.24	9	2.1	2.9×10^2	2.88
512	32	0.5	2.83	1.63	166	44.2	4.7×10^4	99.53	4.56	10	3.6	2.8×10^3	15.7	118.4	2.03	9	2.2	3.0×10^2	2.92
512	32	5.0	3.12	1.79	175	44.6	7.2×10^4	97.05	5.37	9	3.9	4.7×10^3	22.2	60.36	3.08	8	2.7	4.4×10^2	5.04
1024	64	0.5	2.86	1.66	325	80.7	3.3×10^5	242.2	6.32	9	3.7	1.1×10^4	23.4	213.3	2.42	8	2.2	6.1×10^2	3.10
1024	64	1.0	2.97	1.72	355	88.5	4.8×10^5	222.2	6.74	12	4.2	1.7×10^4	32.4	143.9	2.97	11	2.5	7.9×10^2	4.39

When introducing a finite a parameter, new routes may be started anywhere on the lattice which results in a lattice size dependency. To partly compensate for this, the impact of a is normalized by X^2 in (18). The variation of the simulated PTN maps for increasing a and fixed $b = 0.5$ is shown in the second row of Figure 11. For $a < 15$ one observes the formation of a single large cluster with only a few individual routes occurring outside this cluster. Slightly increasing a beyond $a = 15$ one finds a sharp transition to a situation with several (two or more) clusters. For much larger values of a the number of clusters further increases and the situation becomes more and more homogeneous: the routes tend to cover all available lattice space area.

5.3 Statistical characteristics of model PTN

From the above qualitative investigation we conclude that realistic PTN maps are obtained for small or vanishing a

and $b \geq 0.5$. In the following we will fix $a = 0$ and X large enough as discussed above. To quantitatively investigate the behavior of the simulated networks on the remaining parameters including R and S let us now compare their statistical characteristics with those we have empirically obtained for real-world networks. In Table 5 we have chosen to list the same characteristics of the simulated PTNs as selected for the real-world networks in Table 2. To provide for additional checks of the correlations between simulated and real-world networks, we present the characteristics in all L-, P-, and C-spaces. Let us note that our choice of the underlying grid to be a square lattice limits the number of nearest neighbors of a given station in L-space to $k_L \leq 4$. Moreover, as far as no direct links between these neighbors occur, the clustering coefficient in L-space vanishes, $c_L = 0$. Nonetheless, as we discuss below, both characteristics display nontrivial behavior similar to real-world networks when measured for P- and C-spaces.

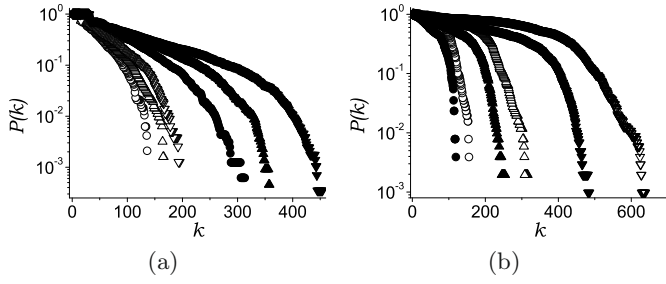


Fig. 12. Cumulative node degree distributions $P(k)$ (4) for several simulated PTNs in (a) P- and (b) C-spaces. $R = 256$, $S = 16$ (\circ), $R = 256$, $S = 32$ (\bullet), $R = 512$, $S = 16$ (\triangle), $R = 512$, $S = 32$ (\blacktriangle), $R = 1024$, $S = 16$ (∇), $R = 1024$, $S = 32$ (\blacktriangledown).

As noted above we choose a vanishing parameter $a = 0$ and $b = 0.5$ and for comparison $b = 5.0$. The data shown in the Table was obtained for simulated PTNs of different numbers of routes, $R = 256, 512, 1024$ and route lengths $S = 16, 32, 64$. In the range of parameters covered in the Table we observe only weak changes of the various characteristics. Natural trends are that with the increase of the number of routes R the maximal and mean shortest path length increases in all spaces. This is most pronounced in L-space, while it is weakest in C-space. A similar increase is observed in L-space when increasing the number of stations S per route. Choosing the values of R in the range $R = 256 \div 1024$ and $S = 16, S = 32$ the average and maximal values of the characteristics studied here are found within the ranges seen for real-world PTNs, see Table 2. More detailed information is contained in the distributions of these characteristics and their correlations.

Let us examine the node degree distributions of some selected PTNs. As explained above, the L-space degrees are restricted by the geometry of the underlying square lattice. Thus we may observe non-trivial distributions only in P-, C-, and B-spaces. The cumulative node degree distributions in P-space are shown in Figure 12a. All these distributions display two regions each governed by an exponential decay with a separate scale. Note, that increasing both S and R leads to an increase of the ranges over which these regions extend. This is in line with the results for real world PTNs found in previous studies [22,26] as well as in Section 3. Within the parameter ranges chosen here the current model does not seem to attain a power law node degree distribution in P-space.

Comparing the C-space node degree distributions for real-world and simulated PTNs (Figs. 3c and 12b, correspondingly) one again finds a definite tendency to an exponential behavior with two different scales in both cases. As can be expected we observe that the scale of the exponential decay increases with the number of routes R while it decreases with the number of stations per route S .

Cumulative harness distributions $P(r, s)$ for two simulated networks with different values of the parameter b ($b = 0.2$, $b = 1.0$) are shown in Figure 13. These appear to reproduce the harness behavior of real world networks as given in Figures 8 and 9. Both exponential and scale-free

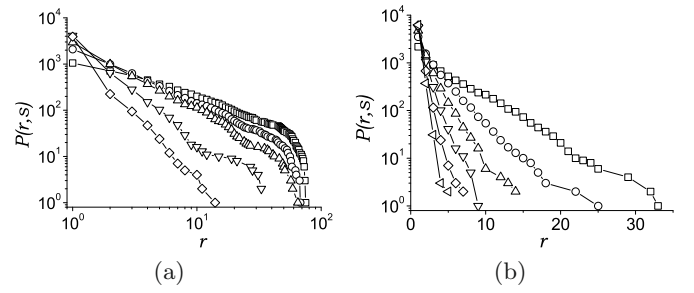


Fig. 13. Cumulative harness distributions $P(r, s)$ for the simulated PTN with $R = 256$, $S = 32$. (a) $a = 0$, $b = 0.2$, $s = 2$ (\square), 4 (\circ), 6 (\triangle), 11 (∇), 16 (\diamond). (b) $a = 0$, $b = 1.0$, $s = 2$ (\square), 3 (\circ), 4 (\triangle), 5 (∇), 6 (\diamond), 7 (\blacktriangleleft). Compare with plots in Figures 8, 9 for the real-world networks.

behavior as observed for the real-world PTNs is found. A prominent feature demonstrated by Figure 13 is that one can tune the decay behavior by changing the parameter b . For small values of b the probability of a route to proceed in parallel with other routes is high. Thus for small b the $P(r, s)$ distribution shows a high probability for the formation of ‘hubs’ of parallel routes as reflected by its power-law decay distribution. For larger b such hubs are suppressed as shown by the exponential decay of their distribution.

Summarizing, the comparison of the statistical characteristics of real world networks with those of simulated ones one can definitely state that the model proposed above captures many essential features of real world PTNs. This is especially evident if one includes into the the comparison different network representations (different spaces) as performed above.

6 Conclusions

This paper was driven by two main objectives towards the analysis of urban public transport networks. First, we wanted to present a systematic survey of statistical properties of PTNs based on the data for cities of so far unexplored network size. The second objective was to present a model that with a small number of simple rules would be capable to reproduce the main properties.

Especially helpful in our analysis was the use of different network representations (different spaces, introduced in Sect. 2). Whereas former PTN studies used some of these, here within a systematic approach we calculate PTN characteristics as they show up in all L-, P-, C-, and B-spaces.

The networks under consideration appear to be strongly correlated small-world structures with high values of clustering coefficients and comparatively low mean shortest path values. Standard network characteristics that we find in these various representations correspond to features a passenger is interested in when using public transport. For example, any two stops in Paris are on the average separated by $\langle \ell_L \rangle - 1 = 5.4$ stations (with a maximal value of 27) and to travel between them one should do

$\langle \ell_P \rangle - 1 = 1.7$ changes on average. The power-law node degree distributions observed for many networks in \mathbb{L} - and for some in \mathbb{P} -space give strong evidence of correlations within these networks. However, for the properties of degree distributions as well as for features of these networks, such as clustering, assortativity and others we find considerable diversity in their expression. Recent work on urban street networks found classifications that discriminate between properties of different classes of city organization. For the present sample of PTNs however, we conclude that there is no simple division of the PTNs we studied into well defined groups as e.g. seen for street and canal networks [63,64] where a division into a few groups was found (however analyzing only small areas of city maps). This result is far from obvious: one might have expected that networks all set up in large urban areas and serving an almost identical purpose would turn out to display strongly aligned properties. However, this diversity is an empirical fact and one that would remain hidden if we had restricted our observations to only a handful of measurements.

Beyond traditional network characteristics there are specific features unique to PTNs and networks with similar construction principles that we have addressed. In particular, public transport routes are often found to proceed in parallel for a sequence of stations. While the very fact that several routes should follow the same path may seem counter-intuitive (why should a route retrace another's path instead of exploring nearby unserved areas?), we have quantified this behavior in terms of the harness distribution and given possible explanations noting costs of infrastructure, and operational advantages such as system redundancy. The harness concept may also be useful for a quantitative description of other embedded networks with real space links such as cables, pipes, or neurons etc.

Moreover, our analysis of the geographical data for Berlin and Paris reveals a self-avoiding walk scaling of PTN routes a fact strongly supported by the empirical study which again appears to be counter-intuitive (should a line not be straight to minimize time of travel). We give a first explanation speculating that this shape of the routes may result from an optimization with respect to total passenger traveling time, area coverage and costs of operation.

The network growth model that we developed captures both of these special features of PTN as well as generating profiles of network characteristics in the various representations which are in line with those found for real world PTNs. By varying only a single parameter one may e.g. discriminate between scale-free and exponential harness distributions, both of which are observed in real cities. The method used, a non equilibrium growth model in terms of attractive self-avoiding walks (SAW) on a square lattice may further be extended to study the effects of geographical constraints e.g. coast-lines, rivers and bridges or disorder. Note in this context that SAW-scaling is unaffected by weak disorder [65].

Obviously, the two objectives in the PTN study we have so far achieved in this paper – the empirical analysis and the modeling – naturally call for an analytic approach.

This will be a task for forthcoming studies. Another natural continuation of this work will be the analysis of different possibly dynamic phenomena that may occur on and with PTNs. Of particular interest is the robustness of PTNs against targeted attacks and random failures [29].

Yu.H. acknowledges support of the Austrian FWF project 19583-PHY.

References

1. R. Albert, A.-L. Barabási, *Rev. Mod. Phys.* **74**, 47 (2002)
2. S.N. Dorogovtsev, J.F.F. Mendes, *Adv. Phys.* **51**, 1079 (2002)
3. M.E.J. Newman, *SIAM Rev.* **45**, 167 (2003)
4. S.N. Dorogovtsev, S.N. Mendes, *Evolution of Networks* (Oxford University Press, Oxford, 2003)
5. Yu. Holovatch, O. Olemskoi, C. von Ferber, T. Holovatch, O. Mryglod, I. Olemskoi, V. Palchykov, *J. Phys. Stud.* **10**, 247 (2006)
6. L.A.N. Amaral, A. Scala, M. Barthélemy, H.E. Stanley, *Proc. Natl. Acad. Sci. USA* **97**, 11149 (2000)
7. R. Guimera, L.A.N. Amaral, *Eur. Phys. J. B* **38**, 381 (2004)
8. R. Guimera, S. Mossa, A. Turtshi, L.A.N. Amaral, *Proc. Nat. Acad. Sci. USA* **102**, 7794 (2005)
9. A. Barrat, M. Barthélemy, R. Pastor-Satorras, A. Vespignani, *Proc. Nat. Acad. Sci. USA* **101**, 3747 (2004)
10. L.-P. Chi, R. Wang, H. Su, X.-P. Xu, J.-S. Zhao, W. Li, X. Cai, *Chin. Phys. Lett.* **20**, 1393 (2003)
11. Y. He, X. Zhu, D.-R. He, *Int. J. Mod. Phys. B* **18**, 2595 (2004)
12. W. Li, X. Cai, *Phys. Rev. E* **69**, 046106 (2004)
13. W. Li, Q.A. Wang, L. Nivonen, A. Le Méhauté, *Physica A* **368**, 262 (2006)
14. P. Sen, S. Dasgupta, A. Chatterjee, P.A. Sreeram, G. Mukherjee, S.S. Manna, *Phys. Rev. E* **67**, 036106 (2003)
15. P. Crucitti, V. Latora, M. Marchiori, *Physica A* **338**, 92 (2004)
16. R. Albert, I. Albert, G.L. Nakarado, *Phys. Rev. E* **69**, 025103 (2004)
17. M. Marchiori, V. Latora, *Physica A* **285**, 539 (2000)
18. V. Latora, M. Marchiori, *Phys. Rev. Lett.* **87**, 198701 (2001)
19. V. Latora, M. Marchiori, *Physica A* **314**, 109 (2002)
20. K.A. Seaton, L.M. Hackett, *Physica A* **339**, 635 (2004)
21. C. von Ferber, Yu. Holovatch, V. Palchykov, *Condens. Matter Phys.* **8**, 225 (2005), e-print [arXiv:cond-mat/0501296](https://arxiv.org/abs/cond-mat/0501296)
22. J. Sienkiewicz, J.A. Holyst, *Phys. Rev. E* **72**, 046127 (2005), e-print [arXiv:physics/0506074](https://arxiv.org/abs/physics/0506074); J. Sienkiewicz, J.A. Holyst, *Acta Phys. Polonica B* **36**, 1771 (2005)
23. P. Angeloudis, D. Fisk, *Physica A* **367**, 553 (2006)
24. P.-P. Zhang, K. Chen, Y. He, T. Zhou, B.-B. Su, Y. Jin, H. Chang, Y.-P. Zhou, L.-C. Sun, B.-H. Wang, D.-R. He, *Physica A* **360**, 599 (2006)
25. C. von Ferber, T. Holovatch, Yu. Holovatch, V. Palchykov, *Physica A* **380**, 585 (2007)
26. X. Xu, J. Hu, F. Liu, L. Liu, *Physica A* **374**, 441 (2007)
27. H. Chang, B.-B. Su, Y.-P. Zhou, D.-R. He, *Physica A* **383**, 687 (2007)

28. C. von Ferber, T. Holovatch, Yu. Holovatch, V. Palchykov, in *Traffic and Granular Flow '07*, edited by C. Appert-Rolland, F. Chevoir et al. (Springer, 2009), e-print [arXiv:0709.3203](https://arxiv.org/abs/0709.3203)
29. C. von Ferber, T. Holovatch, Yu. Holovatch, in *Traffic and Granular Flow '07*, edited by C. Appert-Rolland, F. Chevoir et al. (Springer, 2009), e-print [arXiv:0709.3206](https://arxiv.org/abs/0709.3206)
30. For links see <http://www.apta.com>
31. Due to an updated database numbers may slightly differ from those given in [25]
32. C. von Ferber, T. Holovatch, Yu. Holovatch, V. Palchykov, e-print [arXiv:0803.3514v1](https://arxiv.org/abs/0803.3514v1)
33. J.-L. Guillaume, M. Latapy, *Physica A* **371**, 795 (2006)
34. P. Erdős, A. Rényi, *Publ. Math. (Debrecen)* **6**, 290 (1959); P. Erdős, A. Rényi, *Publ. Math. Inst. Hung. Acad. Sci.* **5**, 17 (1960); P. Erdős, A. Rényi, *Bull. Inst. Int. Stat.* **38**, 343 (1961)
35. B. Bollobás, *Random Graphs* (Academic Press, London, 1985)
36. R. Cohen, D. ben-Avraham, S. Havlin, *Phys. Rev. E* **66**, 036113 (2002)
37. H.A. Simon, *Biometrika* **42**, 425 (1955)
38. D. de S. Price, *J. Amer. Soc. Inform. Sci.* **27**, 292 (1976)
39. A.-L. Barabási, R. Albert, *Science* **286**, 509 (1999); A.-L. Barabási, R. Albert, H. Jeong, *Physica A* **272**, 173 (1999)
40. R.F. i Cancho, R.V. Solé, e-print [arXiv:cond-mat/0111222](https://arxiv.org/abs/cond-mat/0111222); S. Valverde, R.F. i Cancho, R.V. Solé, *Europhys. Lett.* **60**, 512 (2002); R.F. i Cancho, R.V. Solé, in *Statistical mechanics of Complex Networks*, edited by R. Pastor-Satorras, M. Rubi, A. Diaz-Guilera, Lecture Notes in Physics (Springer, Berlin, 2003), Vol. 625, p. 114
41. M.T. Gastner, M.E.J. Newman, *Eur. Phys. J. B* **49**, 247 (2006)
42. N. Mathias, V. Gopal, *Phys. Rev. E* **63**, 021117 (2001)
43. R.F. i Cancho, R.V. Solé, *Proc. Natl. Acad. Sci. USA* **100**, 788 (2003); R.F. i Cancho, *Physica A* **345**, 275 (2005)
44. M.E.J. Newman, *Phys. Rev. Lett.* **89**, 208701 (2002)
45. M.E.J. Newman, *Phys. Rev. E* **67**, 026126 (2003)
46. J.A. Holyst, J. Sienkiewicz, A. Fronczak, P. Fronczak, K. Suchecki, *Phys. Rev. E* **72**, 026108 (2005)
47. A. Fronczak, P. Fronczak, J.A. Holyst, *Phys. Rev. E* **68**, 046126 (2003)
48. S.N. Dorogovtsev, J.F.F. Mendes, J.G. Oliveira, *Phys. Rev. E* **73**, 056122 (2006)
49. U. Brandes, *J. Math. Sociology* **25**, 163 (2001)
50. G. Sabidussi, *Psychometrika* **31**, 581 (1966)
51. P. Hage, F. Harary, *Social Networks* **17**, 57 (1995)
52. A. Shimbel, *Bull. Math. Biophys.* **15**, 501 (1953)
53. L.C. Freeman, *Sociometry* **40**, 35 (1977)
54. L. Benguigui, *J. Phys. I France* **2**, 385 (1992)
55. Maps provided by <http://www.fahrinfo-berlin/Stadtplan>
56. Geocoding application on <http://developer.navteq.com/>
57. B. Nienhuis, *Phys. Rev. Lett.* **49**, 1062 (1982)
58. M.E.J. Newman, S.H. Strogatz, D.J. Watts, *Phys. Rev. E* **64**, 026118 (2001)
59. Z. Liu, Y.-C. Lai, N. Ye, P. Dasgupta, *Phys. Lett. A* **303**, 337 (2002)
60. M.E.J. Newman, *Phys. Rev. E* **64**, 016131 (2001)
61. X. Li, G. Chen, *Physica A* **328**, 274 (2003)
62. J.J. Ramasco, S.N. Dorogovtsev, R. Pastor-Satorras, *Phys. Rev. E* **70**, 036106 (2004)
63. A. Cardillo, S. Scellato, V. Latora, S. Porta, *Phys. Rev. E* **73**, 066107 (2006)
64. D. Volchenkov, P. Blanchard, *Phys. Rev. E* **75**, 026104; D. Volchenkov, *Condens. Matter Phys.* **11**, 331 (2008)
65. A.B. Harris, *Z. Phys. B* **49**, 347 (1983); Y. Kim, *J. Phys. C* **16**, 1345 (1983); V. Blavats'ka, C. von Ferber, Yu. Holovatch, *Phys. Rev. E* **64**, 041102 (2001); C. von Ferber, V. Blavats'ka, R. Folk, Yu. Holovatch, *Phys. Rev. E* **70**, 035104(R) (2004)