

# Walk this way: Spatial grounding for city exploration<sup>1</sup>

Johan Boye<sup>\*</sup>, Morgan Fredriksson<sup>†</sup>, Jana Götze<sup>\*</sup>,  
Joakim Gustafson<sup>\*</sup> and Jürgen Königsmann<sup>†</sup>

<sup>\*</sup>KTH, School of Computer Science and Communication, 100 44 Stockholm, Sweden

<sup>†</sup>Liquid Media, Hammarby allé 34, 120 61 Stockholm, Sweden

## 1 Introduction

Recent years have seen an immense proliferation of smartphones among the general public. Smartphones feature an open computing platform and GPS satellite tracking facilities, and coupled with geographic databases they allow for the creation of spatially-aware applications like routing a pedestrian from A to B (see e.g. Google Navigation (Google 2012a), or providing information of sites and services in the immediate vicinity (see e.g. Google Maps (Google 2012b)). Though some services and experimental systems rely on spoken *output* (Bartie et al. 2006, Krug et al. 2003), so far no such spatially-aware service has been based on spoken *dialogue* (e.g. the possibility for the user to intervene and ask “Should I turn left here?” or “What street am I walking on?” etc.). Furthermore, the advantage of the spoken-dialogue approach over a map-based approach is that many people find interpreting maps on a small screen to be strenuous and confusing (Looije et al, 2007). It is therefore safe to say that well-functioning spoken dialogue would be a valuable contribution to the plethora of spatially-aware mobile applications.

This paper describes an implemented dialogue system for helping a user explore the city of Stockholm. The system can either guide the user to a location of his choice (“I want to go to Odengatan”), or to specific spots chosen by the system, like a statue or an interesting architectural detail on a particular building. The latter setting in particular is interesting as it allows us to investigate various methods for producing referring spatial expressions, in order to help the user find quite small objects in a complex city environment.

In general, the city exploration problem addressed here is challenging since it involves the interpretation and generation of utterances within a rapidly changing spatial context under uncertainty.

## 2 Background

Many researchers within cognitive psychology have investigated how people give route instructions to one another (see e.g. Denis et al., 1999), and what the elements of a good route description are (see e.g. Lovelace et al. 1999, Tom and Denis, 2003). It is however not clear how these results transfer to computational

---

<sup>1</sup> Supported by the European Commission, project *Spacebook*, grant no 270019.

models of route description generation. One finding is that a big portion of such dialogues are devoted to grounding; making sure that the dialogue partner actually sees and understands what is being referred to. Grounding is a well-studied phenomenon also in dialogue systems (see e.g. Traum 1999, Skantze 2007).

The implemented systems for guiding pedestrians have mostly been based on spoken output from the system, with little or no possibility for the user to provide information (Malaka and Zipf 2000, Krug et al. 2003, Jöst et al. 2005, Bartie et al 2006, Zipf and Jöst 2005). Spoken dialogue systems in spatial domain have mostly focused on non-dynamic contexts where the user can ask questions about a static map (e.g. Cai et al. 2003, Wang et al. 2008), on virtual environments such as computer games (e.g. Boye and Gustafson 2005, Boye et al. 2006, Skantze et al. 2006, Striegnitz et al 2011), in indoor environments (Cuayáhuil et al 2011), or on natural-language interfaces to robots (Lemon et al. 2001, MacMahon et al. 2006, Johansson et al. 2011). Few if any researchers have so far addressed the topic of spoken natural-language dialogue with a user in a real, dynamic city environment.

### 3 Uncertainty and grounding

A recurring problem for any pedestrian routing system is to describe to the user how to get from his current position to the next node in the planned route. This has to be done reliably even though the user's position, speed and direction are uncertain due to possible errors in GPS readings. Giving simple instructions like "turn left here" is therefore a risky strategy; such instructions might be nonsensical for the user if he is not quite where the system believes him to be. Furthermore, the interpretation of left and right is not always clear, for instance in parks and open squares, or when the user is standing still without the system knowing which way he is facing. Therefore, before giving directions, it is often preferable that the system first *grounds* the current position with the user by means of reference landmarks in the near vicinity.

Consider for instance the situation depicted in Figure 1a. Here the system seeks to describe the route given by the route planner, first to node A, then (when the user reaches A) to B, then down a flight of stairs to node C, then turn left, etc. Before giving instructions, our system first calculates if there is a clear line of sight from the user's assumed position to a number of reference landmarks. It then selects the most salient landmark, seeks to make the user aware of it, and describes the route relative to it. Here is a sample dialogue:

1. **System:** There is a fountain about 35 metres from here. Can you see it?
2. **User:** Yes.
3. **S:** Good! Please walk to the left of the fountain.  
(*user walks*)
4. **S:** Please turn right and walk to the top of the stairs.
5. **U:** What?
6. **S:** There is a flight of stairs leading down about 25 metres from here. Can you see it?



**Figure 1a, (left):** Generating route instructions. **Figure 1b (right):** GPS drift. Black circles represent actual positions; white circles represent reported GPS positions.

In utterance 1, since there is no good way of describing node A, the system cannot ask directly about it. Instead, the system calculates that there are two describable landmarks visible from the user's presumed position; a fountain and an archway, of which the fountain is considered most salient. When the user confirms (utterance 2), the system gives the next instruction with a reference to the fountain. If the user would have answered in the negative, the system would have proceeded to ask about another visible landmark. If all possibilities are exhausted, the user is asked to simply start walking, so the system can adjust his course if needed.

Determining salience and producing good referential expressions is a difficult problem in general. Salience measures used by our system include rarity (rare objects such as fountains are more salient than entrances to buildings), distance, uniqueness, and familiarity (objects that have been mentioned before in the dialogue are considered more salient, and are easily described, e.g. "the fountain that you passed before").

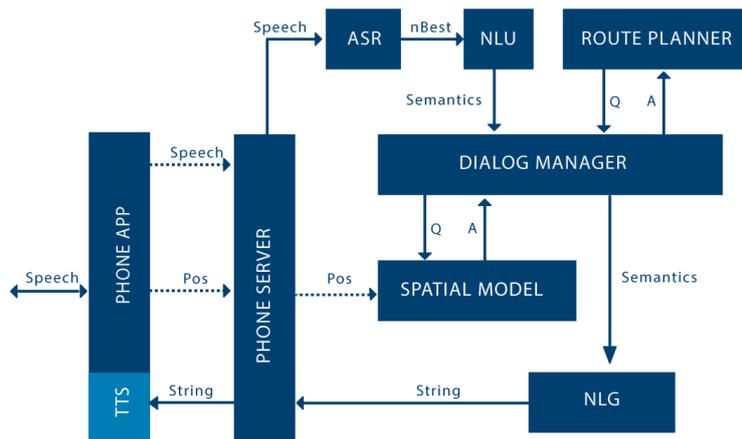
The system knows the user's position by means of the GPS receiver in the user's Android device. Unfortunately, the so-called "canyon effect" (Boriello et al. 2005) can introduce inaccuracies into GPS readings, and these errors can be quite substantial. In our tests, the GPS receiver could report positions 60 metres or more from the correct position. In addition, these errors usually magnify gradually, the reported position "drifting away" from the correct position (see Figure 1b), making errors harder to detect.

In the solution we are currently pursuing, the system is either in 'certain' or 'uncertain' mode, depending on the perceived reliability of the reported GPS coordinates. In each time step, the system computes whether the line connecting the two latest reported positions intersects any polygon representing a building. If it does, the latest coordinate is surely wrong, since the user cannot walk on top of buildings. The system then enters 'uncertain mode', in which the system cannot reliably know when the next waypoint is reached. It can then optionally produce a

grounding-seeking utterance (e.g. “When you reach the next intersection, please say ‘I’m there’.”). As soon as the user gives this feedback, the system will check whether the current GPS readings also show that the user is at the reported position. If they do, the system again enters ‘certain mode’, and will start relying on the GPS readings again.

## 4 System architecture

The system is implemented to work speech-only and “eyes-free” – the user should not need to look down on a map on the screen, but rather be free to experience the city. The architecture described here is used both for the fully automatic system and for a Wizard-of-Oz data collection that preceded it. In the latter case, an operator GUI took the place of the Dialogue Manager. The operator GUI showed the user’s position as a colored dot on a map, and used Google street view to show an approximation of the user’s visual context.



**Figure 2:** System architecture.

The user downloads a client app to his Android device which once started connects to a central phone server. The client app sends the sound stream from the microphone of the Android device to the phone server, and as soon as contact has been established with the GPS satellites, it also starts sending the coordinates (i.e. latitude-longitude pairs) of its current location. The sound stream is sent to speech recognition and parsing, and a semantic expression representing the utterance is sent to the Dialogue Manager (DM). Coordinates are sent directly to the DM. The dialogue manager updates its context model based on the input, and decides what to say to the user, and when to say it. The DM may also call an external planner to compute a route between two points in the city. For user speech input we are currently using a commercial off-the-shelf speech recognizer with a hand-written language model. For speech output, we use the built-in speech synthesizer on the An-

droid device. The architecture also supports the use of server-side speech synthesis streamed to the handset, as well as the speech recognizer to be run on the handset. The latter feature would make it possible to maintain a dialogue in places where the 3G data connection fails.

Coordinates are also sent to the Spatial Model, which is a module that maintains the mapping from the logical representation of the city (in terms of buildings, streets, etc.) to the algebraic representation (in terms of polygons, lines, and coordinates). The DM can also call the Spatial Model to perform visibility calculations to find out whether there is a free line-of-sight between two given points.

The Spatial Model automatically generates its polygon representation of the city from an export from OpenStreetMap (Haklay, 2008), generated by indicating an area on the map. A minimal bounding rectangle is computed for each polygon in order to speed up visibility calculations, as it is faster to compute whether a line intersects a rectangle than an arbitrary polygon. If the Dialogue Manager needs to find out if B is visible from A, a request is sent to the Spatial Model, which first computes whether the line AB intersects any bounding rectangle in the entire city representation. If not, there is a clear line of sight from A to B. If the line intersects a bounding rectangle, a second more expensive calculation is carried out to check whether AB intersects the polygon inside the rectangle.

Such visibility calculations are currently used for three purposes: Firstly, as mentioned in section 3, they are used to detect GPS drift. Secondly, they are used to find candidates for referring expressions (the objects of which have to be visible from the user's assumed position). Thirdly, they are used to produce better route plans. Street objects in OpenStreetMap may contain many nodes very close to each other (in particular in roundabouts or curved streets), and consequently route plans can become very long. By iteratively weeding out any node visible from the preceding node, route plans become more suitable for narration.

## 5 Concluding remarks

The system presented here routes pedestrians to their destination, using spoken dialogue to first ground reference landmarks used in the routing instructions. Ongoing work includes user tests in a part of Stockholm in order to assess and improve the implemented strategies.

## References

- Bartie, P. and Mackaness, W (2006). Development of a speech-based augmented reality system to support exploration of cityscape. *Transactions in GIS*, 10(1):63–86.
- Borriello, G., Chalmers, M., LaMarca, A. and Nixon, P. (2005) Delivering real-world ubiquitous location systems. *Communications of the ACM*, vol 8, issue 3, pp. 36-41.
- Boye, J. and Gustafson, J. (2005) How to do dialogue in a fairy-tale world. *Proceedings of the 6<sup>th</sup> SIGDial workshop on discourse and dialogue*, Lisbon, Portugal.
- Boye, J., Gustafson, J. and Wirén, M. (2006) Robust spoken language understanding in a computer game. *Journal of Speech Communication*, 48, pp. 335-353.

- Cai, G. Wang, H. and MacEachren, A. (2003) Communicating Vague Spatial Concepts in Human-GIS Interactions: A Collaborative Dialogue Approach . In *Spatial Information Theory: Foundations of geographic information science*, LNCS Volume 2825/2003, pp 287-300.
- Cuayáhuítl, H. and Dethlefs, N. (2011) Spatially-Aware Dialogue Control Using Hierarchical Reinforcement Learning. In *ACM Trans. on Speech and Language Processing (Special Issue on Machine Learning for Adaptive Spoken Dialogue Systems)*, vol. 7, no. 3, pp. 5:1-5:26
- Denis, M., Pazzaglia, F., Cornoldi, C. and Bertolo, L. (1999) Spatial discourse and navigation: an analysis of route directions in the city of Venice. *Applied cognitive psychology*, vol 13, no 2.
- Google Inc. (2012a) Google Maps Navigation. <http://www.google.com/mobile/navigation/>.
- Google Inc. (2012b) Google Maps for mobile. <http://www.google.com/mobile/maps>.
- Gustafson, J, Bell, L, Beskow, J, Boye, J, Carlson, R, Edlund, J, Granström, B, House, D and Wirén M (2000) AdApt – a multimodal conversational dialogue system in an apartment domain, *Proceedings of ICSLP 00*.
- Haklay, M. (2008) OpenStreetMap: User-generated street maps. *Pervasive computing IEEE*, vol. 7, issue 4, pp. 12-18.
- Johansson, M, Skantze, G, and Gustafson, J. (2011) Understanding Route Directions in Human-Robot Dialogue. *Proc SemDial 2011: Proceedings of the 15th Workshop on the Semantics and Pragmatics of Dialogue*, pages 19–27.
- Jöst, M., Häußler, J., Merdes, M. and Malaka, R. (2005) Multimodal interaction for pedestrians: an evaluation study. In *IUI '05: Proceedings of the 10th international conference on Intelligent user interfaces*, pp 59–66.
- Krug, K., Mountain, D. and Phan, D. (2003) Webpark: Location-based services for mobile users in protected areas. *GeoInformatics*, pp. 26–29.
- Lemon, O., Bracy, A., Gruenstein, A., and Peters, S.(2001) A Multi-Modal Dialogue System for Human-Robot Conversation , In *Proc. NAACL*.
- Looije, R., te Brake, G. and Neerincx, M. (2007) Usability engineering for mobile maps. In *Proceedings of Mobility'07, 4<sup>th</sup> int. conference on mobile technology, applications, and systems*.
- Lovelace, K., Hegarty, M. and Montello, D. (1999) Elements of good route descriptions in familiar and unfamiliar environments. *Spatial Information Theory. Cognitive and computational foundations of geographic information science*, LNCS, 1661/1991, Springer-Verlag.
- MacMahon, M., Stankiewicz, B. and Kuijpers, B. (2006) Walk the Talk: Connecting Language, Knowledge, Action in Route Instructions. *National Conf on Artificial Intelligence (AAAI-06)*.
- Malaka, R. and Zipf, A.. (2000) Deep map – challenging IT research in the framework of a tourist information system. *Information and communication technologies in tourism*, Springer.
- Skantze, G. (2007). Making grounding decisions: Data-driven estimation of dialogue costs and confidence thresholds. In *Proceedings of SigDial* (pp. 206-210). Antwerp, Belgium.
- Skantze, G., Edlund, J., & Carlson, R. (2006). Talking with Higgins: Research challenges in a spoken dialogue system. *Perception and Interactive Technologies*, pp. 193-196. Springer
- Striegnitz, K., Denis, A., Gargett, A. Garoufi, K., Koller, A. and Theune, M. (2011) Report on the Second Challenge on Generating Instructions in Virtual Environments (GIVE-2.5). In *Proceedings of the Generation Challenges Session at the 13th European Workshop on Natural Language Generation (ENLG)*.
- Tom, A. and Denis, M. (2003) Referring to landmark or street information in route directions: What difference does it make?. *Spatial Information Theory: Foundations of geographic science*. LNCS 2825/2003, pp. 362-374. Springer-Verlag.
- Traum, D. (1999) *Computational Models of Grounding in Collaborative Systems*. AAAI Technical Report FS-99-03
- Wang, H., Cai, G. and MacEachren, A. (2008) GeoDialogue: A Software Agent Enabling Collaborative Dialogues between a User and a Conversational GIS . In *Tools with Artificial Intelligence, 2008. ICTAI '08. 20th IEEE*.
- A Zipf and M. Jöst. (2005) Implementing adaptive mobile GI services based on ontologies - examples for pedestrian navigation support. *Computers, Environment and Urban Systems*, Special Issue on LBS and UbiGIS., 2005.