# A Method for Parsing Route Descriptions using Sidewalk Network Databases

Kouzou Noaki and Masatoshi Arikawa

*Center for Spatial Information Science, at the University of Tokyo* 4-6-1 Komaba, Meguro-ku, Tokyo, 153-8904 Japan
E-mail:    {noaki, arikawa}@csis.u-tokyo.ac.jp

**Abstract** People use descriptions referring to places through their daily life. These descriptions correspond to the region of the real world. We call such descriptions geo-referenced descriptions. We have studied a method of converting such descriptions like addresses and place names into their corresponding coordinates, that is, tuples of longitude and latitude. The process of converting descriptions into coordinates is called *geocoding*. In this paper, we focus on natural route descriptions as a new type of target to geocode. We first explain a core schema of sidewalk network databases on the basis of a characteristic of natural route descriptions, and then propose *Formal Route Statement* (*FRS*) to represent and process natural route descriptions by means of a computer. Also, we present our prototype system to geocode natural route descriptions using sidewalk network databases based on our proposed framework.

**Keyword**   Web and Internet,    Spatial DB,    Data Visualization,  GIS

## 歩道ネットワークを用いた経路記述を対象とした構文解析手法

野秋　浩三　　　　有川　正俊

東京大学空間情報科学研究センター　〒153-8904　東京都目黒区駒場四丁目 6 番 1 号
E-mail:    {noaki, arikawa}@csis.u-tokyo.ac.jp

**あらまし**　我々は日常の多くのシーンにおいて現実空間の位置を参照する記述を使用している. 多様なデジタルデータを現実空間の位置で検索・管理することは情報の活用可能性を広げ, 利用を高度化させる. 住所や地名文字列を解釈し, 対応する位置の座標値に変換する手法を総称してジオコーディングという. 本研究では, 計算機上で自然言語による場所記述をジオコーディングするための基礎研究として, 都市空間における経路記述を対象としたジオコーディング手法を提案する. 具体的には, 経路記述を解析した結果を地図上に線分表示する. 歩道ネットワーク・データベースを利用し, この処理を実現するために経路記述の幾何記述である Formal Route Statement とジオコーディングアルゴリズムを開発し, この枠組みに基づいてプロトタイプシステムの実装を行った.

**キーワード**　Web とインターネット, 空間 DB, データの可視化,  GIS

## 1. Introduction

Retrieving information on the Internet has become indispensable in our daily life. When we want to have a tasty meal with our friends, we get information of restaurants from the Internet using search engines for some keywords. However, it is not easy to get exact information because there are billions of documents and the number of them is getting to increase. It may cause us inconvenient when we are finding information using computers, especially local information like restaurants.

We assume that the reason of the inconvenience is the difference of the framework of managing information between our brains and the current computer systems. One of people's methods of memorizing information is being associated with locations. Examples of this are "Mr. Suzuki who lives in Tokyo", "a good Chinese noodle shop under the elevated rail" and "an important file in the second drawer of my desk". On the other words, locations remind people of their memory. We define these descriptions which correspond to the region in the real world as *geo-referenced descriptions*. Our motivation for the research is considering the framework to dealing with information using geo-referenced descriptions as important one of keys referring to the information.

We have studied a method of converting geo-referenced de-

scriptions like addresses and place names into their corresponding geographic coordinates [1]. The process of converting descriptions into coordinates is called geocoding. In this paper, we focus on natural route descriptions as a new type of target to geocode.

Section 2 defines geocoding problem. Section 3 explains sidewalk network databases and a prototype system for managing the database. In section 4, we propose *Formal Route Statement* (*FRS*) to represent and process route descriptions in natural language by means of a computer. In section 5, we mention a method for geocoding natural route descriptions and a prototype system. Finally, we state conclusion and future work in section 6.

## 2. Definition of Geocoding Problem

Geocoding problem is a problem for inquiring the region in the real world of the geo-referenced description. We define geocoding problem as follows:

*Input data*:    a text string $g$ of geo-referenced description

*Output data*:    a region $r$ corresponding to $g$ in the real world

Geocoding process $F$ is defined as a mapping function which meets the following equation.

$$r = F(g)$$

This paper provides a solving method of geocoding problem using sidewalk network databases. We go into details of sidewalk network databases in the next section.

## 3. Sidewalk Network Database

Sidewalk network databases store underground walks, footbridges and cross walks for pedestrians. Sidewalk network databases are simply structured as nodes and links. They are provided as commercial products by Shobunsha Publications Inc. [2]. The commercial sidewalk network databases presently cover major cities in Japan. We redesign the schema of sidewalk network databases in order to increase flexibility in advanced applications. In this section, we explain the schema of sidewalk network databases and prototype system for managing it.

### 3.1. Schema of Sidewalk Network Database

A sidewalk network database is a directed graph representing sidewalks in urban environment. It consists of

1.    A set $N$ of nodes and

2.    A binary relation $L$ on $N$. We call $L$ the set of links of the directed graph. Links are thus pairs of nodes and also vectors between the pairs of nodes.

We define spatial objects as data units corresponding to entities in the world (e.g., footbridges, crosswalks and intersections). An example of sidewalk network databases is shown in Fig. 1. Each node is represented by a circle. In Fig.1, a set $N$ is the set of node data $n$.

$$N = \{n_1, n_2, n_3, \ldots\}$$

A node data $n$ has the following structure.

$$n = (id, spatial\_anchor, class, m, img, in\_link, out\_link)$$

The elements of  $n$  are the followings:

| | |
|---|---|
| *id* | the identifier of the node |
| *spatial_anchor* | a name of a instance of spatial object |
| *class* | a class of spatial object |
| $m = (x_m, y_m)$ | geographic coordinates |
| *img* | picture data |
| *in_link* | the identifiers of the incoming links to $n$ |
| *out_link* | the identifiers of the outgoing links from $n$ |

A *spatial_anchor* is a text data referring to the instance of spatial object. Nodes composing the specific instance of spatial object have *spatial_anchor* value only if the place has a name. Spatial object in the database can be classified into four main classes according to their role in the city or geographic features. Each class has its own characteristics and restrictions [3]. The classes categorized spatial objects of the same type such as "station exit", "intersection", "street" and "slope street". For example, the *spatial_anchor* and the *class* value of the node $n_2$ are "$\alpha$ station exit" and "station exit". Composing the $\beta$ intersection on the map, the nodes $n_3$, $n_4$, $n_5$ and $n_6$ have

" $\beta$ intersection" as *spatial_anchor* and "intersection" as *class*.

Next, in Fig.1, a set $L$ is the set of link data $l$.

$$L = \{l_1, l_2, l_3, \ldots\}$$

A link data $l$ has the following structure;

$$l = (start\_node, end\_node)$$

The elements of $l$ are the following;

    $start\_node$      the *id* of the node which $l$ starts from

    $end\_node$      the *id* of the node which $l$ arrives at

Each link in $L$ is represented by an arrow from *start_node* to *end_node*. For example, $n_3 \rightarrow n_6$ is the link $l_7$ of Fig. 1 and $n_6 \rightarrow n_3$ is the link $l_{110}$ of Fig. 1. It is a notable fact that there is not the link $n_2 \rightarrow n_1$ because sidewalk network databases do not include sidewalks leading inside buildings.
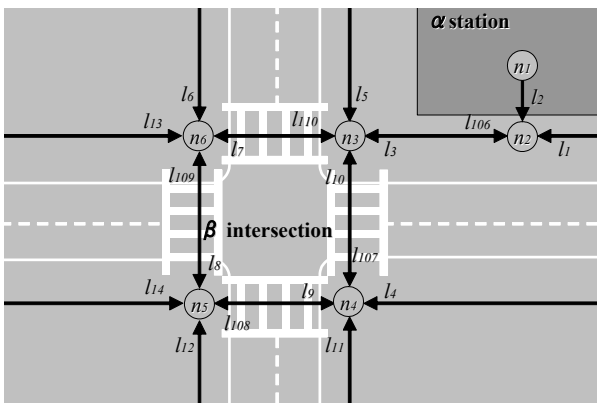


**Fig.1.** Example of a sidewalk network database

## 3.2. Prototype of Database Management System

We have developed a naive management system for sidewalk network database. We explain each component in the user interface (Fig. 2) as follows:

(A) Menu button

Users can change the operation mode by the menu buttons. Main functions are (1) loading and saving network data which is XML format and (2) adding, erasing and moving both nodes and links. Furthermore, we select functions of referring to node information (e.g., a name of a spatial object) and filling it using the entry form.

(B) Map area

Sidewalk network databases and a map image in the selected area are overlapped and visualized. Figure 3 shows a dialog box for an entry form of node information, which allows users to add a name and a class for a spatial object. This window appears when users push the entry button in menu buttons and click a node on map area.

(C) Display of node information

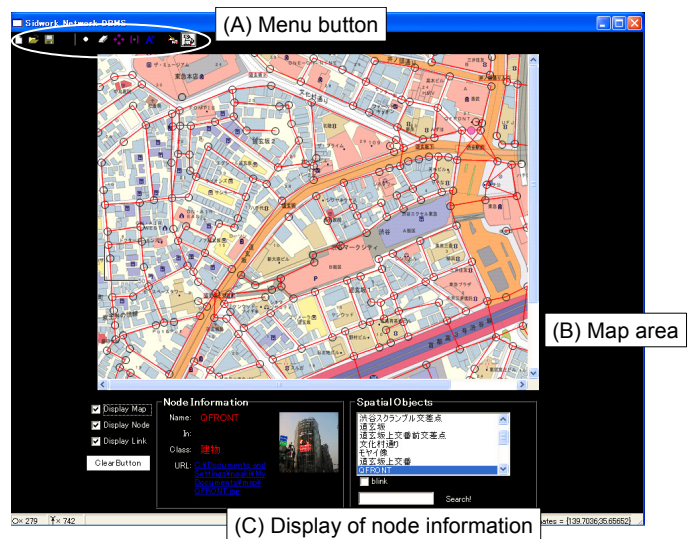This area allows users to see the name, class and picture image of a spatial entity on the map area.



**Fig.2.** Graphical user interface of the prototype system for managing sidewalk network databases
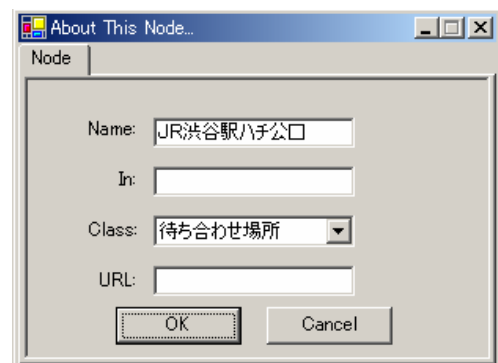


**Fig. 3.** Dialog box of entry form for node information

## 4. Formal Route Statement

On the assumption that all of route descriptions in natural language can be expressed with nodes and links, we propose Formal Route Statement (FRS) to represent and process natural route descriptions. FRS also works as a query language for the sidewalk network databases. Using FRS, a route description is represented as sub-graph of directed graph of sidewalk networks (Fig. 4).

Figure 5 shows the grammar of FRS. Generalization tables are indispensable to converting various casual descriptions into regular ones, one of which is FRS (Table 1). A use case of the generalization tables is to make an instance of the spatial relationship as a value of the attribute "*link.connect*" in Fig. 3. The attribute "*link.connect*" plays an important role to find spatial object when a name referring to next place is omitted.
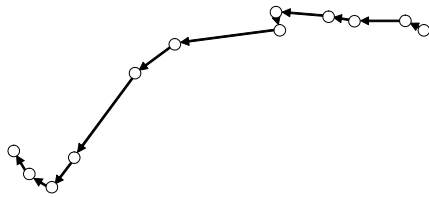


**Fig.4.** Sub-graph *G* which is the result of *FRS*(*route_desc*), where *route_desc* is a text string of a natural route description

| | |
|---|---|
| *FRS* ::= | *node_desc*(0)(:*link_desc*(*i*):*node_desc* (*i*+1))* [*i*={0,...,*n*}]; |
| *node_desc*(*i*) ::= | *node*(*i*).*node_attribute_list* |
| *node_attribute_list* ::= | *none* \| *node_attribute_value* (&*node_attribute_value*)* |
| *node_attibute_value* ::= | *node_attribute* = *value* |
| *node_attribute* ::= | *id* \| *name* \| *coordinate* \| *in* \| *class* \| *status* |
| *value* ::= | *numerical_value* \| *stirng_value* \| *url* \| *status_values* \| *connect_values* |
| *status_values* ::= | *start* \| *end* \| *via* |
| *connect_values*::= | *straight* \| *right* \| *left* |
| *link_desc*(*i*) ::= | *link*(*i*).*link_attribute_list* |
| *link_attribute_list* ::= | *none* \| *link_attribute_value* (&*link_attribute_value*)* |
| *link_attibute_value* ::= | *link_attribute* = *value* |
| *link_attribute* ::= | *id* \| *start_node*(*id*) \| *end_node*(*id*) \| *direction* \| *connect* \| *distance* |

**Fig 5.** Grammar of Formal Route Statement (*FRS*)

**Table 1.** Example of a generalization table for descriptions of spatial relationship and values of the attribute "*link.connect*"

| Specialized descriptions for spatial relationships | Generalized descriptions for spatial relationship (values of *link.connect*) |
|---|---|
| go forward, go ahead, advance | straight |
| turn to the right, on the right, on one's right | right |
| turn to the left, on the left, on one's left | left |

## 5. Geocoding for Natural Route Description

An *address description*, that is a kind of geo-referenced description, includes both some *numbers* for representing the locations of buildings, and some *names* for representing administrative areas where someone lives or works. An address description corresponds to a closed region on a large-scale map or a point on a small-scale map. On the other hand, a *route description* includes both some *nouns* referring to place names and *prepositions* for representing spatial relationships between places. A route description usually corresponds to a polyline on a map. The process is not only converting place names in a route description into geographic coordinates, but also verifying spatial validity between places.

The geocoding for natural route descriptions consists of the following two processes. We go into details of them in the following sections (5.1) and (5.2).

- Converting Natural Route Description into Formal Route Statement (5.1)
- Validating Formal Route Statement (5.2)

Formal Route Statement (FRS) is a geometric model for representing natural route descriptions. Computers can deal with natural route descriptions indirectly through FRS as an intermediate description.

## 5.1. Converting Natural Route Description into Formal Route Statement

The procedure to convert natural route descriptions into FRS can be classified into two main processes. In this stage, FRS need not have geographic coordinates.

(a) Separating spatial anchors and spatial relationship descriptions in a natural route description

In the case that the input text is "ハチ公口を出て道玄坂を上り交差点を右へ", the text is separated as follows:

$node\_desc(0)$ = "ハチ公口"

$link\_desc(0)$ = "を出て"

$node\_desc(1)$ = "道玄坂"

$link\_desc(1)$ = "を上り"

$node\_desc(2)$ = "交差点"

$link\_desc(2)$ = "を右へ"

(b) Generalizing spatial anchors and spatial relationship descriptions

In the above example, $node\_desc(0)$, $node\_desc(2)$ and $link\_desc(0)$ are inappropriate descriptions for matching with sidewalk network databases. They are generalized as follow:

$node\_desc(0)$ = "JR 渋谷駅ハチ公口"

$node\_desc(2)$ = "道玄坂上交番前交差点"

$link\_desc(2)$ = "を右へ曲がる"

In the above process of (a), there are already significant achievements in the fields of natural language processing [4]. In the above (b), we plan to realize the functions to complement the incomplete name of a spatial anchor and correct inappropriate natural route descriptions using sophisticated geographic thesaurus and generalization rules for names. However (b) is the future work and we do not focus on it in the current stage of our research.

## 5.2. Validating Formal Route Statement

This stage verifies the spatial validity between places which are referred by $node\_desc$s of Formal Route Statement. The procedure to validate Formal Route Statement can be classified into two subprocesses for $node\_desc$ and $link\_desc$.

A $node\_desc(k)$ is needed to match with the set of valid spatial anchors. The set of valid spatial anchors is the result from the following function.

$valid\_sa(passedlink, link\_desc(k\text{-}1))$      $k$=1…$n$

In the function $valid\_sa$, $passedlink$ is the last link which is matched at the time. For example, in Fig. 6, the plane is divided in quarters using the angle of $passedlink$. Depending on the meaning of $link\_desc(k\text{-}1)$, the target of matching with $node\_desc(k)$ is changed. Figure 6 shows that three target areas for matching a set of nodes with $node\_desc(k)$ as follows:

(1)  within the plane of Ⅱ when $link\_desc(k\text{-}1)$ is "right".

(2)  within the plane of Ⅲ when $link\_desc(k\text{-}1)$ is "left".

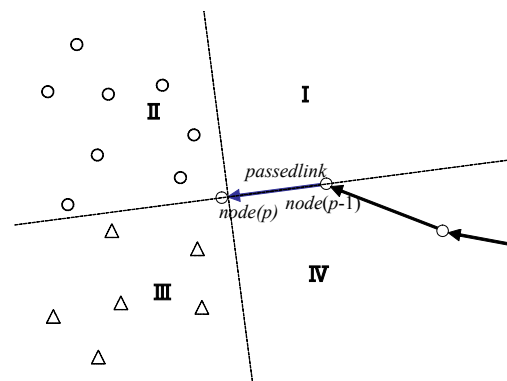(3)  within the plane of Ⅱ+Ⅲ when $link\_desc(k\text{-}1)$ is "straight".



**Fig.6.** Division the plane for matching with valid spatial anchors

The $node\_desc(0)$, which is a head of $node\_desc$s, is matched with the set of all spatial anchors because there is no $passedlink$.

Furthermore, there are the following two types of matching with $node\_desc$. Both matchings are executed simultaneously.

**Spatial anchor point matching**: The results of the queries are a set of nodes corresponding to the spatial anchors.

**Path matching**: This step decides the sequence order of places in the route using the rule which minimizes each of the moving costs between neighbor places.

A *link_desc*(*m*) is needed to match with the set of valid spatial relationship descriptions. The set of valid spatial relationship descriptions is the result from the following function.

*valid_sr*(*passedlink*, *node*(*p*).*out_link*)

Depending on an angle formed by *passedlink* and *node*(*p*).*out_link*, the target of matching with *link_desc*(*m*) is changed. Figure 7 shows that three target areas for matching a set of links with *link_desc*(*m*) is as follows:

(1)  "right" is the target when one or more link data of *node*(*p*).*out_link* are within the region Ⅰ.

(2)  "straight" is the target when one or more link data of *node*(*p*).*out_link* are within the region Ⅱ.

(3)  "left" is the target when one or more link data of *node*(*p*).*out_link* are within the region Ⅲ.
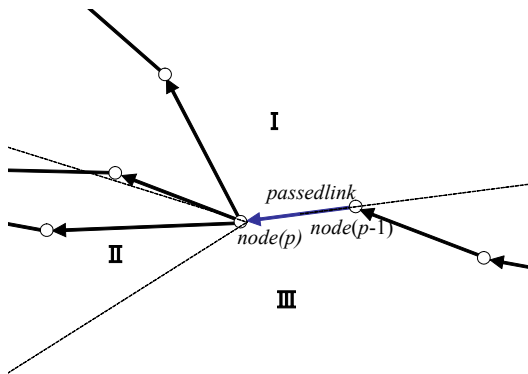


**Fig. 7.** Division the plane for matching with *link.connect*

## 5.3. Prototype System

Our proposed framework has been verified through developing a prototype system. Figure 8 shows the user interface of prototype system.

### 5.3.1. Overview of Prototype System

We have developed a prototype system which processes a route description in Japanese and then visualize it as a polyline on the map using sidewalk network databases. Each component in the user interface (Fig. 8) is as follows:

(A) Input text form

Users can input a natural route description using this text form.

(B) Output text form

A result of separating input text and validating separated elements is displayed in this form.

(C) Output map area

A route is visualized on the sidewalk network as a result of geocoding a natural route description.
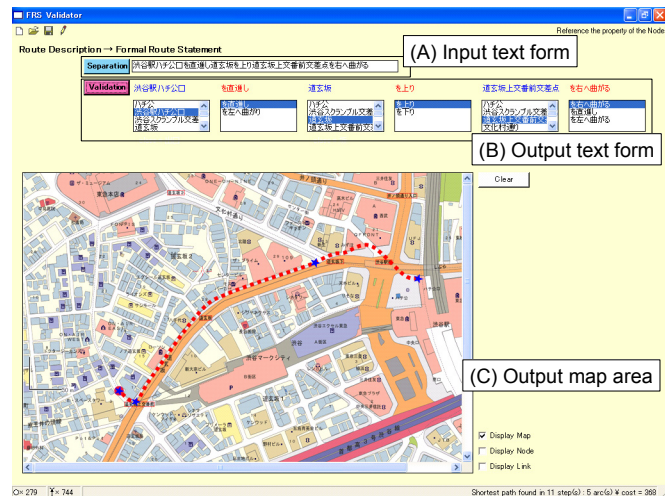


**Fig. 8.** Graphical user interface of the prototype system

### 5.3.2. Experimental Demonstration

Using the example of input text in section 5.1, figure 9 shows behavior of processing. The captions of the figures explain the details of the behavior.
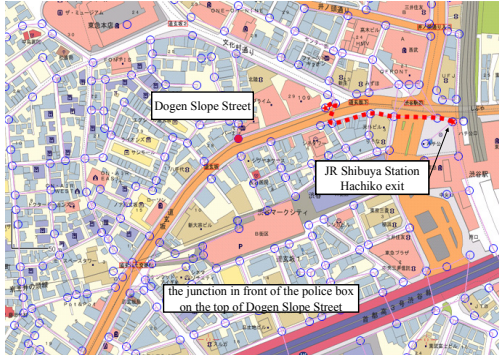
**Fig.9(a).** Our prototype system searches the shortest path (the dashed line) from the node matched *node_desc*(0) to the nearest one of the nodes making up *node_desc*(1). This figure shows the result of processing *node_desc*(0)+ *link_desc*(0)+ *node_desc*(1).
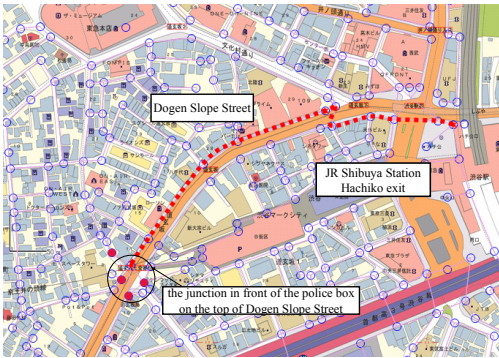


**Fig.9(b).** It searches the shortest path to the nearest one of nodes of *node_desc*(2). This figure shows the result of processing *node_desc*(0)+ *link_desc*(0)+ *node_desc*(1) + *link_desc*(1)+ *node_desc*(2).
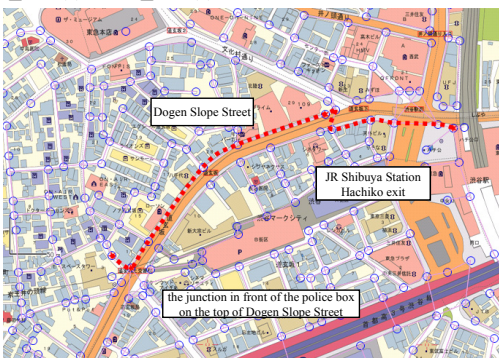


**Fig.9(c).** The end node is deduced from the description of *link_desc*(2) and the direction of its incoming link. This figure shows the result of processing *node_desc*(0)+ *link_desc*(0)+ *node_desc*(1) + *link_desc*(1)+ *node_desc*(2) + *link_desc*(2).

## 6. Summary, Conclusion and Future Work

In this research, we proposed a basic framework to geocode natural route descriptions for walkers by means of sidewalk network databases. On the basis of the structure of route description, we proposed three frameworks are as follows:

- Formal Route Statement

- Schema of sidewalk network database

- A method to validate Formal Route Statement

Formal Route Statement is able to represent a natural route description as a formal statement. We design the schema of sidewalk network databases in order to process Formal Route Statement using computers. This database is indispensable to validate Formal Route Statement. Validating methods for FRS verify the spatial validity between places. The prototype system has been developed to prove the validity of our proposed framework. The primary function of the prototype system is a naive database management system for sidewalk network databases. Also, its secondary functions include both FRS Validator and FRS Visualizer using our validating method.

Our research is converting a route description in natural language into Formal Route Statement which is a geometric description using network database. It is distinguish from the previous research of generating route descriptions using network database in car/human navigation service.

The system is able to process natural route descriptions composed of spatial anchors and spatial relationship descriptions which are stored in the database. In other words, the current system cannot recognize a lot of different ways. Furthermore the expressions of distance (e.g., "go about 100 meters") and direction using absolute indicator (e.g. "go up north", "go to the direction of Harajuku") are also tasks that exceed above the system's ability.

We plan to extent the sidewalk network database, the processible expressions in the system and implement the other functions. In the sidewalk network database, the database needs to store z-coordinates. For example, descriptions after spatial anchors can be anticipated by using the knowledge of not only the

class of spatial anchors but also geographic features stored in the sidewalk network databases. In walking starting from a lowest place, it is clear that we cannot go down there. Furthermore, with street grade, width and paving as attributes of link, the system presents the route in consideration of user's preference or physical ability. The system can present women who do not want to walk through poorly-lighted street at midnight and elderly people who also want to avoid hard slope and long stairway with information of the best route for them.

We aim for the realization of the robust system. For example, the system corrects route descriptions when invalid route descriptions are detected. In addition, there is a problem about ambiguity of natural language. In that case, multiple solutions, which are possible routes to a destination, should be ranked by some criteria of quality of the geocoding results.

This study proposed one of the methods of geocoding to convert geo-reference text data into tuples of longitude and latitude. As a result of establishment of these methods, we will be able to get local information using the spatial representation among ever-increasing digital data.

## References

[1] Takeshi Sagara: Studies for advanced practical use of non-structured and semi-structured data (in Japanese), Doctorial dissertation (2003).

[2] Shobunsha Publications Inc., http://www.mapple.co.jp/

[3] Annette Herskovits: Language and spatial cognition, Cambridge University Press, Melbourne (1986).

[4] Makoto Nagao, Satoshi Sato, Sadao Kurahashi, Tatsuhiko Tsunoda: Natural language processing (in Japanese), Makoto Nagao (edit), Iwanami Shoten, Tokyo (1996).