# Toward a Measurement-based Geographic Location Service

Artur Ziviani[1,2], Serge Fdida[1], José F. de Rezende[2], and
Otto Carlos M. B. Duarte[2]

[1] Laboratoire d'Informatique de Paris 6 (LIP6)
Université Pierre et Marie Curie (Paris 6)
Paris, France
{Artur.Ziviani,Serge.Fdida}@lip6.fr

[2] Grupo de Teleinformática e Automação (GTA)
COPPE/Poli – Universidade Federal do Rio de Janeiro
Rio de Janeiro, Brazil
{rezende,otto}@gta.ufrj.br

**Abstract.** Location-aware applications require a geographic location service of Internet hosts. We focus on a measurement-based service for the geographic location of Internet hosts. Host locations are inferred by comparing delay patterns of geographically distributed landmarks, which are hosts with a known geographic location, with the delay pattern of the target host to be located. Results show a significant correlation between geographic distance and network delay that can be exploited for a coarse-grained geographic location of Internet hosts.

## 1 Introduction

Location-aware applications take into account where the users are accessing from and thereby can offer novel functionalities in the Internet. Examples of these novel location-aware applications are: targeted advertising on web pages, automatic selection of a language to first display the content, accounting the incoming users based on their positions, restricted content delivery following regional policies, and authorization of transactions only when performed from pre-established locations. In peer-to-peer networks, location-aware construction of overlay networks can avoid unnecessary high latency hops, thus improving routing performance [1]. Multimedia delivery systems, such as Content Distribution Networks (CDNs), can also benefit from knowing the location of their clients [2]. For example, benefits include the indication of close servers to clients or the adaptation of multimedia content based on the location of clients. Lakhina *et al.* [3] investigate the geographic location of Internet components to create a base for the development of geographically-driven topology generation methods. In the current Internet, however, there is no direct relationship between a host identification and its physical location. The novel location-aware applications then require the deployment of a geographic location service for Internet hosts.

We focus on a measurement-based geographic location service of Internet hosts to support location-aware applications. We build upon GeoPing [4] that adopts an empirical approach based on the observation that hosts sharing similar delays to other fixed hosts tend to be near each other geographically. In a previous work [5], we have evaluated different similarity models to compare the delay patterns gathered from different reference hosts. In this paper, we carry out live experiments to evaluate the correlation between geographic distance and network delay as well as the achieved distance accuracy. Our findings indicate that contrary to conventional wisdom there is a significant level of correlation between distance and delay. This correlation becomes stronger as connectivity within the network becomes richer. Moreover, such a correlation can be exploited to perform a coarse-grained geographic location of Internet hosts.

This paper is organized as follows. Section 2 briefly reviews schemes to determine the geographic location of Internet hosts. In Section 3, we formalize the measurement-based geographic location of Internet hosts. We introduce the adopted measures of similarity in Section 4. Section 5 presents our experiments and results. Finally, in Section 6 we conclude and discuss future work.

## 2 Related Work

A DNS-based approach to provide a geographic location service of Internet hosts is proposed in RFC 1876 [6]. Nevertheless, the adoption of the DNS-based approach is restricted since it requires changes in the DNS records and administrators have no motivation to register new location records. Tools such as IP2LL [7] and NetGeo [8] query Whois databases in order to obtain the location information recorded therein to infer the geographic location of a host. Such an information, however, may be inaccurate and stale. Moreover, if a large and geographically dispersed block of IP addresses is allocated to a single entity, the Whois databases may contain just a single entry for the entire block. As a consequence, a query onto the Whois databases provides the registered location of the entity that controls the block of IP addresses, although the concerned hosts may be geographically dispersed.

Padmanabhan and Subramanian [4] investigate three important techniques to infer the geographic location of an Internet host. The first technique infers the location of a host based on the DNS name of the host or another nearby node. This technique is the base of GeoTrack [4], VisualRoute [9], and GTrace [10]. Quite often network operators assign names to routers that have some geographic meaning, presumably for administrative convenience. For example, the name `bcr1-so-2-0-0.Paris.cw.net` indicates a router located in Paris, France. Nevertheless, not all names contain an indication of location. Since there is no standard, operators commonly develop their own rules for naming their routers even if the names are geographically meaningful. Therefore, the parsing rules to recognize a location from a node name must be specific to each operator. The creation and management of such rules is a challenging task as there is no standard to follow. As the position of the last recognizable router in the path toward

the host to be located is used to estimate the position of such a host, a lack of accuracy is also expected. The second technique splits the IP address space into clusters such that all hosts with an IP address within a cluster are likely to be co-located. Knowing the location of some hosts in the cluster and assuming they are in agreement, the technique infers the location of the entire cluster. An example of such a technique is GeoCluster [4]. This technique, however, relies on information that is partial and possibly inaccurate. The information is partial because it comprises location information for a relatively small subset of the IP address space. Moreover, such an information may be inaccurate because the databases rely on data provided by users, which may be unreliable to provide correct location information. The third technique is based on delay measurements and the exploitation of a possible correlation between geographic distance and network delay. Such a technique is the base of GeoPing [4]. The location estimation of a host is based on the assumption that hosts with similar network delays to some fixed probe machines tend to be located near each other. Therefore, given a set of landmarks with a well known geographic location, the location estimation for a target host to be located is the location of the landmark presenting the most similar delay pattern to the one observed for the target host.

## 3 Measurement-based Geographic Location Service

We formalize the problem of inferring a host location from delay measurements as follows. Consider a set $\mathcal{L} = \{L_1, L_2, \ldots, L_K\}$ of $K$ landmarks. Landmarks are reference hosts with a well known geographic location. Consider a set $\mathcal{P} = \{P_1, P_2, \ldots, P_N\}$ of $N$ probe machines. Fig. 1 illustrates the steps in inferring a host location from delay measurements, which are detailed along this section. The probe machines periodically determine the network delay, which is actually the minimum RTT of several measurements, to each landmark (Fig. 3). Therefore, each probe machine $P_x \in \mathcal{P}$ keeps a delay vector $\mathbf{d}_x = [d_{1x}, d_{2x}, \ldots, d_{Kx}]^T$, where $d_{ix}$ is the delay between the probe machine $P_x$ and the landmark $L_i \in \mathcal{L}$. Suppose one wants to determine the geographic location of a given target host $T$. A location server that knows the landmark set $\mathcal{L}$ and the probe machine set $\mathcal{P}$ is then contacted. The location server asks the $N$ probe machines to measure the delay to host $T$ (Fig. 3). Each probe machine $P_x \in \mathcal{P}$ returns to the location server a delay vector $\mathbf{d}'_x = [d_{1x}, d_{2x}, \ldots, d_{Kx}, d_{Tx}]^T$, i.e., the delay vector $\mathbf{d}_x$ plus the just measured delay to host $T$ (Fig. 3). After receiving the delay vectors from the $N$ probe machines, the location server is able to construct a delay matrix $\mathbf{D}$ with dimensions $(K+1) \times N$:

$$\mathbf{D} = \begin{bmatrix} d_{11} & d_{12} & \ldots & d_{1N} \\ d_{21} & d_{22} & \ldots & d_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ d_{K1} & d_{K2} & \ldots & d_{KN} \\ d_{T1} & d_{T2} & \ldots & d_{TN} \end{bmatrix} \tag{1}$$
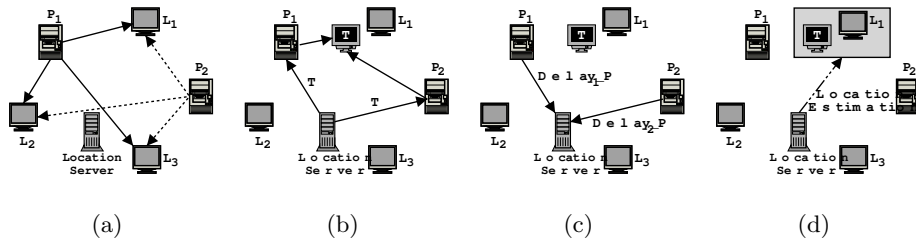
**Fig. 1.** Inferring a host location from delay measurements.

The delay vectors gathered by the demanding location server from the probe machines correspond to the columns of the delay matrix $\mathbf{D}$. The location server then compares the lines of the delay matrix $\mathbf{D}$ to estimate the location of host $T$. The delay matrix $\mathbf{D}$ combined with the knowledge of the location of the landmarks of the set $\mathcal{L}$ compose a delay map recording the relationship between network delay and geographic location.

## 4 Measuring the Similarity between Delay Patterns

In this section, we investigate how to best measure the similarity between the delay pattern of each landmark and the one observed for the target host. The delay patterns result from the partial viewpoints gathered by the distributed probe machines. The landmark that presents the most similar delay pattern with respect to the one of the target host provides the location estimation of that host. Measuring the similarity of the concerned delay patterns is thus a key point for the accuracy of the host location from delay measurements.

The function $\mathcal{S}(\mathbf{x}, \mathbf{y})$ is defined to measure the degree of dissimilarity between two delay patterns $\mathbf{x}$ and $\mathbf{y}$ of size $N$, where $N$ is the number of adopted probe machines. These delay patterns are gathered by the probe machines from each landmark and from the target host to be located. To formalize the dissimilarity evaluation, we also define a line vector $\mathbf{1}_i$ of size $K + 1$ that has all elements equal to 0, except for the $i^{th}$ element that has a value of 1. The landmark $L$ that provides the location estimation of the target host $T$ is the landmark that gives the minimum dissimilarity

$$\mathcal{S}_{\min} = \arg \min_{i=1,\dots,K} \mathcal{S}(\mathbf{1}_i \mathbf{D}, \mathbf{1}_{K+1} \mathbf{D}). \tag{2}$$

We first consider distance-based measures of dissimilarity [11] to compare delay patterns. The generalized form to represent a distance metric is given by

$$\mathcal{S}_\gamma(\mathbf{x}, \mathbf{y}) = \left[ \sum_{i=1}^{N} |x_i - y_i|^\gamma \right]^{\frac{1}{\gamma}}, \quad \gamma > 0. \tag{3}$$

When $\gamma = 1$, we have the Manhattan or city-block distance. In contrast, for $\gamma = 2$, we have the Euclidean distance. It is shown that the Chebyshev distance, i.e. $\gamma = \infty$, can be expressed as

$$\mathcal{S}_\infty(\mathbf{x}, \mathbf{y}) = \lim_{\gamma \to \infty} \mathcal{S}_\gamma(\mathbf{x}, \mathbf{y}) = \max_i |x_i - y_i|.$$

We also consider the Canberra distance given by

$$\mathcal{S}_{\text{canb}}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{N} \frac{|x_i - y_i|}{x_i + y_i}. \tag{4}$$

If both $x_i$ and $y_i$ are zero the ratio of the difference to the sum is taken to be zero. The Canberra distance is suitable for variables taking non-negative values and is sensitive to small changes close to zero values [11].

The two delay patterns $\mathbf{x}$ and $\mathbf{y}$ can also be thought of as two vectors in a $N$-dimensional delay space. The similarity $\mathcal{S}_{\cos}(\mathbf{x}, \mathbf{y})$ between them is measured by computing the cosine of the angle $\theta$ between these two vectors. The cosine of the angle $\theta$ between the delay vectors $\mathbf{x}$ and $\mathbf{y}$ is computed by

$$\cos\theta = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\|\|\mathbf{y}\|} \tag{5}$$

where "$\cdot$" denotes the dot-product of the two vectors and $\|\mathbf{x}\|$ is the Euclidean size of vector $\mathbf{x} \in \mathbb{R}^N$, i.e. $\|\mathbf{x}\| = \sqrt{\sum_{i=1}^{N} x_i^2}$.

An alternative measure of similarity is to compute the coefficient of correlation between the two delay patterns $\mathbf{x}$ and $\mathbf{y}$. This correlation-based similarity model is denoted by $\mathcal{S}_{\text{cor}}(\mathbf{x}, \mathbf{y})$. The coefficient of correlation is defined as

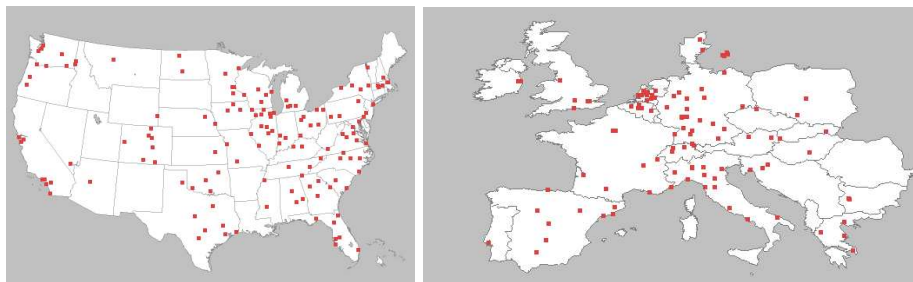$$\text{corr}(\mathbf{x}, \mathbf{y}) = \frac{\sigma_{\mathbf{xy}}^2}{\sigma_{\mathbf{x}} \sigma_{\mathbf{y}}}, \tag{6}$$

where $\sigma_{\mathbf{xy}}^2$ denotes the covariance between delay patterns $\mathbf{x}$ and $\mathbf{y}$, and $\sigma_{\mathbf{x}}$ is the standard deviation of $\mathbf{x}$.

## 5 Experimental Results

In this paper, we analyze the results of live experiments to evaluate basic properties of a measurement-based geographic location service of Internet hosts.

### 5.1 Experimental Setup

For the experiments, we use 9 measurement boxes from the NIMI (National Internet Measurement Infrastructure) project [12] as our probe machines. They are geographically distributed as follows: 5 in Western Europe, 3 in the U.S., and 1 in Japan. Recent works [4,13] indicate that 7 to 9 dimensions provide sufficient network distance representation. The experimental set of landmarks comes from two datasets:

(a) USA                                      (b) Western Europe

**Fig. 2.** Geographic location of landmarks.

- LibWeb – this set of hosts is mainly composed of university sites extracted from library web (LibWeb) servers around the world [14].
- RIPE – these hosts are part of the Test Traffic Measurements (TTM) project of the RIPE network [15]. All hosts on the RIPE network are equipped with a GPS card, thus allowing their exact geographic position to be known.

The resulting experimental dataset totals to 397 landmarks that are sparsely distributed worldwide. The geographic distribution of these landmarks is as follows: 199 in North America (U.S. and Canada), 156 in Western Europe, 19 in Eastern Europe, 13 in Latin America, 9 in Asia and Oceania, and 1 in the Middle East. This distribution is intended to at least roughly reflect the distribution of users (hosts) to be located. In a previous work [16], we propose the demographic placement of landmarks to better represent the location of users (hosts). It should be noted that landmarks are unsuspecting participants in the procedure since a landmark may be any host, with a known geographic location, able to echo `ping` messages. Figure 2 shows the geographic location of the landmarks in the U.S. and in Western Europe, which are regions likely to have rich connectivity and host most users to be located.

The probe machines measure the delay toward the set of landmarks. The delay metric is actually the minimum of several RTT measurements to avoid taking into account congestion delays. The measurements toward each landmark from the different probe machines are enough spaced to avoid triggering detection systems against DDoS (Distributed Denial of Service) attacks.

### 5.2 Correlation between Geographic Distance and Network Delay

In this section, we evaluate the correlation between geographic distance and network delay. Until recently, common sense claimed that there is a weak correlation between distance and delay within the network. Claffy [17] mentions this
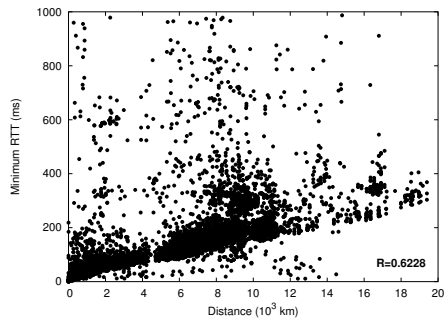
as one of some myths about Internet data. Our results show that few landmarks present very large delays, probably due to poor connectivity. To avoid taking into account these outliers in our evaluation, we consider data within the $98^{th}$, $95^{th}$, and $90^{th}$ percentiles of the measured network delay. The observed correlation between geographic distance and network delay is moderate to strong in these cases, resulting in R=0.6229, R=0.8093, and R=0.8767, respectively. Figures 5.2, 5.2, and 5.2 present the corresponding scatter plots for these results that cover landmarks located worldwide.

We also observe that poor connectivity weakens the correlation between geographic distance and network delay. We then identify the landmarks located in North America and Western Europe. These regions are likely to offer the richest connectivity linking their hosts. We observe an even stronger correlation on these well connected regions, indicating that the correlation becomes stronger as connectivity becomes richer. The coefficients of correlation for the data within the $98^{th}$, $95^{th}$, and $90^{th}$ percentiles in these well connected regions are R=0.7126, R=0.8595, and R=0.8958, respectively. The corresponding scatter plots for these results are shown in Figures 5.2, 5.2, and 5.2 for landmarks located in North America and Western Europe (NA-WE). Recent findings [3,18] indicate a strong correlation between population and router density in economically developed countries. Moreover, most users, and consequently most hosts to be located, are likely to be in these regions with richer connectivity, whereby a stronger correlation between distance and delay holds.
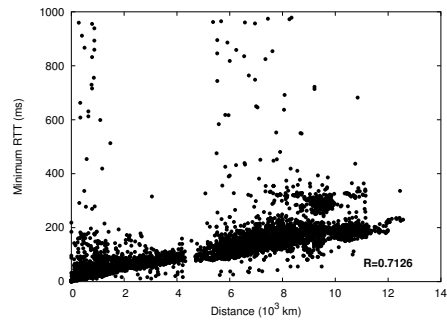
### 5.3 Distance Accuracy

In this section, we consider the whole set of landmarks distributed worldwide, in a total of 397 landmarks. To evaluate the distance accuracy, we take one landmark as a target host and use the remaining landmarks to infer a location estimation for this target. The distance accuracy is measured by the error distance from the location estimation to the location of the target host. We apply the different measures of dissimilarity presented in Section 3 to compare the delay patterns gathered by the probe machines from the landmarks.
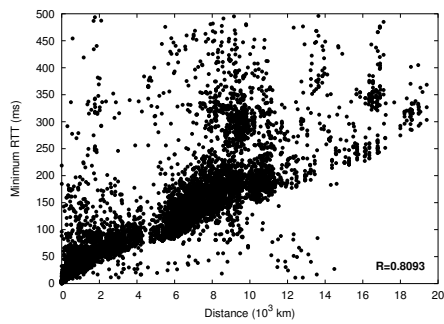
Figure 4 shows the probability density function (pdf) of the error distance worldwide for the different measures of dissimilarity. The Canberra distance performs slightly better than the others, providing smaller error distances to more hosts. This distance measure is known to be suitable for non-negative values, such as network delay, and more sensitive to values near zero. This favors a more accurate location of some hosts in comparison with the other measures of dissimilarity since eight out of the nine probe machines are in the U.S. or in Western Europe. For the Canberra distance, we observe that the median value of the error distance is 314 km with a kurtosis of 40.79, showing that the observed distribution of the error distance is heavy-tailed. This is because, for some target hosts, even if the elected landmark is the geographically closest landmark to the target, not necessarily it is nearby the target host. These results indicate that delay measurements can indeed be exploited to determine the geographic location of Internet hosts, although at a coarse granularity.
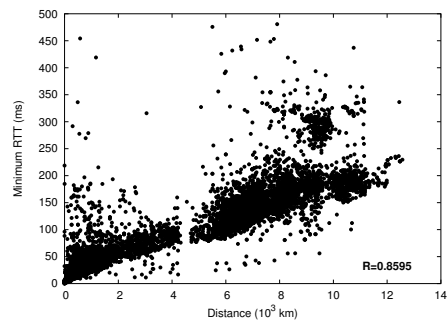
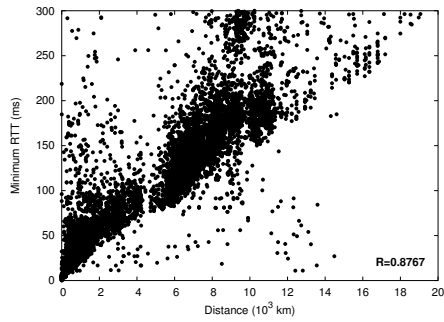(a) Worldwide − $98^{th}$ delay percentile
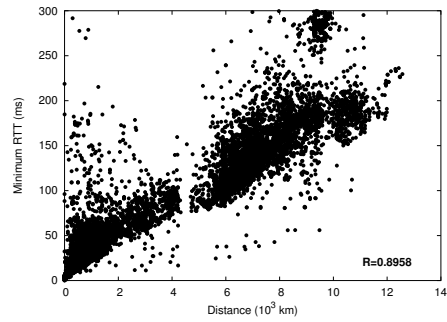
(b) NA-WE − $98^{th}$ delay percentile

(c) Worldwide − $95^{th}$ delay percentile
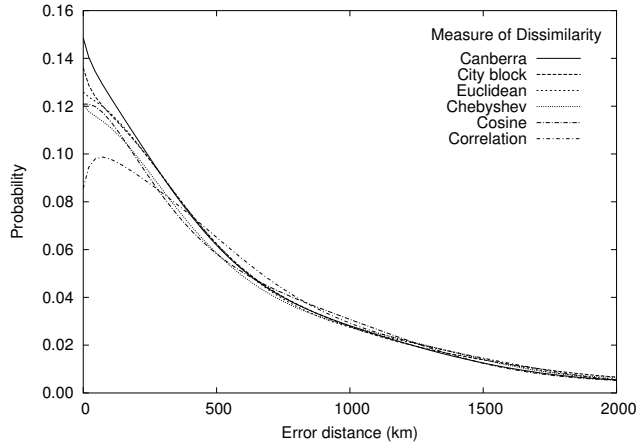
(d) NA-WE − $95^{th}$ delay percentile

(e) Worldwide − $90^{th}$ delay percentile

(f) NA-WE − $90^{th}$ delay percentile

**Fig. 3.** Correlation between geographic distance and network delay, both worldwide and within the North America and Western Europe (NA-WE) regions.

**Fig. 4.** PDF of the error distance.

## 6 Conclusion

This paper investigates some key properties toward a measurement-based geographic location service of Internet hosts. Such a service can be viewed as an underlying infrastructure for the deployment of novel location-aware applications in the Internet. Live experiments have been carried out to evaluate the correlation between geographic distance and network delay as well as the achieved distance accuracy for different measures of dissimilarity. Our findings indicate that contrary to conventional wisdom there is a significant correlation between geographic distance and network delay. We show that this correlation can be exploited to provide a coarse-grained geographic location of Internet hosts. The location estimation of a host is the location of the landmark presenting the most similar delay pattern with respect to the target host. This poses a fundamental limit: the system has a discrete space of answers since the number of possible answers correspond to the number of landmarks adopted.

As future work, we intend to investigate methods to adopt a continuous space of answers instead of a discrete one. Recent works [13,19,20] propose to infer network proximity without direct measurements by embedding network distances such as delay into a coordinate system of reduced dimensions. Similar concepts can be applied to the measurement-based geographic location of Internet hosts to provide more accurate estimations using fewer measurements.

## Acknowledgment

# References

1. Ratnasamy, S., Handley, M., Karp, R., Shenker, S.: Topologically-aware overlay construction and server selection. In: Proc. of the IEEE INFOCOM'2002, New York, NY, USA (2002)
2. Sripanidkulchai, K., Maggs, B., Zhang, H.: Efficient content location using interest-based locality in peer-to-peer systems. In: Proc. of the IEEE INFOCOM'2003, San Francisco, CA, USA (2003)
3. Lakhina, A., Byers, J.W., Crovella, M., Matta, I.: On the geographic location of Internet resources. IEEE Journal on Selected Areas in Communications **21** (2003) 934–948
4. Padmanabhan, V.N., Subramanian, L.: An investigation of geographic mapping techniques for Internet hosts. In: Proc. of the ACM SIGCOMM'2001, San Diego, CA, USA (2001)
5. Ziviani, A., Fdida, S., de Rezende, J.F., Duarte, O.C.M.B.: Similarity models for Internet host location. In: Proc. of the IEEE International Conference on Networks - ICON'2003, Sydney, Australia (2003) 81–86
6. Davis, C., Vixie, P., Goowin, T., Dickinson, I.: A means for expressing location information in the domain name system. Internet RFC 1876 (1996)
7. University of Illinois at Urbana-Champaign: (IP Address to Latitude/Longitude) http://cello.cs.uiuc.edu/cgi-bin/slamm/ip2ll/.
8. Moore, D., Periakaruppan, R., Donohoe, J., Claffy, K.: Where in the world is netgeo.caida.org? In: Proc. of the INET'2000, Yokohama, Japan (2000)
9. Visualware Inc.: (VisualRoute) http://www.visualware.com/visualroute/.
10. CAIDA: (GTrace) http://www.caida.org/tools/visualization/gtrace/.
11. Gordon, A.D.: Classification: Methods for the Exploratory Analysis of Multivariate Data. Chapman and Hall (1981)
12. Paxson, V., Mahdavi, J., Adams, A., Mathis, M.: An architecture for large-scale Internet measurement. IEEE Communications Magazine **36** (1998) 48–54
13. Tang, L., Crovella, M.: Virtual landmarks for the Internet. In: ACM Internet Measurement Conference 2003, Miami, FL, USA (2003)
14. : (LibWeb) http://sunsite.berkeley.edu/Libweb.
15. : (RIPE Test Traffic Measurements) http://www.ripe.net/ttm/.
16. Ziviani, A., Fdida, S., de Rezende, J.F., Duarte, O.C.M.B.: Demographic placement for Internet host location. In: Proc. of the IEEE Global Communications Conference - GLOBECOM'2003, San Francisco, CA, USA (2003)
17. Claffy, K.: Internet measurement: myths about Internet data. Talk at NANOG24 Meeting (2002) http://www.caida.org/outreach/presentations/Myths2002/.
18. Yook, S.H., Jeong, H., Barabási, A.L.: Modeling the Internet's large-scale topology. Proc. of the National Academy of Sciences (PNAS) **99** (2002) 13382–13386
19. Ng, T.S.E., Zhang, H.: Predicting Internet network distance with coordinates-based approaches. In: Proc. of the IEEE INFOCOM'2002, New York, NY, USA (2002)
20. Lim, H., Hou, J.C., Choi, C.H.: Constructing Internet coordinate system based on delay measurement. In: ACM Internet Measurement Conference 2003, Miami, FL, USA (2003)