# Development of map matching algorithm for low frequency probe data

Tomio Miwa [a,*], Daisuke Kiuchi [b], Toshiyuki Yamamoto [a], Takayuki Morikawa [c]

[a] EcoTopia Science Institute, Nagoya University, Furo-cho, Chikusa-ku, Nagoya 464-8603, Japan
[b] Natural Gas Business Division, Mitsubishi Corporation, 3-1, Marunouchi 2-Chome, Chiyoda-ku, Tokyo 100-8086, Japan
[c] Graduate School of Environmental Studies, Nagoya University, Furo-cho, Chikusa-ku, Nagoya 464-8603, Japan

### A R T I C L E   I N F O

### A B S T R A C T

In order to lower the operating costs of a large scale probe vehicle system, countermeasures for decreasing operating cost of the system, such as lowering of data polling frequency and use an existing fleet management system, are necessary. Such countermeasures, however, reduce the accuracy of the traffic information that is generated from the collected probe vehicle data. In this study, the authors developed several map matching algorithms that can be applied to low frequency and little information probe vehicle data. These map matching algorithms were verified using actual probe vehicle data collected in the area around Nagoya, Japan. The results show that the data can be map matched with a high degree of accuracy by combining an appropriate link cost, generation of reasonable candidate routes, evaluation of the routes, and introducing the concept of a driver's route choice.

© 2012 Elsevier Ltd. All rights reserved.

## 1. Introduction

Probe vehicle systems that make use of vehicles traveling through the traffic network can collect traffic data more cheaply and cover a wider area than can stationary roadside equipment, such as loop detectors and ultrasonic detectors. Therefore, in recent years, attempts have been made in many parts of the world to collect traffic data using probe vehicles (Sermons and Koppelman, 1996; Quiroga and Bullock, 1998; Murakami and Wagner, 1999; Hellinga and Fu, 2002; Hellinga et al., 2008). The data obtained from probe vehicles consists of consecutive vehicle position coordinates. Traffic information is generated from this data using a process called "map matching," which derives the cruising route of the probe vehicle and the vehicle's position on the road. Since the accuracy of this process directly influences the accuracy of the generated traffic information, the agent of the traffic information service must develop a map matching algorithm with the highest possible accuracy. Moreover, since a probe vehicle system can provide data only for those sections of the road that are actually traveled, it is necessary to organize a large number of probe vehicles in order to collect traffic data from widespread sections of road. Additionally, a map matching algorithm based on simple procedure is needed for processing the huge amount of data from such a large scale probe vehicle system.

A large scale probe vehicle system might impose an enormous financial burden on the system administration. For this reason, measures to decrease operating costs, such as by reducing the data polling frequency and limiting the data size (deleting some information items), are required (Liu et al., 2007a). The resulting traffic information, however, will be less accurate. If the polling frequency is low, for example, 1 min or 500 m intervals, the probe vehicle's cruising route will be difficult to identify, especially in areas with a highly dense road network (Liu et al., 2006). To maintain the required level of accuracy in the generated traffic information, a map matching algorithm that can accurately process the low polling frequency and little information probe data is needed. Such a map matching algorithm might also make use of data from

---

* Corresponding author. Tel.: +81 52 789 5018; fax: +81 52 789 5728.
  *E-mail address:* miwa@nagoya-u.jp (T. Miwa).

existing fleet management systems, such as taxi and truck management systems, as a by-product. Although this type of fleet management system includes little information (only vehicle position and a time stamp) and has a low polling frequency (an interval of 1 min or more or a few hundred meters), a map matching algorithm that can deal with it would help to realize a probe vehicle system of overwhelming superiority.

In this study, the authors attempt to develop a map matching algorithm for probe vehicle data that has a low polling frequency and little information.

## 2. Review of existing map matching algorithms

Generally, "map matching" means either estimating the vehicle's state (position in the road space) based on detailed information such as vehicle heading from measuring instruments (for example, Kim and Kim, 2001), or guessing the route and the vehicle's position along the route based on relatively simple information, such as GPS data (for example, White et al., 2000). This paper focuses on the latter and defines "map matching" as the latter technique.

### 2.1. Classification of map matching algorithms

Many map matching algorithms for probe vehicle data have been proposed over the last 10 years. These can be classified into two major types according to the processing frame. The first type processes probe data after the trip is over and finds the overall route (for example, Yin and Wolfson, 2005), while the second type processes data during a trip (for example, White et al., 2000). We call the former an "off-line process" and the latter an "on-line process."

An off-line algorithm processes either data that has been divided into sections corresponding to significant changes in the probe vehicle's direction of travel (Makimura et al., 2002) or data for individual trips (Miwa et al., 2004). Since data for relatively long sections is available, more information about the route is available and GPS errors have a smaller impact on map matching accuracy. The process is cumbersome, however, because the data must be divided and map matched repeatedly. Moreover, when a huge amount of accumulated data is processed, processing costs become high. On the other hand, on-line map matching involves the processing of small amounts of data (only the data being gathered) to identify vehicle route and position on the road. In this case, GPS errors may affect the accuracy of the map matching results, but because processing is repeated each time new data is collected, the procedure is concise. Since our goal is to process the huge amount of data gathered by a large scale probe vehicle system, in this work we attempt to develop a new on-line algorithm.

### 2.2. Existing map matching algorithms

Several on-line map matching algorithms have been proposed over the last 10 years, arguing the need for a map matching method applicable to low polling frequency and/or low accuracy data and for an algorithm that does not depend on an in-vehicle system (Hellinga et al., 2003). With the prospect that our developed algorithm will be applied to huge volumes of low frequency and little information probe vehicle data, the assumption is that it should be executed by a simple procedure using only the vehicle's position coordinates and polling time stamps in urban areas with a dense road network. This subsection reviews some on-line map matching algorithms and the problems encountered when processing low frequency and little information probe data in large scale system.

White et al. (2000) pointed out the problems of easily matching a GPS point (hereinafter called a "plot") to the nearest node or shape point in a network and proposed some algorithms for matching a plot to the appropriate link. However, the proposed algorithms utilize the information of vehicle heading, and high frequency data for each link on the route is assumed to include at least one plot that can be used. Tradisauskas et al. (2009) developed an on-line map matching algorithm for an intelligent speed adaptation system and proposed seven weighting functions for evaluating candidate links. However, the proposed algorithm utilizes odometer speed information and targets high frequency data. Greenfeld (2002) proposed a method that does not use information about vehicle heading and speed, but uses only the positional relationships between plots and links. This method identifies a vehicle's position on the link by evaluating the angle and distance between a link and a line connecting two plots. However, in this method, high polling frequency data is assumed available. Yang et al. (2005) proposed an algorithm that sets a "node buffer," which is a certain constant distance from an intersection node, and matches a plot near the intersection node to the node itself so as to prevent mis-matching around the intersection. In the algorithm, plots not near the intersection are matched to the nearest link. If a plot has several candidate links nearby, matching of the plot is suspended, and it is matched to the candidate link constituting the shortest route between the link matched to the previous plot and the link matched to the subsequent plot. This algorithm immediately matches the collected data to the road network and produces few matching errors around intersections. However, the cumbersome procedure repeatedly segments the data and must search for routes when the urban road network has many candidate links. Moreover, the plots matched to a node buffer must again be matched to the route so that the link travel time can be accurately calculated. Brakatsoulas et al. (2005) developed an algorithm for relatively low frequency data (polled every 30 s). They developed on-line and off-line versions and compared them. The on-line algorithm utilizes an error ellipse formed from the measurement error and sampling error to search for candidate links. Additionally, they introduced a "look-ahead" technique for dealing with link skipping situations and to improve matching accuracy. However, the process is not simple and the look-ahead

technique may cause processing difficulties with lower frequency data in urban areas with a dense road network. Wenk et al. (2006) attempted to localize (process incrementally) the off-line matching method developed by Brakatsoulas et al. (2005). The algorithm utilizes an error ellipse and determines a matching route based on the weak Fréchet distance. The results show that incremental processing reduces the map matching accuracy. Lou et al. (2009) developed an algorithm specialized for low frequency data. The developed algorithm is not specifically for on-line processing but may be applied. This algorithm retrieves candidate links for each plot and connects two candidate links by minimum distance route. After that, it finds the most reasonable sequence of the candidate routes. In the algorithm, since the minimum distance route is applied, if the distance between consecutive plots is long, the route between plots may not be evaluated adequately. However, Lou et al. presented two important concepts. The first one is consideration of anteroposterior relationship on the plot's trajectory, and the second is consideration of road speed constraints. These two concepts should be considered in our study.

## 3. Data and precision indices for map matching

### 3.1. Data

The probe vehicle data used in this study was collected as part of the Nagoya Probe Taxi Project (October 2002–March 2003) (Yamamoto, 2005). This project was set up to develop an information technology system infrastructure. It involved experiments using more than 1500 taxis in cooperation with 41 member companies of the Nagoya Taxi Association. This study uses data with high frequency polling (5 s or 50 m). Note that many past investigations took 1-s data to be high frequency data. Due to the lack of 1-s data, data with a 5-s or 50-m interval is treated as high frequency data in this study.

Data considered adequate for the purposes of this study was selected and manually matched to actual (true) routes. Low frequency data was created by thinning out the plots. Table 1 gives an outline of the data used. The criteria for selecting adequate data are (1) at least one use of an expressway, (2) data does not include too many outliers, and (3) data includes meandering drive without passenger. Table 2 shows the types of roads involved and their percentages of the distance traveled in each data set. The percentage for expressways and urban expressways is much higher in data set 2 than in data set 1.

From the distance-based high frequency data, 10 kinds of data of different frequency were created by taking data points in 50 m increments from 50 m to 500 m (50, 100, 150, 200, 250, 300, 350, 400, 450, and 500 m). In the same manner, from the time-based high frequency data, 10 kinds of data were created by taking data points at 5 s or 15 s increments from 5 s to 90 s (5, 10, 15, 20, 25, 30, 45, 60, 75, and 90 s). Of course, it would be possible to create data with lower frequency than these 10 sets. However, Liu et al. (2007b) showed that even if a cruising route can be matched correctly, the accuracy of the generated

**Table 1**
Data used in study.

| Data set no. | Vehicle no. | Polling interval | Date | Time | Distance traveled (km) |
|---|---|---|---|---|---|
| 1 | 1 | 50 m | 2002/11/15 | 8:00–17:00 | 102.3 |
| | 2 | 50 m | 2002/10/31 | 13:00–23:00 | 97.1 |
| 2 | 3 | 5 s | 2002/12/25 | 11:00–16:00 | 70.1 |
| | 4 | 5 s | 2002/10/29 | 19:00–1:00 | 60.2 |

**Table 2**
Road type percentages in data.

| Data set no. | Expressways and urban expressways (%) | National roads (%) | Prefectural roads (%) | Other roads (%) |
|---|---|---|---|---|
| 1 | 2 | 12 | 66 | 20 |
| 2 | 17 | 9 | 59 | 15 |

**Table 3**
Data frequency of created data.

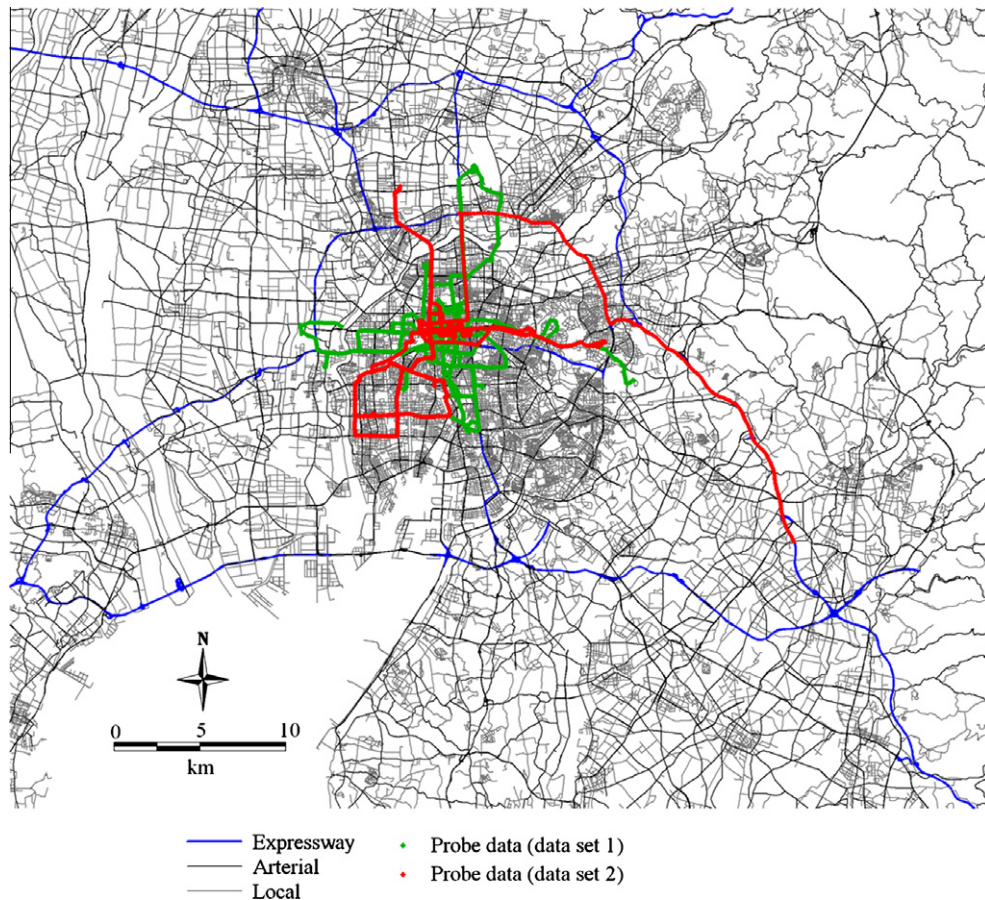| Data set no. | Polling interval | Distance interval (m) | | Time interval (s) | |
|---|---|---|---|---|---|
| | | Mean | SD | Mean | SD |
| 1 | 100 m | 115.6 | 12.6 | 25.8 | 31.9 |
| | 200 m | 224.1 | 14.2 | 50.0 | 46.1 |
| | 450 m | 476.6 | 17.1 | 106.0 | 72.4 |
| 2 | 20 s | 105.2 | 114.9 | 20.1 | 2.8 |
| | 45 s | 236.8 | 235.2 | 45.2 | 3.8 |
| | 90 s | 471.9 | 420.5 | 90.5 | 5.5 |

**Fig. 1.** Road network and probe vehicle data.

**Table 4**
Road densities and usage of areas by probe vehicles.

| | Area (km$^2$) | Total road length ($\geqslant$5.5 m width) (km) | Road density per unit area (km) | Usage by probe vehicles (%) |
|---|---|---|---|---|
| Whole study area | 2634 | 21,446 | 8 | 100 |
| Nagoya City area | 327 | 6861 | 21 | 96 |
| Central urban area | 26 | 720 | 28 | 49 |

travel time information becomes quite low for low frequency data, e.g. with a 60-s interval. Therefore, this study does not investigate lower frequency data. Table 3 gives statistics for several of these created data. Because the selected plots come after the predetermined polling interval, the polling intervals of the created data are somewhat longer than the predetermined intervals. Note that the standard deviation of time interval in distance-based data is larger than that of time-based data, vice versa.

Road network data used in this study covers the 50 km $\times$ 50 km area around Nagoya and consists of 121,290 links and 45,861 nodes. This is the network of links that are more than 5.5 m wide. Fig. 1 shows the whole network and the probe vehicle data used in this study while Table 4 summarizes the road densities. The figure clearly shows that probe vehicles, especially those in data set 1, cruise the central urban area most of the time.

### 3.2. GPS error distribution and movement threshold

To determine whether a probe vehicle is stationary or moving, the distribution of GPS errors should be known. Fig. 2 shows the cumulative distribution of the orthogonal distance from the plot to a link on the correct route. The figure shows
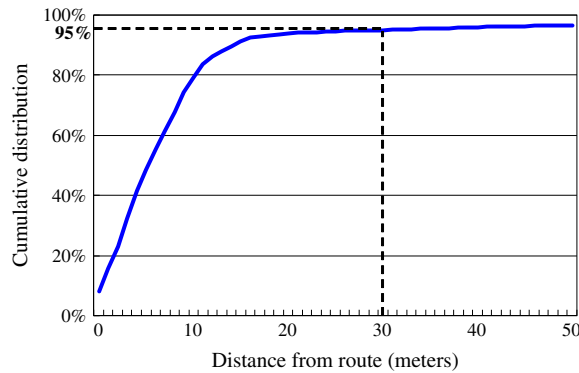
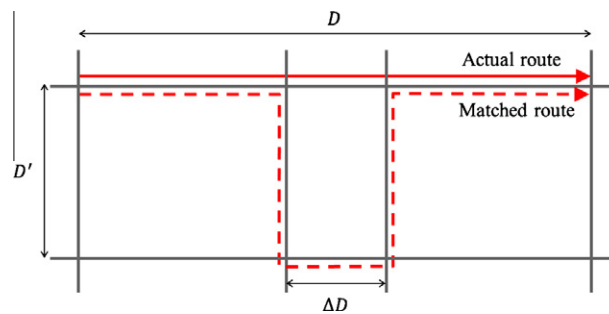**Fig. 2.** Cumulative distribution of orthogonal distance from plot to link.



**Fig. 3.** Illustration of map matching accuracy indices.

that 95% of plots fall within 30 m of the route. This means that if the distance between two consecutive plots is greater than 30 m, there is a 95% probability that the vehicle has moved. Therefore, we set the movement threshold value to 30 m. It should be noted that, since we do not consider road width and the error distribution is calculated from the orthogonal distance from a plot to the route, the data in Fig. 2 and the movement threshold value are not strictly identical to the GPS error.

### 3.3. Indices of map matching accuracy

The most commonly used index for expressing map matching accuracy is the ratio of plots correctly matched to links. This is the ARP index (accuracy ratio of plot matched). Note that with low frequency probe vehicle data, when the route between two consecutive plots is incorrect, the derived traffic information is also incorrect, even if the plots are matched to the correct links. Therefore, the two accuracy indexes, ARR (accuracy ratio of length of route identified) and IARR (inaccuracy ratio of length of route identified), are defined as follows:

$$ARR = \text{length of correctly matched route/total length of correct route} \tag{1}$$

$$IARR = \text{length of incorrectly matched route/total length of matched route} \tag{2}$$

Note that ARR is 1.0 if the matched route includes all of the links of the correct route, even if incorrect links are included in the matched route. On the other hand, IARR is the ratio of the length of incorrect links to the matched route. Therefore, if the correct route is perfectly matched, ARR is 1.0 and IARR is 0.0. Fig. 3 shows an example of the relationship between ARR and IARR. If $\Delta D$ closes to 0 in the figure, ARR closes to 1.0 (perfect). However, IARR does not close to 0 (perfect) on the condition that $D'$ exists.

## 4. Development of on-line map matching algorithm

This section describes the base algorithm. The problems encountered when applying the base algorithm to low frequency and little information probe vehicle data are also highlighted. Finally, improved algorithms with better map matching accuracy are described.
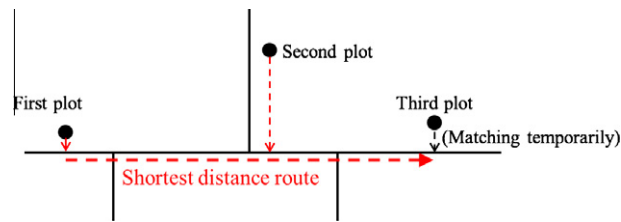
**Fig. 4.** An example of base algorithm.

### 4.1. Base algorithm

The base algorithm, which is outlined below, is almost the same as the algorithm proposed by Kojima and Hato (2004) and can be executed using a very simple procedure.

*Step 1*: Match the first plot to the nearest link. The link to which the first plot is matched is called the "first link."
*Step 2*: Obtain the next plot for which the distance from the first plot is greater than the movement threshold (30 m). This is called the "second plot." Any plots between the first and second plots are called "sub-second plots," and their distances from the first plot are shorter than the movement threshold.
*Step 3*: Obtain the next plot for which distance from the second plot is larger than the movement threshold. This is called the "third plot." Plots between the second plot and third plot are called "sub-third plots."
*Step 4*: Temporarily match the third plot to the nearest link whose direction is consistent with that from the second plot to the third plot. This is called the "third link."
*Step 5*: Search for the route with the shortest distance from the first link to the third link. Match the second and sub-second plots to the nearest links on the route.
*Step 6:* Replace the first link with the link matched to the second plot in Step 5. Replace the second and sub-second plots with the third and sub-third plots, respectively. Go to Step 3.

The difference between the above algorithm and that of Kojima and Hato (2004) is the introduction of the movement threshold which was discussed in Section 3.2. Fig. 4 illustrates the key process of the base algorithm. Even though the second plot is closer to the vertical link in the figure, it should be matched to the horizontal link because the third plot is located in horizontal direction and matching the second plot to the vertical link produces an unrealistic detour. This process is based on the same concept as the first one presented by Lou et al. (2009). Note that this algorithm processes not the latest (third) plot but one or more previous plots (the second plot and sub-second plots). In distance-based data, especially, the resulting time lag becomes large, e.g. 72.4 s in 450 m interval data (see Table 3). Although this means real-time processing cannot be achieved, information about a vehicle's movement direction is obtained with a higher confidence level, leading to greater accuracy in map matching.

In Step 4 of the above algorithm, the angle between the link matched third plot and the direction from the second to third plots must be less than 90°. The last sub-third plot is not used for calculating the angle because a vehicle's direction of movement can be miscalculated if an outlier is present within a short distance to the third plot (Quddus et al., 2003).

Fig. 5 shows the map matching accuracy of this base algorithm by solid lines. Because two different data collection methods are considered (distance-based and time-based), the horizontal scale is the average distance between successive plots in each data set. The figures show that both ARP and ARR decrease as the polling frequency decreases, and the base algorithm is unable to correctly match (identify the correct route for) the low frequency data. On the other hand, IARR tends to be lower (higher accuracy) for low frequency data because high frequency data deteriorate map matching in areas of high road density because of GPS error. Additionally, a route that is matched based on the minimum distance tends to be shorter than the actual route. Moreover, data set 2 has better indices than data set 1 because, as shown in Table 2, the ratio of expressway and urban expressway travel is much higher in this set. In other words, since data set 2 has fewer turning movements and a lower road density around the route, it is relatively easier to map match. Therefore, the map matching ease and the usability of the two different polling schemes (time-based scheme and distance-based scheme) cannot be compared.

For almost all polling intervals, ARR is lower than the ARP because the ARR index evaluates the accuracy of the route between consecutive plots. In other words, improving the ARR inevitably improves the ARP. Therefore, map matching accuracy is evaluated in this work using the ARR and IARR indices.

The dotted lines in Fig. 5 show the accuracy of the map matching results by not using third plot (using only two consecutive plots). In these cases, the second plot is matched to the nearest link that is in a consistent direction from the first to second plots and the shortest distance route is searched from the first link to the second plot matched link. The sub-second plots are matched to the nearest links on the route. From these figures, we find that if the third plot is not used, the map matching accuracy is reduced in all cases. These figures show the effectiveness of the usage of the third plot.
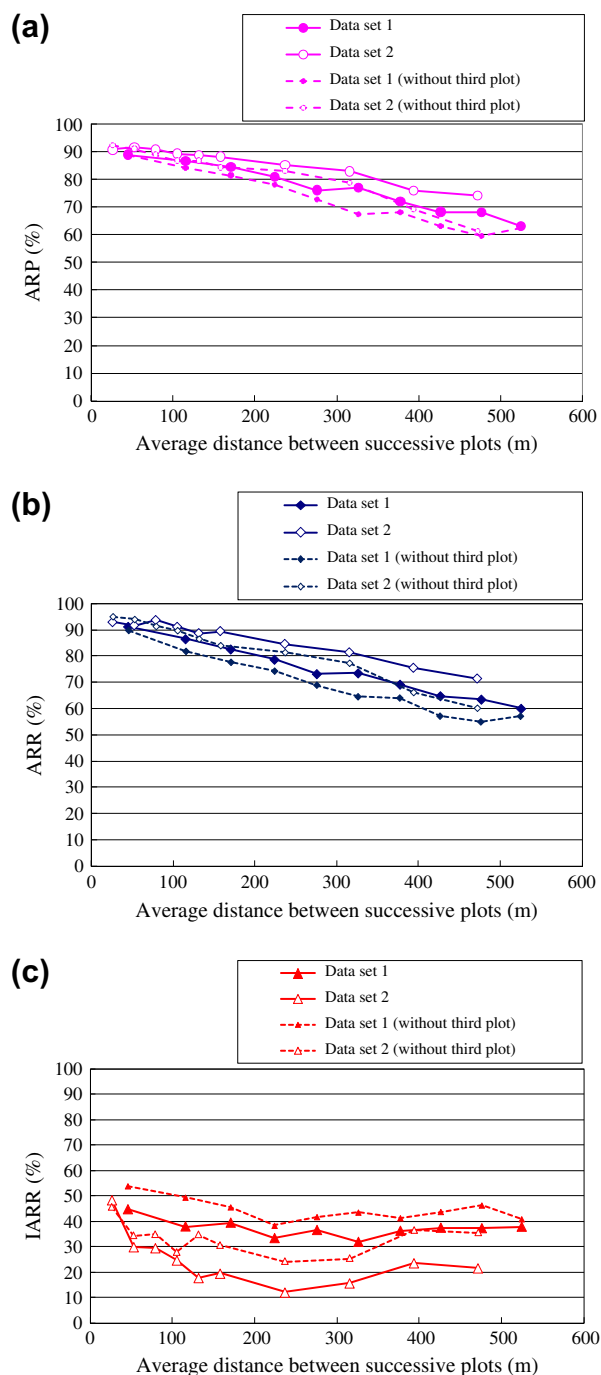
**Fig. 5.** Map matching accuracy of base algorithm: (a) Change in $A_P$ with polling interval, (b) change in $A_L$ with polling interval, (c) change in IARR with polling interval.

## 4.2. Map matching errors of base algorithm

Fig. 6 shows examples of mis-matching that are frequently observed in the results obtained with the base algorithm. Fig. 6a shows a matched route (red[1] line) distant from the second plot. This is the most often observed type of error and is caused by applying a minimum distance route in the matching process. Fig. 6b shows a matched route for an unrealistic

---

[1] For interpretation of color in Fig. 6, the reader is referred to the web version of this article.
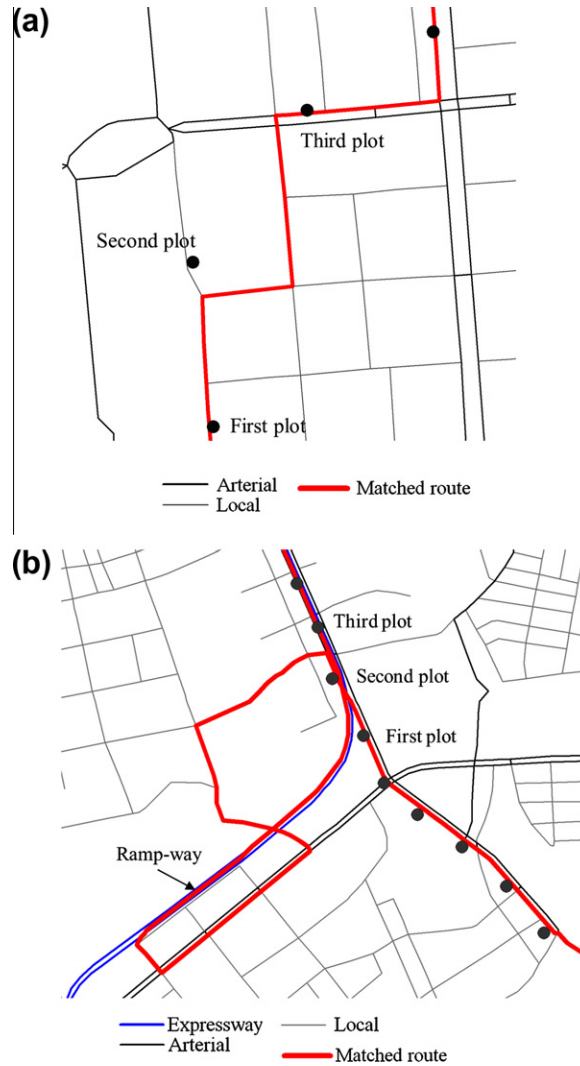
**Fig. 6.** Example of map matching errors: (a) example of mis-matching due to using only the minimum distance route, (b) example of mis-matching around expressways.

detour via an expressway. This error is due to temporary matching of the third plot to the nearest link. Since the third plot is matched to an expressway link (blue line), which is the nearest to the third plot disregarding the current route (on an arterial road), the matched route (red line) has to make a detour so as to enter an expressway through a ramp-way. This type of error is often seen near elevated urban expressways. The following subsections propose methods for reducing these types of errors.

### 4.3. Application of link cost

Map matching errors due to applying the minimum distance route (Fig. 6a) can be reduced by defining an appropriate link cost for the route search in Step 5. In the base algorithm, the link length is set to the link cost. Fig. 7 shows the link cost concept. In the situation shown in the figure, the shortest distance route from the first plot to the third plot is the lower one. However, second plot is near to the upper route. In order to search the route which is near to the second plot, the link cost should be modified using the distances from the second plot to each link ($d_{i,2}$, $d_{j,2}$, $d_{m,2}$, $d_{n,2}$). In this study, the following two alternative link costs are considered.

Link cost 1:

$$cost_i = d_{i,2} \times l_i \tag{3}$$

Link cost 2:

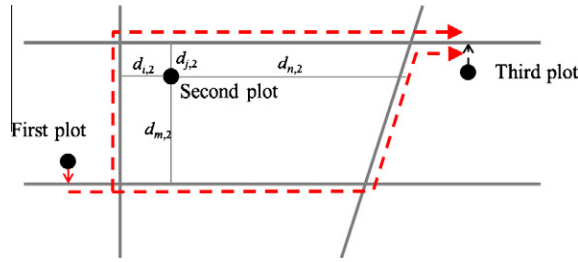$$cost_i = -1 \times ln(Prob_{i,2}) \times l_i \tag{4.a}$$
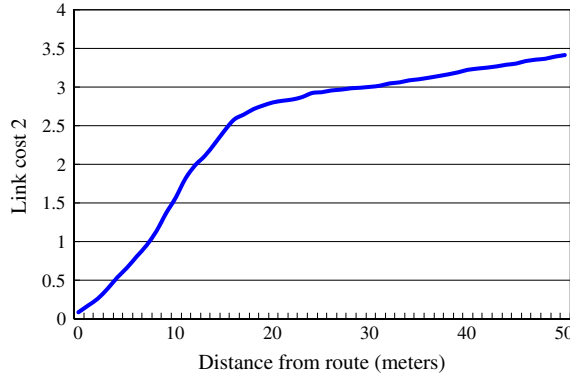
**Fig. 7.** Concept of link cost.



**Fig. 8.** Value of link cost 2 (for link length of 1 m).

$$Prob_{i,2} = 1 - CDF(d_{i,2}) \tag{4.b}$$

where for link $i$, $cost_i$ is the link cost, $l_i$ is the link length, $d_{i,2}$ is the distance to the nearest second or sub-second plot, and $Prob_{i,2}$ is the probability measure that the nearest second or sub-second plot is located on link $i$ and this measure is calculated by Eq. (4.b). And *CDF* (cumulative distribution function) can be determined from Fig. 2. Past studies have shown that link cost 1 will have an effect (Miwa et al., 2004), and link cost 1 is intuitive and easy to use. On the other hand, link cost 2 shows the likelihood that the nearest plot is on the link and expresses the statistical influence of GPS errors. Fig. 8 shows the value of link cost 2 in the case that the link length is 1 m. From this figure, we find that link cost 1 changes linearly in accordance with the distance from route, whereas link cost 2 increases relatively rapidly until 15 m and moderately after that.

### 4.4. Selection of multi-candidate routes

Map matching errors due to inappropriate temporary matching of the third plot can be reduced by considering more than one link for temporary matching in Step 4. In this study, we matched the third plot with up to three links that are in a consistent direction from the second to third plots. Minimum cost routes from the first link to each third link are map matching route candidates. Increasing the number of temporary matching links increases the possibility that the correct route is included in the set of candidate routes. However, processing time also will increase if too many links are considered for temporary matching. Three links are adopted for temporary matching because, if the direction from the second to third plots is considered, a maximum of three possible links — the nearest surface road, expressway, and expressway ramp — should be considered as candidates.

However, the set of candidate routes may include an unrealistic detour like the case shown in Fig. 6b. Therefore, the minimum required travel time for each candidate route $k$ from the first to third plots is defined as $T_{min,k}$. If the actual travel time, $T$, is shorter than $T_{min,k}$, the route is excluded from the set of candidates. This procedure is similar to the second concept presented by Lou et al. (2009). The minimum possible travel time $T_{min,k}$ can then be expressed as follows.

$$T_{min,k} = \frac{L_{e,k}}{v_e} + \frac{L_{g,k}}{v_g} \tag{5}$$

where for a certain candidate route $k$, $L_{e,k}$ is the length of the expressway/urban expressway section, $L_{g,k}$ is the length of the surface road section. $v_e$ is the maximum possible speed of travel on the expressway/urban expressway, and $v_g$ is the maximum possible speed of travel on the surface road. In this study, $v_e$ and $v_g$ are set to 120 km/h and 80 km/h, respectively, in order to account for GPS errors (in Japan, the regulation speeds are 100 km/h and 40–60 km/h, respectively).

### 4.5. Evaluation of candidate routes

#### 4.5.1. Evaluation based on positional relationship
The following two evaluation formulas using a positional relationship are considered for identifying a map matching route.

Evaluation formula 1:

$$\min E_k = \sum_{j2} d_{j_2,k} + d_{2,k} + \sum_{j3} d_{j_3,k} + d_{3,k} \tag{6}$$

Evaluation formula 2:

$$\max E_k = \sum_{j_2} ln(Prob_{j_2,k}) + ln(Prob_{2,k}) + \sum_{j_3} ln(Prob_{j_3,k}) + ln(Prob_{3,k}) \tag{7}$$

where for a certain candidate route $k$, $E_k$ is the evaluation value, and $d_{2,k}$ ($d_{j_2,k}$) is the distance from the second plot ($j$th sub-second plot) to the nearest link on the route. $Prob_{2,k}$ ($Prob_{j_2,k}$) is the probability that the second plot ($j$th sub-second plot) is located on the route and is calculated in the same manner as Eq. (4.b).

#### 4.5.2. Evaluation based on driver's route choice concept
If the distance between consecutive plots is long, several routes might be proximate to the first, second, and third plots. This situation often arises in cases where low frequency data are used and/or in areas with a highly dense road network. In these cases, determining the route from the positional relationship of plots and candidate routes is clearly difficult. In this subsection, a function that expresses the likelihood of a probe vehicle's taking a candidate route and that can be used to determine the traveled route is proposed.

In this study, we assume that the likelihood of a target probe vehicle taking a particular candidate route can be expressed as follows:

$$g_k = \sum_i \left\{ \left( \beta_{i,0} + \sum_j \beta_{i,j} \bar{x}_j \right) x_{k,i} \right\} + \varepsilon_k \tag{8}$$

where $g_k$ is the likelihood value for candidate route $k$, $x_{k,i}$ is the $i$th attribute of candidate route $k$, and $\bar{x}_j$ is the $j$th attribute of the route traveled by the target probe vehicle prior to the first link. In other words, in this function, the attributes of the route traveled prior to the first link are assumed to influence the choice of the subsequent route. $\beta_{i,j}$ is the unknown parameter that expresses the influence of $x_{k,i}$ and $\bar{x}_j$ on the likelihood. Finally, $\varepsilon_k$ is a random term with zero mean.

If random term $\varepsilon_k$ follows a Gumbel distribution, then the probability that a target probe vehicle used candidate route $k$ can be expressed as the logit model represented by Eq. (9). The unknown parameters in Eq. (8) can be estimated using the maximum likelihood estimation method (Ben-Akiva and Lerman, 1985).

$$P_k = \frac{exp(\theta V_k)}{\sum_k^{\prime} \exp(\theta V'_k)} \tag{9}$$

where $P_k$ is the probability that the probe vehicle chooses route $k$, $\theta$ is the scale parameter for the Gumbel distribution, $V_k$ is the systematic term for the utility of choosing route $k$, which is expressed by the first term on the right side of Eq. (8). Furthermore, the variance of the random term may change according to the route length (Gliebe et al., 1999). This heteroscedasticity can be relaxed by introducing the following scale parameter structured according to the length of the route from the first link to the third link, $D_{1,3}$ (Morikawa and Miwa, 2006).

$$\theta = \frac{\pi}{\sqrt{6}D_{1,3}} \tag{10}$$

In this study, the function represented by the first term on the right side of Eq. (8) is called the "route determination model." This model is applied to the evaluation of candidate routes which are generated by applying link cost and route selection. The candidate route with the maximum value calculated by the route determination model is adopted.

## 5. Validation of map matching accuracy

### 5.1. Accuracy of map matching using positional evaluation formulas

The preceding section proposed two link cost formulas, minimum possible travel time, and three evaluation formulas for possible use in map matching. Of the three evaluation formulas, those based on a positional relationship are evaluated in this subsection. The map matching results for all combinations of two link costs and two positional evaluation formulas show that link cost 1 produces slightly more accurate results than link cost 2, with little difference in the results given by the
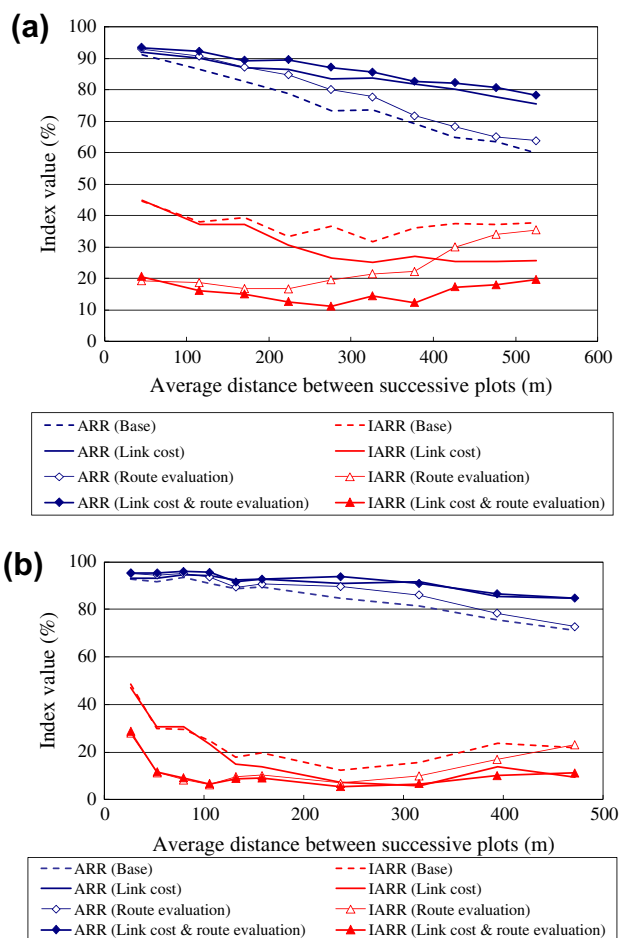
**Fig. 9.** Map matching accuracy after applying link cost and route selection: (a) data set 1, (b) data set 2.

two evaluation formulas. The results of applying link cost 1, minimum required travel time, and evaluation formula 2 are shown in Fig. 9.

These figures show that the improved algorithm has better map matching accuracy than the base algorithm, with ARR increasing from 60% to 80% for data set 1 and from 70% to 85% for data set 2, especially for low frequency data. Moreover, the IARR of the improved algorithm decreases from 40% to 20% for data set 1 and from 20% to 10% for data set 2, also for low frequency data. The ARR for both data sets is further increased (improved) by applying the link cost, which considers the relationship between the candidate route and the second and sub-second plots. On the other hand, the IARR can be decreased (improved) further by applying the link cost to the low frequency data and multi-candidate routes to the high frequency data. These results indicate that applying only the link cost does not prevent map matching errors.

### 5.2. Accuracy of map matching using driver's route choice behavior concept

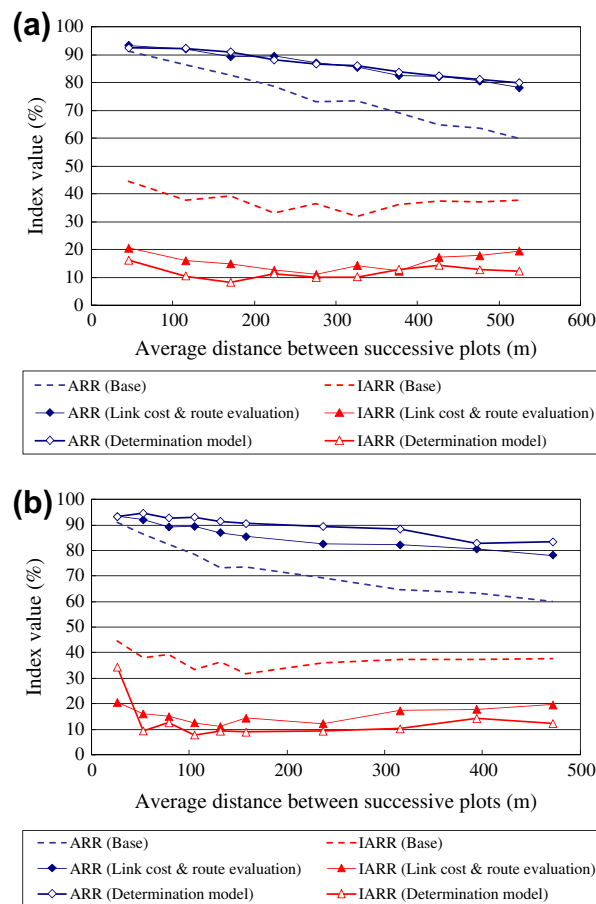#### 5.2.1. Parameter estimation

Table 5 shows the estimated results for the unknown parameters. In this case, since the number of samples is sufficient for statistical evaluation, a parameter has a statistical significance with 95% confidence level if the $t$-value is larger than 1.96. The parameters shown in the table are estimated from all data sets. To carry out this estimation, the candidate route set comprises candidate routes from the first link to third links as well as the correct route. The travel time for the unit distance is the actual travel time from the first plot to the third plot divided by each candidate route distance from the first link to the third link. This value is a measure of how reasonable the traffic information generated by the map matching process is.

Table 5 also shows that the parameters for $\bar{x}$ are estimated with statistical significance. This means that the attributes of the route before and after an arbitrary point (the first plot) are related. Specifically, the parameter for the arterial road dummy is a positive value with statistical significance and is part of the parameter for the length of the surface road. This means that if the route traveled is an arterial road, the probe vehicle tends to continue traveling on arterial roads. Moreover, the positive parameter for the local street dummy and the negative parameter for the expressway dummy, which both have

**Table 5**
Results of estimates using the route determination model.

| Explanatory variables | | Estimates (t-value) | |
|---|---|---|---|
| x | x̄ | | |
| Length of expressway (m) | – | −0.500 | (−6.9) |
| Length of surface arterial road (m) | Constant | −0.157 | (−2.0) |
| | Arterial road dummy | 0.481 | (4.4) |
| Number of turns (local street → local street) | – | −59.2 | (−5.2) |
| Number of turns (local street → arterial road) | – | −32.2 | (−2.9) |
| Number of road type changes (arterial road ↔ local street) | – | −64.9 | (−4.7) |
| Travel Time (seconds/m) | Constant | −4.41 | (−6.9) |
| | Night dummy | 1.73 | (2.9) |
| | Local street dummy | 1.90 | (2.8) |
| | Expressway dummy | −10.1 | (−2.0) |
| Number of samples | | 427 | |
| Adjusted Rho-squared | | 0.608 | |
| Hit ratio | | 0.867 | |



**Fig. 10.** Map matching accuracy after applying the route determination model: (a) data set 1, (b) data set 2.

statistical significance and are parts of the parameter for travel time for unit distance, mean that if the probe vehicle is on a local street, it tends to travel at a low speed thereafter, while if it is on an expressway, it travels at a high speed thereafter.

The parameters shown in this table are probably generic for taxi probe vehicles in the study area, so they can be applied to a taxi-based map matching system in Nagoya. However, there are differences in route choice behavior among various vehicle types, such as tuck and private vehicles. If this method is to be applied to a probe vehicle system using other types of vehicle, the parameters in the route identification function should be re-estimated.

### 5.2.2. Map matching accuracy

In this subsection, the availability of the route determination model is verified. This determination model is applied for selecting route from candidate routes which are made using link cost 1 and minimum required travel time.

As mentioned above, parameters shown in Table 5 are generic for taxis. However, since the parameters were estimated using all of the data sets, these parameters should not be applied again to the same data sets in order to validate the methodology. Therefore, two new data sets are constructed; one for vehicles Nos. 1 and 4 (data in the first row and last row in Table 1), and the other for vehicles Nos. 2 and 3 (data in the second row and third row in Table 1). The parameters estimated from each new data set are applied to the other in a procedure known as "cross-validation." The values for the two sets of estimated parameters are similar to those for all of the data sets (Table 5).

Fig. 10 shows the map matching accuracy. These figures show that applying the route choice concept can further improve map matching accuracy. ARRs have not improved in data set 1, but have improved in data set 2. The improvements in data set 2 are due to correctly choosing the expressway ramp. Additionally, although there is some variation in the results, IARRs have improved in both data sets. This is because the route determination model is able to exclude routes that repeat unrealistic movements, such as traveling on local streets and/or turning at minor intersections. These results mean that understanding a driver's route choice behavior may enable us to develop a more accurate map matching algorithm.

## 6. Conclusions and future research

In this study, the authors have developed an on-line map matching algorithm that can be applied to low frequency and little information probe vehicle data, such as that with a low polling frequency and/or containing little information about the vehicle state. After investigating cases in which the base algorithm failed to achieve map matching, it was shown that there are two typical types of mis-matchings and map matching accuracy can be improved by introducing several additional techniques and parameters. For a matched route distant from the plot, link cost that multiplies link length by the distance from the link to the nearest plot produces slightly more accurate results than link cost that considers probabilistic nagure of GPS error. For inappropriate temporary matching, introducing multi-candidate routes and minimum required travel time for candidate routes are effective. Moreover, introducing a driver's route choice concept and considering the relationship between the route already traveled and the route that will be traveled thereafter has the possibility for further improvement of map matching accuracy.

Additionally, it was shown that if the third plot is not used, the map matching accuracy is reduced. Such reduction of the accuracy happens in all methods proposed in this study. These results mean that the existing on-line map matching methods designed for low frequency probe data which utilize only two consecutive plots could be improved by using third plot.

Although the map matching program coded in Fortran in this study is not necessarily optimized for processing speed, it was able to process all of the data in Table 1 (about 30 h of data) in as little as 7 s (low frequency data) and 40 s (high frequency data) using a conventional laptop PC (cpu: Intel core 2 2 GB; memory: 2 GB). Applying a more effective route search algorithm, such as the A* algorithm of Hart et al. (1968), would realize a speedier process. Therefore, it is thought that the developed algorithm can be applied to large scale probe vehicle systems. Finally, this study does not analyze the differences in the availability of polling schemes. In the future, a map matching algorithm with higher accuracy should be developed based on further research into the map matching accuracy of the data from various polling schemes.

### References

Ben-Akiva, M., Lerman, S., 1985. Discrete Choice Analysis. MIT Press, Cambridge, MA.
Brakatsoulas, S., Pfoser, D., Salas, R., Wenk, C., 2005. On map-matching vehicle tracking data. In: Proceedings of the 31st International Conference on Very Large Databases, pp. 853–864.
Gliebe, J.P., Koppelman, F.S., Ziliaskopoulos, A., 1999. Route choice using a paired combinatorial logit model. In: 78th Annual Meeting of the Transportation Research Board, Washington, DC, USA.
Greenfeld, J.S., 2002. Matching GPS observations to locations on a digital map. In: 81st Annual Meeting of the Transportation Research Board, Washington, DC, USA.
Hart, P.E., Nilsson, N.J., Raphael, B., 1968. A formal basis for the heuristic determination of minimum cost paths. IEEE Transactions on Systems Science and Cybernetics 4, 100–107.
Hellinga, B.R., Fu, L., 2002. Reducing bias in probe based arterial link travel time estimates. Transportation Research Part C 10, 257–273.
Hellinga, B., Fu, L., Takada, H., 2003. Obtaining traveller information via mobile phone location referencing – challenges and opportunities. In: Proceedings of the Annual Conference of the Transportation Association of Canada, pp. 21–24.
Hellinga, B., Izadpanah, P., Takada, H., Fu, L., 2008. Decomposing travel times measured by probe-based traffic monitoring systems to individual road segments. Transportation Research Part C 16, 768–782.
Kim, S., Kim, J.H., 2001. Adaptive fuzzy-network-based C-measure map-matching algorithm for car navigation system. IEEE Transactions on Industrial Electronics 48, 432–441.

Kojima, E., Hato E., 2004. Online matching algorithm using probe person data. In: Proceedings of the 29th Annual Conference of the Infrastructure Planning, Kobe, Japan (in Japanese).

Liu, K., Yamamoto, T., Morikawa, T., 2006. An analysis of the cost efficiency of probe vehicle data at different transmission frequencies. International Journal of ITS Research 4, 21–28.

Liu, K., Yamamoto, T., Li, Q., Morikawa, T., 2007a. Cost-effectiveness in probe vehicle systems. Journal of Eastern Asia Society for Transportation Studies 7, 116–127.

Liu, K., Yamamoto, T., Morikawa, T., 2007b. Comparison of time/space polling schemes for a probe vehicle system. In: Proceedings of the 14th World Congress on Intelligent Transport Systems, Beijing, China.

Lou, Y., Zhang, C., Zheng, Y., Xie, X., Wang, W., Huang, Y., 2009. Map-matching for low-sampling-rate GPS trajectories. In: Proceedings of the 17th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems, pp. 352–361.

Makimura, K., Kikuchi, H., Tada, S., Nakajima, Y., Ishida, H., Hyodo, T., 2002. Performance indicator measurement using car navigation systems. In: 81st Annual Meeting of the Transportation Research Board, Washington, DC, USA.

Miwa, T., Sakai, T., Morikawa, T., 2004. Route identification and travel time prediction using probe-car data. International Journal of ITS Research 2 (1), 21–28.

Morikawa, T., Miwa, T., 2006. Preliminary analysis on dynamic route choice behavior using probe-vehicle data. Journal of Advanced Transportation 40 (2), 141–163.

Murakami, E., Wagner, D.P., 1999. Can using global positioning system (GPS) improve trip reporting? Transportation Research Part C 7, 149–165.

Quddus, M.A., Ochieng, W.Y., Zhao, L., Noland, R.B., 2003. A general map matching algorithm for transport telematics applications. GPS Solutions Journal 7 (3), 157–167.

Quiroga, C.A., Bullock, D., 1998. Travel time studies with global positioning and geographic information systems: an integrated methodology. Transportation Research Part C 6, 101–127.

Sermons, M.W., Koppelman, F.S., 1996. Use of vehicle positioning data for arterial incident detection. Transportation Research Part C 4 (2), 87–96.

Tradisauskas, N., Juhl, J., Lahrmann, H., Jensen, C.S., 2009. Map matching for intelligent speed adaptation. IET Intelligent Transport Systems 3 (1), 57–66.

Wenk, C., Salas, R., Pfoser, D., 2006. Addressing the need for map-matching speed: localizing global curve-matching algorithm. In: 18th International Conference on Scientific and Statistical Database Management, pp. 379–388.

White, C.E., Bernstein, D., Kornhauser, A.L., 2000. Some map matching algorithms for personal navigation assistants. Transportation Research Part C 8, 91–108.

Yamamoto, T., 2005. P-DRGS and Eco-point TDM project in the second stage at the Aichi Expo. Journal of International Association of Traffic and Safety Sciences 29 (1), 110–113.

Yang, J., Kang, S., Chon, K., 2005. The map matching algorithm of GPS data with relatively long polling time intervals. Journal of Eastern Asia Society for Transportation Studies 6, 2561–2573.

Yin, H., Wolfson, O., 2005. A weight-based map matching method in moving objects databases. In: Proceedings of the 16th International Conference on Scientific and Statistical Database Management, pp. 437–438.