



# A Speech-Centric Perspective for Human-Computer Interface: A Case Study

LI DENG AND DONG YU

Microsoft Research, One Microsoft Way, Redmond, WA 98052, USA

Received xx, xx; Accepted xx, xx

Published online: xx xx

**Abstract.** Speech technology has been playing a central role in enhancing human-machine interactions, especially for small devices for which graphical user interface has obvious limitations. The speech-centric perspective for human-computer interface advanced in this paper derives from the view that speech is the only natural and expressive modality to enable people to access information from and to interact with any device. In this paper, we describe some recent work conducted at Microsoft Research, aimed at the development of enabling technologies for speech-centric multimodal human-computer interaction. In particular, we present a case study of a prototype system, called MapPointS, which is a speech-centric multimodal map-query application for North America. This prototype navigation system provides rich functionalities that allow users to obtain map-related information through speech, text, and pointing devices. Users can verbally query for state maps, city maps, directions, places, nearby businesses and other useful information within North America. They can also verbally control applications such as changing the map size and panning the map moving interactively through speech. In the current system, the results of the queries are presented back to users through graphical user interface. An overview and major components of the MapPointS system will be presented in detail first. This will be followed by software design engineering principles and considerations adopted in developing the MapPointS system, and by a description of some key robust speech processing technologies underlying general speech-centric human-computer interaction systems.

**Keywords:** human-computer interaction, speech-centric multimodal interface, robust speech processing, MapPointS, speech-driven mobile navigation system

## 1. Introduction

Speech recognition technology enables a computer to automatically convert an acoustic signal uttered by users into textual words, freeing them from the constraints of the standard desktop-style interface (such as mouse pointer, menu, icon, and window etc.). The technology has been playing a key role in enabling and enhancing human-machine communications. In combination with multimedia and multimodal processing technologies, speech processing will in the future also contribute, in a significant way, to facilitating human-human interactions. In applications

such as distributed meetings, audio-visual browsing, and multimedia annotations, automatic processing of natural, spontaneous speech will collaborate with automatic audio-visual object tracking and other multimedia processing techniques to complete full end-to-end systems. In addition to the multimedia applications, the most important role that speech can play is in a full range of the devices that demand efficient human inputs. Since speech is the only natural and expressive modality for information access from and interaction with any device, we highlight the speech-centric view of human-machine interface (HCI).

*Deng and Yu*

49 Speaking is the most natural form of human-to-  
50 human communication. We learn how to speak in the  
51 childhood, and we all exercise our speaking communi-  
52 cation skills on a daily basis. The possibility to translate  
53 this naturalness of communication into the capability  
54 of a computer is our natural expectation, since a com-  
55 puter is indeed equipped with huge computing and  
56 storage capacities. However, the expectation that com-  
57 puters should be good at speech has not been a reality,  
58 at least not yet. One important reason for this is that  
59 speech input is prone to error due to imperfection of  
60 the technology in dealing with variabilities from the  
61 speaker, speaking style, and the acoustic environment.  
62 The imperfection, in addition to a number of social and  
63 other reasons, raises the issue that speech alone is not  
64 sufficient as the most desirable input to computers. Use  
65 of multimodal inputs in an HCI system, which fuses  
66 two or more input modalities (speech, pen, mouse, etc.)  
67 to overcome imperfection of speech technology in its  
68 robustness as well as to complement speech input in  
69 other ways, is becoming an increasingly more impor-  
70 tant research direction in HCI.

71 Major HCI modalities in addition to speech are  
72 related to graphic user interface (GUI). GUI is based  
73 primarily on the exploitation of visual information,  
74 and has significantly improved HCI by using intuitive  
75 real-world metaphors. However, it is far from the ulti-  
76 mate goal of allowing users to interact with computers  
77 without training. In particular, GUI relies heavily  
78 on a sizeable screen, keyboard, and pointing device,  
79 which are not always available. In addition, with more  
80 and more computers designed for mobile usages and  
81 hence subject to the physical size and hands-busy or  
82 eyes-busy constraints, the traditional GUI faces an  
83 even greater challenge. Multimodal interface enabled  
84 by speech is widely believed to be capable of dramati-  
85 cally enhancing the usability of computers because  
86 GUI and speech have complementary strengths.  
87 While speech has the potential to provide a natural  
88 interaction model, the ambiguity of speech and the  
89 memory burden of using speech as output modality  
90 on the user have so far prevented it from becoming the  
91 choice of mainstream interface. Multimodal Intelligent  
92 Personal Assistant Device, or MiPad, was one of our  
93 earlier attempts in overcoming such difficulties by  
94 developing enabling technologies for speech-centric  
95 multimodal interface. MiPad is a prototype of wireless  
96 Personal Digital Assistant (PDA) that enables users to  
97 accomplish many common tasks using a multimodal

spoken language interface (speech + pen + display). 98  
MiPad, as a case study for speech-centric multimodal 99  
HCI application, has been described in detail in our 100  
recent publication [2]. In this paper, we will present a 101  
second case study based on a new system built within 102  
our research group more recently, called MapPointS. 103

104 During past several years, many different methods 104  
of integrating multiple modalities (voice, visual, and 105  
others) in HCI have been proposed and implemented, 106  
and some key issues have been discussed [10–13, 16]. 107  
Many prototype systems have also been built based on 108  
the use of multiple modalities [1, 2, 7, 9, 14], most 109  
of which have focused on the special advantage of 110  
the speech input for mobile or wireless computing as 111  
in multimodal PDA's. Both of our prototype systems, 112  
MiPad and MapPointS, have such mobile computing 113  
in the special design consideration. Their design also 114  
takes the speech-centric perspective — fully exploiting 115  
the efficiency of the speech input where other modalities 116  
have special difficulties. 117

118 The focus of this paper, the prototype MapPointS, is 118  
a speech-centric, multimodal, location-related, map- 119  
query application for North America. The unique 120  
advantage of the system is its full and direct exploi- 121  
tation of the frequently updated backend database 122  
provided by the existing Microsoft product, Map- 123  
Point (<http://mappoint.msn.com>). MapPointS 124  
essentially adds the “Speech” modality and its interface into 125  
MapPoint, and hence MapPointS. MapPointS provides 126  
rich functionalities to allow the users to obtain map- 127  
related information through speech, text, and pointing 128  
devices. (MapPoint provides the same functionalities 129  
with the inputs of text and pointing devices only). 130  
With MapPointS, the users can verbally query for 131  
state maps, city maps, directions, places (e.g., school 132  
names), nearby businesses, and many other useful in- 133  
formation. They can also verbally control applications 134  
such as changing the map size and panning the map 135  
moving interactively through speech. In the current 136  
system, the results of the queries are presented back to 137  
users through GUI. An overview and the major com- 138  
ponents of the MapPointS system will be presented 139  
in detail in this paper first. Following this presenta- 140  
tion, we will describe several key software design 141  
engineering principles and considerations in devel- 142  
oping MapPointS. Finally we will present some key 143  
speech processing technologies underlying the gen- 144  
eral speech-centric HCI systems including MiPad and 145  
MapPointS. 146

## Speech-Centric Perspective for Human-Computer Interface

## 147 2. System Overview and Functionality 148 of Mappoints

149 MapPointS is a map query application that supports  
150 a large set of map query commands through speech,  
151 text, and pointing devices. These commands can be  
152 classified into the following five categories:

- 153 1. *Application Control*: Application control com-  
154 mands are used to control MapPointS applications.  
155 For example, a user can use speech (as well as other  
156 modalities) to quit the application, to pan the map  
157 towards eight directions, to zoom the maps, or to  
158 open and save the map.
- 159 2. *Location Query*: Location queries are used to search  
160 for the map of a specific location. For example, a  
161 user can query for a map with city names, state  
162 names, joint city and state names, place names (e.g.,  
163 Seattle University), or referenced locations (e.g.,  
164 here; this place; and this area, etc., which are indi-  
165 cated by the mouse click rather than by the speech  
166 input.
- 167 3. *Route Query*: Route queries are used to obtain  
168 directions from one location to another. There  
169 are two types of such queries. The first type  
170 contains both “from” and “to” information. For  
171 example, a user can say “How do I get from  
172 <startlocation> to <endlocation>” to obtain direc-  
173 tions from <startlocation> to <endlocation>. The  
174 <startlocation> and <endlocation> can be any lo-  
175 cation type specified in location query. The second  
176 type of queries contains information about “to lo-  
177 cation” only. “How may I go to <location>” is an  
178 example of such queries. When a query with “to  
179 location” only is submitted by a user, the system  
180 will infer the most probable from location based on  
181 the user’s dialog context.
- 182 4. *Nearest Query*: “Nearest” queries are used to find  
183 the closest or the nearest instance of a specific type  
184 of places to the current location. MapPointS sup-  
185 ports about 50 types of locations including bank,  
186 gas station, airport, ATM machine, restaurant, and  
187 school. For instance, a user can query for the near-  
188 est school, Chinese restaurant, etc. When such a  
189 query is made, MapPointS will infer the most prob-  
190 able current reference location based on the dialog  
191 context.
- 192 5. *Nearby Query*: “Nearby” queries are similar to the  
193 “nearest” queries above. The difference is that all  
194 nearby instances of a type of places, instead of only

one, are displayed in the nearby queries. For ex- 195  
ample, a user can query for all nearby gas stations. 196  
Similar to the situation of the nearest query, Map- 197  
PointS needs to infer the most probable reference 198  
location before executing the query. 199

Examples of the above five types of queries are pro- 200  
vided now. Figure 1 is a screen shot where a map of 201  
Seattle is displayed as a result of speech command used 202  
in the location query: “show me a map of Seattle”. A 203  
typical map of Seattle with its surroundings is imme- 204  
diately displayed. All cities in the U.S. can be queries 205  
in the same manner. 206

Figure 2 gives a multimodal interaction example 207  
where the user makes a location query by selecting 208  
an area with mouse and zooming the picture to just 209  
that part of the map while using the following simulta- 210  
neous speech command: “show me this area”. The 211  
portion of the map selected by the user is displayed in 212  
response to such a multimodal query. 213

In Fig. 3 is another multimodal interaction example 214  
for the nearest location query. In this case, the user 215  
clicks on a location, and more or less simultaneously 216  
issues the command: “Show me the nearest school” 217  
with speech. MapPointS displays “Seattle University” 218  
as the result based on the location that the user just 219  
clicked on. 220

In Fig. 4 we show an example of the route query to 221  
find the direction from Seattle to Boston, with a speech 222  
utterance such as “Show me directions from Seattle to 223  
Boston”, or “How may I go from Seattle to Boston”, 224  
etc. If the immediately previous location is Seattle, 225  
then saying just “How may I go to Boson” will give 226  
the identical display as the response to the query. 227

We provide a further example in Fig. 5 of query- 228  
ing nearby restaurants by speaking to MapPointS with 229  
“show me all nearby restaurants”. The system assumes 230  
the current location of the user based on the previous 231  
interactions, and is hence able to display all nearby 232  
restaurants without the need for the user to specify 233  
where he currently is. 234

For the system functionalities illustrated in the above 235  
description and examples, MapPointS demonstrates 236  
the following four specific features: 237

1. *Multi-Modal Human-Computer Interaction*: As we 238  
discussed in Introduction section, one of the trends 239  
of HCI is the integration of multi-modal inputs, 240  
through which speech recognition is integrated with 241  
various other modalities such as keyboard and 242

Deng and Yu

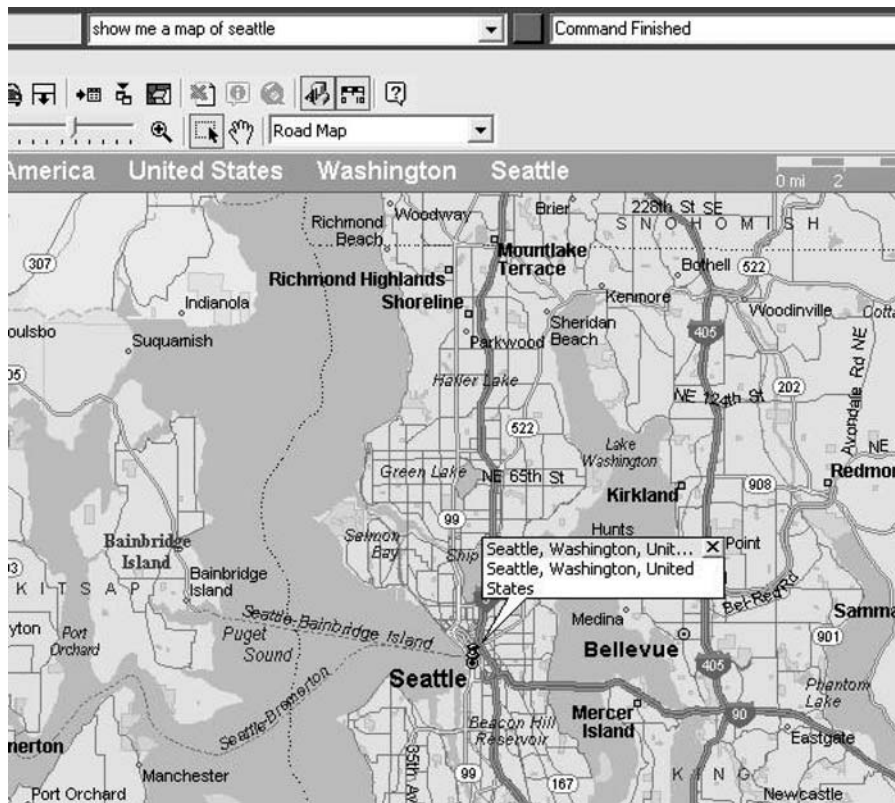


Figure 1. Navigation using voice command: “show me a map of Seattle”.

243 mouse inputs. MapPointS is a good show case for  
 244 this capability since it includes both location search  
 245 (via the name) and location pointing/selection. The  
 246 former is most naturally accomplished using voice  
 247 command because it is difficult to use a mouse or  
 248 a pen to search for one of a very large number of  
 249 items (cities, etc). The latter, location pointing and  
 250 selection, on the other hand, is relatively easy to  
 251 be fulfilled with mouse clicks. For example, a user  
 252 may ask the system to “show me a map of Seattle”.  
 253 The user can then use the mouse to click on a spe-  
 254 cific location or to select a specific area. He/she can  
 255 then or simultaneously issue the command “Show  
 256 me the nearest school around here” with speech as  
 257 the input.  
 258 2. *Integrated Interface for Speech and Text:* In the  
 259 MapPointS, a user not only can use speech to query  
 260 the application but also can use a natural text input  
 261 to ask for the same thing. For example, the user  
 262 can say “Where is the University of Washington”  
 263 to have the University of Washington be identified

in the map. Alternatively, the user can just type  
 in “Where is the University of Washington” in the  
 command bar and obtain the same result.  
 3. *Recognition of a Large Quantity of Names:* As  
 we have mentioned, MapPointS allows its users to  
 query for all cities and places in the US. Accurate  
 recognition of all these names is difficult since there  
 are too many names to be potential candidates. For  
 example, there are more than 30,000 distinct city  
 names in the US, and the total number of valid  
 combinations of “city, state” alone is already larger  
 than 100,000, not to mention all the school names,  
 airport names, etc. in all cities.  
 4. *Inference of Missing Information:* When a user  
 queries information, he/she may not specify full  
 information. For example, when a user submits a  
 query “How may I get to Seattle University”, Map-  
 PointS needs to infer the most probable location that  
 the user is currently at. This inference is automati-  
 cally performed based on the previous interactions  
 between the user and MapPointS.

Speech-Centric Perspective for Human-Computer Interface



Figure 2. User's mouse selection is seamlessly integrated into the speech command: "Show me this area".

285 **3. System Architecture and Components**  
 286 **of Mappoints**

287 The major system components of MapPointS are  
 288 depicted in Fig. 6. The raw signals generated by the  
 289 user are first processed by a semantic parser into the  
 290 "surface semantics" representation. For the speech  
 291 input, the speech recognizer first converts the raw  
 292 signal into a text sequence, with the help from the  
 293 Language Model component, before semantic parsing.  
 294 Each possible modality, speech or otherwise, has its  
 295 separate corresponding semantic parser. However, the  
 296 resulting surface semantics are represented in common  
 297 Semantic Markup Language (SML) format and is thus  
 298 independent of the modality. With this approach, the  
 299 input methods become separated from the rest of the  
 300 system. The surface semantics from all the input media  
 301 are then merged by the Discourse Manager component  
 302 into the "discourse semantics" representation. When

generating the discourse semantics, the discourse manager  
 integrates the environment information (provided  
 by the Environment Manager and Semantic Model  
 components) which includes: (1) dialog context;  
 (2) domain knowledge; (3) user's information, and  
 (4) user's usage history. Such important environment  
 information is used to adapt the Language Model,  
 which improves the speech recognition accuracy and  
 enhance the Semantic Parsers for either the speech  
 or text input. (Semantic Model is the component  
 that provides rules to convert the surface semantics  
 into actionable commands and to resolve possible  
 confusibility.) The discourse semantics is then fed into  
 the Response Manager component to communicate  
 back to the user. The Response Manager synthesizes  
 the proper responses, based on the discourse semantics  
 and the capabilities of the user interface, and plays  
 the response back to the user. In this process, Behavior  
 model provides rules to carry out the required actions.

Deng and Yu

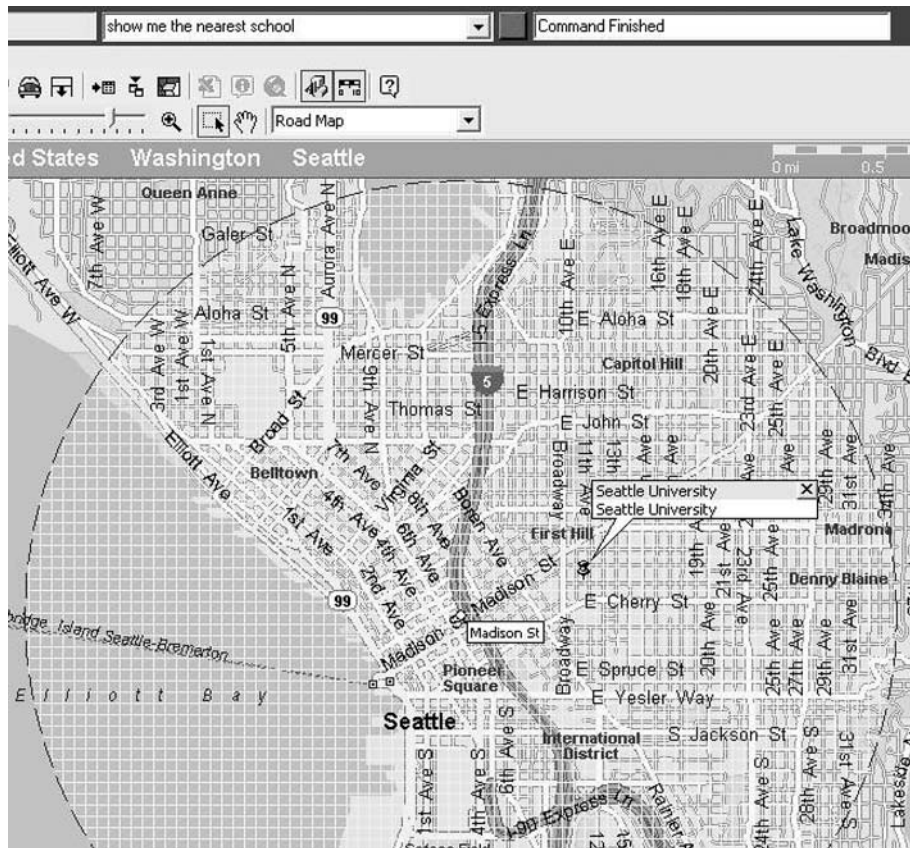


Figure 3. User's latest mouse click input is referenced by voice command: "Show me the nearest school".

322 We have already introduced some components of the  
 323 above main architecture in some of our earlier publi-  
 324 cations (e.g., [2]). In this paper, we focus on two novel  
 325 components of the architecture: Language Model (LM)  
 326 and Environment Manager. The design of these two  
 327 components has been specific to the MapPointS sys-  
 328 tem.

329 As we pointed out in the previous section, one of  
 330 the major difficulties of the task is the recognition of  
 331 the very large quantity of names. Including all names  
 332 in the grammar is infeasible because the total number  
 333 of names is so large that the confusability between  
 334 these names is extremely high and the computation for  
 335 speech recognition search is very expensive.

336 The speech recognition task is conducted as an  
 337 optimization problem to maximize the posterior  
 338 probability:

$$\hat{w} = \arg \max_w P(A | w)P(w),$$

where  $w$  is a candidate word sequence, and  $P(w)$  is  
 the prior probability for the word sequence (or LM  
 probability). This suggests that we can reduce the  
 search effort through controlling the language model  
 so that only the most probable names are kept in the  
 search space. One of the approaches used to better  
 estimate  $P(w)$  is to exploit the user information,  
 especially the user's home address, usage history,  
 and current location. In other words, we can simplify  
 the speech recognition search task by optimizing the  
 following posterior probability:

$$\hat{w} = \arg \max_w P(A | w)P(w | E),$$

where the general LM  $P(w)$  is now refined (i.e.,  
 adapted) to the Environment-specific LM  $P(w | E)$ ,  
 which has a much lower perplexity than the otherwise  
 generic LM. (This environment-specific LM is  
 closely related to topic-dependent LM or user-adapted  
 LM in the literature.) How to exploit the user

Speech-Centric Perspective for Human-Computer Interface

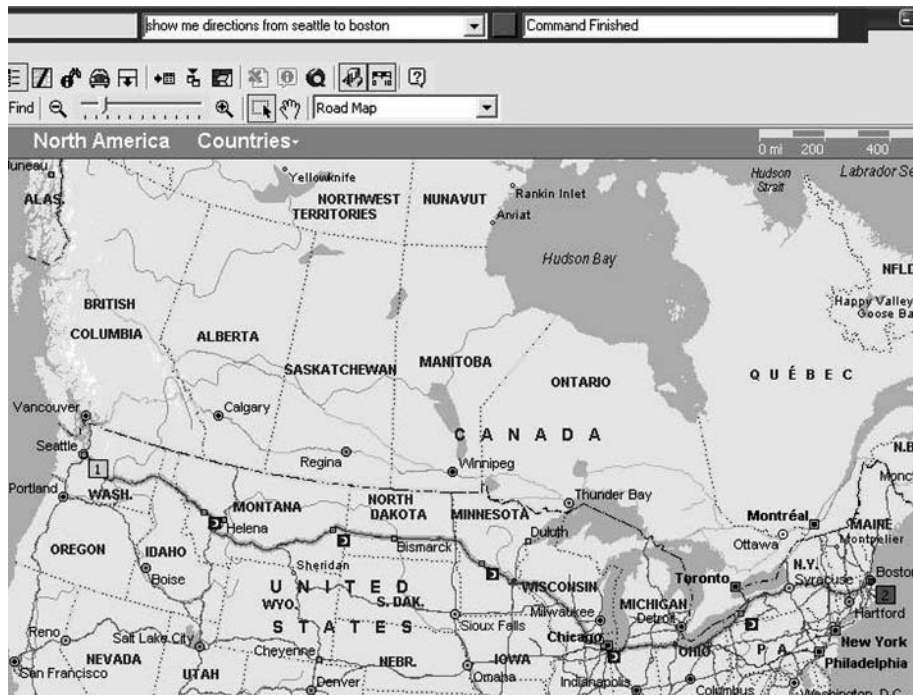


Figure 4. Route query to find direction from Seattle to Boston by speaking to MapPointS: “How may I go from Seattle to Boston”, or just “How may I go to Boston” if the current location is Seattle.

356 “environment” information to adapt the LM is the job  
 357 of the “Environment Manager” component in Fig. 1,  
 358 which we describe in detail in the remainder of this  
 359 section.

360 In the current MapPointS system, the PCFG (Prob-  
 361 abilistic Context Free Grammar) is used as the  
 362 LM. Some examples of the CFG rules are shown  
 below:

```

<query> → <app_query> | <pan_query> | <zoom_query> |
          <location_query> | <route_query> |
          <nearest_query> | <nearby_query> | ...
<location_query> → show me <location> | show me a map
                  of <location> | where is <location> | ...
<location> → <pointed_location> | <named_location> | ...
<pointed_location> → here | this point | this | this place | ...
<named_location> → <city> | <state> | <city_state> | <well-
                  known_place> | ...
<city> → New York City | Seattle | Dallas | ...
    
```

363  
 364 In order to build the environment-adapted LM based  
 365 on the PCFG grammar, the LM probability  $P(w | E)$   
 366 is decomposed into the product of the words that make  
 367 up the word sequence  $w$  and that follow the grammar  
 368 at the same time. The majority of the words which

are relevant to LM in our MapPointS system are the  
 names or name phrases such as “New York City” in  
 the above CRG rules. (Many non-name words in the  
 grammar are provided with uniform LM probabilities  
 and hence they become irrelevant in speech recognition  
 and semantic parsing.)

We now describe how the conditional probability of  
 a name or name phrase given the environment (user  
 information is computed by the Environment Manager  
 component of MapPointS. Several related conditional  
 probabilities are computed in advance based on well  
 motivated heuristics pertaining to the MapPointS task.  
 First, it is noted that users tend to move to a city before  
 querying for small and less-known locations inside  
 that city. On the other hand, they often move between  
 cities and well-known places at any time. In other  
 words, small places (such as restaurants) in a city,  
 except for the city that the user is looking at currently,  
 have very small prior probabilities. Cities, well-known  
 places, and small places in the currently visited city,  
 in contrast, have much higher prior probabilities. For this  
 reason, we organize all names into two categories: the  
 global level and the local level. The global-level name  
 list contains state names, city names, City+State,  
 and well-known places such as Yellowstone National

Deng and Yu

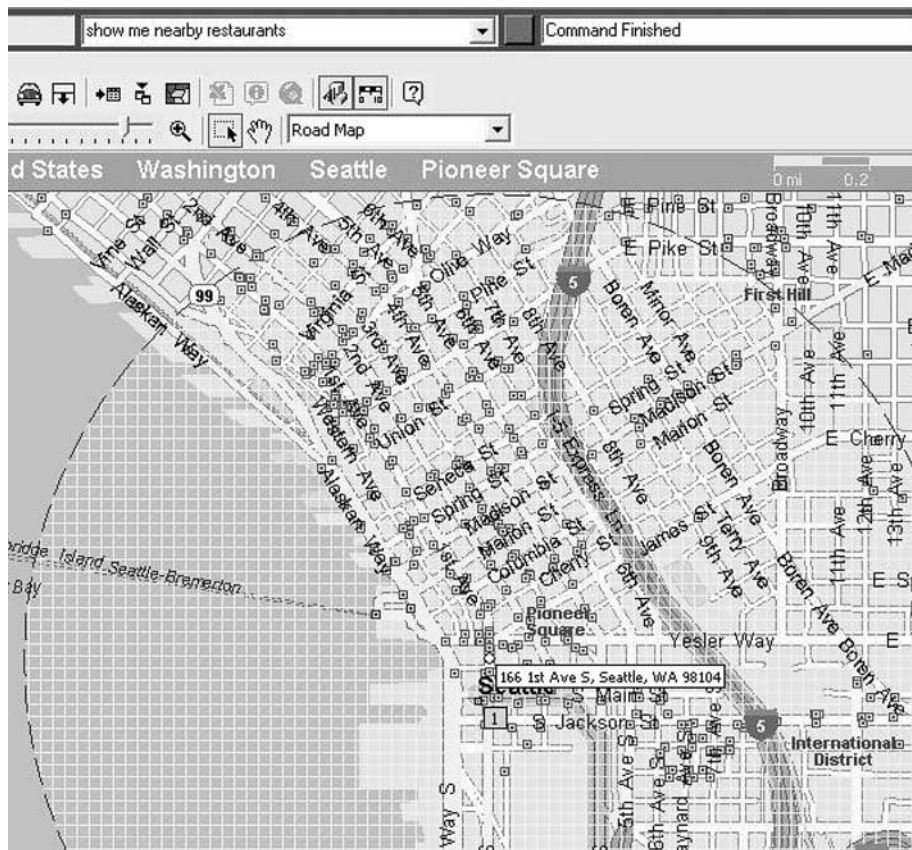


Figure 5. Display of MapPointS in response to the “Nearby Restaurants” query.

394 park. This global-level name list is included in the  
 395 recognition grammar at all times. The local-level  
 396 name list, on the other hand, contains detailed location  
 397 information about a city or a well-known place. When  
 398 the current city is changed, the local-level name list is  
 399 changed accordingly.

400 To speed up the loading of the local-level name list,  
 401 we pre-built the local list for each of the 2000 major  
 402 cities. This is needed because there are usually many  
 403 place names in large cities and these lists are slow to  
 404 build. For local-name lists of small cities, we build  
 405 them when the city is firstly visited and cache the lists  
 406 in the hard drive in order to speed up the process when  
 407 it is visited again.

408 Even after adopting this approach, the number of  
 409 names is still large. The majority of the names in the  
 410 global-level name list are for cite and state combination  
 411 (City+State). The simplest way to include these names  
 412 in the grammar would be to list them all one by one.  
 413 This, however, requires more than 100,000 distinct

entries in the grammar. Typical recognition engines 414  
 can not handle the grammars of such a size efficiently 415  
 and effectively. We thus take a further approach to 416  
 arrange the cities and states in separate lists and allow 417  
 for combinations of them. This approach greatly 418  
 reduces the grammar size since we only need 30,000 419  
 cities and 50 states. Unfortunately, this will provide 420  
 invalid combinations such as “Seattle, California”. 421  
 It also increases the name confusability since now 422  
 there are more than  $30,000 \times 50 = 1,500,000$  possible 423  
 combinations. To overcome this difficulty, we choose 424  
 to list only valid City+State combinations. To accom- 425  
 plish this, we prefix the grammar so that all names 426  
 are organized based on the city names, and each city 427  
 name can only follow the valid subset of the 50 state 428  
 names. The prefixed grammar can be processed by 429  
 recognition engines rather efficiently. For some slow 430  
 systems where the speed and accuracy may still be in- 431  
 adequate, we further pruned the number of City+State 432  
 combinations. 433



Speech-Centric Perspective for Human-Computer Interface

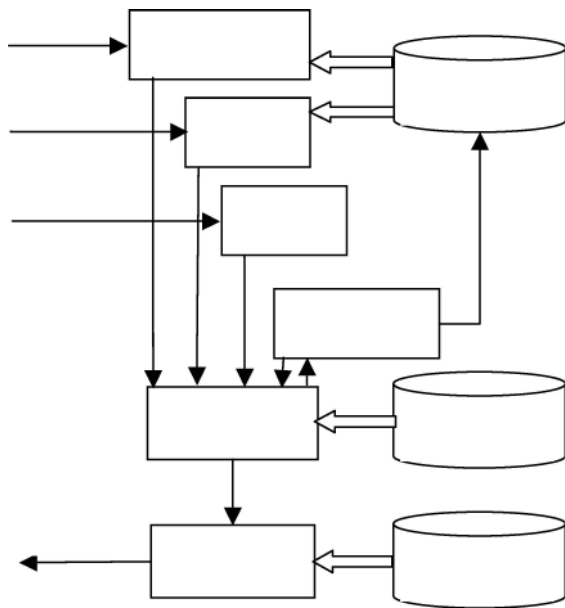


Figure 6. Major system architecture and components in MapPointS.

Table 1. Full list of location classes in MapPointS.

Class ID	Class Type
1	State
2	City
3	Well-known Places
4	Galleries
5	ATMs and banks
6	Gas stations
7	Hospitals
8	Hotels and motels
9	Landmarks
10	Libraries
11	Marinas
12	Museums
13	Nightclubs and taverns
14	Park and rides
15	Police stations
16	Post offices
17	Rental car agencies
18	Rest areas
19	Restaurants—Asian
20	Restaurants—Chinese
21	Restaurants—delis
22	Restaurants—French
23	Restaurants—Greek
24	Restaurants—Indian
25	Restaurants—Italian
26	Restaurants—Japanese
27	Restaurants—Mexican
28	Restaurants—pizza
29	Restaurants—pizza
30	Restaurants—seafood
31	Restaurants—Thai
32	Schools
33	Shopping
34	Casinos
35	Stadiums and arenas
36	Subway stations
37	Theaters
38	Airports
39	Zoos

434 The second heuristic adopted by the MapPointS  
 435 system is motivated by the intuition that if a user queries  
 436 restaurants a lot, the probability that he/she will query  
 437 new restaurants should be high even though they have  
 438 not been queried before. With this heuristic, we or-  
 439 ganize all names into about 40 classes including gas  
 440 stations, schools, restaurants, airports, etc. A complete  
 441 list of the classes can be found in Table 1.

442 We denote the probability that a class of names is  
 443 queried as  $P([Class]|History)$  or  $P([C]|H)$ . The esti-  
 444 mate for this probability is provided as in the Map-  
 445 PointS system:

$$P([C_i] | H) = \frac{\sum_k \exp(-\lambda_h(T - t_{ik}))}{\sum_j \sum_k \exp(-\lambda_h(T - t_{jk}))}$$

446 where  $t_{ik}$  is the time the names in class  $C_i$  was queried  
 447 the  $k$ -th time (as the “History” information),  $T$  is the  
 448 current time, and  $\lambda_h$  is the forgetting factor. We further  
 449 assume that  $[C_i]$  is independent of other factors in the  
 450 environment. This particular form of the probability  
 451 we have adopted says that the further away a past class  
 452 query is, the less it will contribute to the probability of  
 453 the current class query.

454 The third heuristic we have adopted is motivated  
 455 by the intuition that even though names in the global-  
 456 level name list are likely to be queried by users, the  
 457 probabilities that each name would be queried will be

different. For example, large cities such as San 458  
 Francisco and Boston are more likely to be queried 459  
 than small cities such as Renton. For this reason, 460  
 we estimated the prior probabilities of all cities and 461

Deng and Yu

462 well-known places in advance. The estimation is based  
 463 on the MapPoint.NET (<http://mappoint.msn.com/>)  
 464 IIS (Internet Information Server) log data. The IIS  
 465 log records raw queries users of the MapPoint.NET  
 466 submitted (The log, however, does not contain any  
 467 user identification information).

468 We processed more than 40GB of the log data  
 469 to obtain statistics of states, cities, and well-known  
 470 places that users have queried. We found that for the  
 471 cities, the probability computed by the log data is quite  
 472 similar to that estimated based on the city population.  
 473 We denote the probability for each name in the  
 474 class given the class label as  $P(N|[C])$ ; examples are  
 475  $P(\text{Name}|\text{[Class]}=\text{'City'})$  and  $P(\text{Name}|\text{[Class]}=\text{'Well-}$   
 476  $\text{KnownPlace'})$ . For local-level names, we assume a  
 477 uniform distribution for  $P(N|[C])$ . Tables 2 and 3  
 478 show the most frequently queried 10 States and cities  
 479 respectively:

480 The fourth heuristic implemented in the MapPointS  
 481 system uses the intuition that location names related  
 482 to the user are more likely to be queried than other  
 483 names. For example, if a user lives in the Seattle, he/she  
 484 is more likely to query locations in or close to the  
 485 Seattle. We calculate this probability class by class.  
 486 We denote this probability as  $P(\text{Name}|\text{[Class]},\text{User})$  or  
 487 simply  $P(N|[C],U)$  and estimate it according to:

$$P(N_i | [C_k], U) = \frac{S(N_i | [C_k], U)}{\sum_{j:N_j \in [C_k]} S(N_j | [C_k], U)}$$

488 where

$$S(N_i | [C_k], U) = \exp(-\lambda_u d_{iU})P(N_i | [C_k]),$$

Table 2. Top 10 States queried by users of MapPoint.NET and their estimated probabilities.

Top no.	Name	Occurrence in IIS log	Relative frequency
1	California	2950295	0.127832
2	Texas	1791478	0.009605
3	Florida	1512045	0.065515
4	New York City	1117964	0.048440
5	Pennsylvania	1074052	0.046537
6	Illinois	1024543	0.044392
7	Ohio	1006874	0.043626
8	New Jersey	782871	0.033920
9	Michigan	776841	0.033660
10	Georgia	738660	0.032005

Table 3. Top 10 cities queried by users of MapPoint.NET and their estimated probabilities.

Top #	Name	Occurrence in IIS log	Relative Frequency
1	Houston, Texas	309246	0.014637
2	Chicago, Illinois	202948	0.009605
3	Dallas, Texas	169710	0.008032
4	Los Angeles, California	166005	0.007857
5	San Diego, California	141622	0.006656
6	Atlanta, Georgia	140637	0.006656
7	Orlando, Florida	135911	0.006433
8	San Antonio, Texas	122723	0.005809
9	Seattle, Washington	115550	0.005469
10	Las Vegas, Nevada	113927	0.005392

and  $d_{iU}$  is the distance between  $N_i \in C_k$  and the user's home.  $\lambda_u$  is the corresponding decaying parameter.

The fifth heuristic uses the intuition that locations close to the currently visited city are more likely to be queried than other locations. Following the same example, if the user lives in Seattle, not only is he/she more likely to query Bellevue than Springfield, but he/she is also more likely to query for "Everett, Washington" than "Everett, Massachusetts". We denote this probability as  $P(\text{Name}|\text{[C]},\text{CurrentLocation})$  or simply  $P(N|[C],L)$  and estimate it as:

$$P(N_i | [C_k], L) = \frac{S(N_i | [C_k], L)}{\sum_{j:N_j \in C_k} S(N_j | [C_k], L)}$$

where

$$S(N_i | [C_k], L) = \exp(-\lambda_l d_{iL})P(N_i | [C_k]),$$

and  $d_{iL}$  is the distance between  $N_i \in C_k$  and the current location.  $\lambda_l$  is the corresponding decaying factor.

The final, sixth heuristic we adopted is based on the intuition that if a user queries a location often recently, he/she is likely to query the same location again in the near future. For example, if the user lives in Seattle, but he/she queried for "Everett, Massachusetts" several times recently, we would expect that he will more likely to query for "Everett, Massachusetts" than "Everett, Washington" even though Everett, Washington" is more close to his home. We denote

Speech-Centric Perspective for Human-Computer Interface

514 this probability as  $P(\text{Name} | [C], \text{History})$  or simply  
 515  $P(N | [C], H)$  and estimate it as:

$$P(N_i | [C_n], H) = \frac{S(N_i | [C_n], H)}{\sum_{j: N_j \in C_n} S(N_j | [C_n], H)}$$

516 where

$$S(N_i | [C_n], H) = \sum_k \exp(-\lambda_h(T - t_{ik})) P(N_i | [C_n])$$

517 and  $t_{ik}$  is the time when the name  $N_i \in C_n$  was queried  
 518 the  $k$ -th time.  $T$  is the current time, and  $\lambda_h$  is the  
 519 forgetting factor.

520 With the above assumptions and heuristics based  
 521 on well founded intuitions, we obtain the conditional  
 522 probability  $P(\text{Name} | \text{Environment})$  as:

$$\begin{aligned} P(N_i | E) &= \sum_{C_n} P(N_i | [C_n], E) P([C_n] | E) \\ &= \sum_{C_n} P(N_i | [C_n], U, L, H) P([C_n] | H) \\ &= \sum_{C_{ni}} \frac{P(N_i, U, L, H | [C_n])}{P(U, L, H | [C_n])} P([C_n] | H) \\ &= \sum_{C_{ni}} \frac{P(U, L, H | N_i, [C_n]) P(N_i | [C_n])}{P(U, L, H | [C_n])} \\ &\quad \times P([C_n] | H) \end{aligned}$$

523 We further assume that  $U, L$ , and  $H$  are independent  
 524 of each other. This leads to the approximation of

$$\begin{aligned} P(N_i | E) &\approx \sum_{C_{ni}} \frac{P(U | N_i, [C_n]) P(L | N_i, [C_n]) P(H | N_i, [C_n]) P(N_i | [C_n])}{P(U | [C_n]) P(L | [C_n]) P(H | [C_n])} P([C_n] | H) \\ &= \sum_{C_{ni}} \frac{P(N_i | U, [C_n]) P(N_i | L, [C_n]) P(N_i | H, [C_n])}{P^2(N_i | [C_n])} P([C_n] | H) \end{aligned}$$

525 We can further simplify the above equation by as-  
 526 suming that each name belongs to one class. This is  
 527 accomplished by using the location in the map—the  
 528 semantic meaning of the name as the unique identi-  
 529 fier of the name. For example, Everett can mean “Ev-  
 530 erett, Washington”, “Everett, Massachusetts”, “Everett  
 531 Cinema”, and somewhere else. In our MapPointS sys-  
 532 tem’s grammar, we allow for several different kinds of  
 533 Everett’s; each of them, however, is mapped to a dif-  
 534 ferent location in the semantic model with a different  
 535 probability. This treatment removes the class summa-  
 536 tion in the above and we have the final expression of

the environment-specific name probability of: 537

$$\begin{aligned} P(N_i | E) &= \frac{P(N_i | U, [C_n]) P(N_i | L, [C_n]) P(N_i | H, [C_n])}{P^2(N_i | [C_n])} \\ &\quad \times P([C_n] | H), \end{aligned}$$

where  $N_i \in C_n$  and where all the probabilities at the  
 right hand side of the equation have been made avail-  
 able using the several heuristics described above. 539

In the previous discussion, we normalize probabili-  
 ties for each individual conditional probability in the  
 above equations. However, the normalization can be  
 done at the last step. We also noted that the system  
 is not sensitive to small changes of the probabilities.  
 With this in mind, in the MapPointS implementation,  
 we only updated the probabilities when the probability  
 change becomes large. For example, when the current  
 location is 10 miles away to the previous location, or  
 there are 20 new queries in the history. For the same rea-  
 son, the decaying parameters and forgetting parameters  
 are determined heuristically based on the observations  
 from the IIS log. 553

Another important issue in the MapPointS system’s  
 LM computation is smoothing of the probabilities since  
 the training data is sparse. In the current system im-  
 plementation, the probabilities are simply backed up  
 to the uniform distribution when no sufficient amounts  
 of training data are available. 559

With all the above environment or user-specific  
 LM implementation techniques provided by the  
 561  
 562  
 563

Environment Manager component in the MapPointS  
 system, most ambiguities encountered by the sys-  
 tem can be resolved. For example, when a user asks:  
 “Where is Everett”, the system will infer the most prob-  
 able Everett based on the different LM probabilities for  
 the different Everett’s. In most cases, the most probable  
 Everett is either the closest Everett or the frequently  
 visited Everett. In case the system’s guess is incorrect,  
 the user can submit a new query which contains more  
 detailed information in the query. For example, he/she  
 can say “Where is Everett, Washington”. 575

Deng and Yu

Table 4. Four conditions under which the LM of the MapPointS system is constructed and the LM perplexity associated with each condition.

Conditions	LM perplexity
Uniform probability for all city/place names	5748528
Two-level structure for cities and places, but using uniform probabilities for city names	98810
Same as above but using prior probabilities of city names	5426
Same as above but including user-specific information	241

576 Further, in addition to providing useful environmen-  
577 tal or user information to infer the probabilities of  
578 queries in LM, the Environment Manager component  
579 of MapPointS also permits the inference of missing ele-  
580 ments in users' queries to obtain the complete discourse  
581 semantic information. This aspect has been discussed  
582 in [17] in detail and will not be described here.

583 We now present some quantitative results to show  
584 how the user modeling strategy discussed so far in this  
585 section has contributed to the drastic improvement of  
586 the LM. In Table 4, we list the perplexity numbers of  
587 the LM with and without the use of the user-specific  
588 information. These perplexity numbers are based on  
589 four ways of constructing the MapPointS system with  
590 and without using the probabilities and using user  
591 modeling. A lower perplexity of the LM indicates  
592 a higher quality of the LM, which leads to a lower  
593 ambiguity and higher accuracy for speech recognition.  
594 We observe from here that the system utilizing  
595 the user-specific information gives a much lower  
596 perplexity and better LM quality than that otherwise.

#### 597 4. Software Engineering Considerations 598 in MapPoints System Design

599 MapPointS involves its input from multiple modalities,  
600 its output in map presentation, and a large set of data  
601 for training the various system components we have  
602 just described. Without carefully architecting the sys-  
603 tem, the application would be inefficient and difficult  
604 to develop. In designing the MapPointS system, we  
605 have followed several design principles and software  
606 engineering considerations. In this section, we briefly  
607 describe these principles and considerations.

608 The first principle and consideration is *separation*  
609 *of interface and implementation*. Following this princi-

ple, we isolated components by hiding implementation 610  
611 details. Different components interact with each other  
612 through interfaces that have been well defined in ad-  
613 vance. This allowed us to develop and test the system  
614 by refining components one by one. It also allowed us  
615 to hook MapPointS to different ASR engines without  
616 substantially changing the system.

The second principle and consideration is *separa-* 617  
618 *tion of data and code*. MapPointS can be considered as  
619 a system whose design is driven by data and grammar.  
620 In the system design, we separated data from code and  
621 stored the data in the file system. The size of the data  
622 stored is huge since we need to maintain all the city  
623 names, place names, and their associated prior proba-  
624 bilities. By isolating the data from the code, we freely  
625 converted the system from one language to another by  
626 a mere change of the grammar, the place names, and  
627 the ASR engine for a new language.

The third principle and consideration is *separation* 628  
629 *of modalities*. We separated modalities of the speech  
630 input, text input, and the mouse input by representing  
631 their underlying semantic information in a common  
632 SML format. This allowed us to debug modalities one  
633 by one, and also allowed us to integrate more modal-  
634 ities in the future for possible system expansion by  
635 simply hooking the existing system to a new semantic  
636 parser.

The fourth principle and consideration is *full ex-* 637  
638 *ploitation of detailed user feedback*. MapPointS pro-  
639 vides detailed feedback to users in all steps that are  
640 carried out in processing the users' requests. In doing  
641 so, the users become able to know whether the sys-  
642 tem is listening to them and whether the ASR engine  
643 recognizes their requests correctly.

The final principle and consideration is *efficient de-* 644  
645 *sign of the application grammar*. One of the signif-  
646 icant problems of a large system like MapPointS is  
647 the creation of the specific application grammar, or  
648 grammar authoring. A good structured grammar can  
649 significantly reduce the effort in interpreting the re-  
650 sults of speech recognition. In our implementation, we  
651 organized the grammar so that the semantic representa-  
652 tion of the speech recognition results can be interpreted  
653 recursively.

#### 5. Robust Processing Techniques 654 for Speech-Centric HCI Systems 655

Robustness to acoustic environment, which allows 656  
657 speech recognition to achieve immunity to noise and

## Speech-Centric Perspective for Human-Computer Interface

658 channel distortion, is one key aspect of any speech-  
 659 centric HCI system design considerations. For exam-  
 660 ple, for the MiPad and MapPointS systems to be ac-  
 661 ceptable to the general public, it is desirable to remove  
 662 the need for a close-talking microphone in capturing  
 663 speech. The potential mobile application of MapPointS  
 664 for navigation while traveling presents an even greater  
 665 challenge to noise robustness. Although close-talking  
 666 microphones pick up relatively little background noise  
 667 and allow speech recognizers to achieve high accuracy  
 668 for the MiPad-domain or MapPointS-domain tasks, it  
 669 is found that users much prefer built-in microphones  
 670 even if there is minor accuracy degradation. With the  
 671 convenience of using built-in microphones, noise ro-  
 672 bustness becomes a key challenge to maintaining de-  
 673 sirable speech recognition and understanding perfor-  
 674 mance. Our recent work on speech processing aspects  
 675 of speech-centric HCI systems has focused on this  
 676 noise-robustness challenge in the framework of dis-  
 677 tributed speech recognition (DSR).

678 There has recently been a great deal of interest  
 679 in standardizing DSR applications for a plain phone,  
 680 PDA, or a smart phone where speech recognition is  
 681 carried out at a remote server. To overcome bandwidth  
 682 and infrastructure cost limitations, one possibility is  
 683 to use a standard codec on the device to transmit the  
 684 speech to the server where it is subsequently decom-  
 685 pressed and recognized. However, since speech recog-  
 686 nizers only need some features of the speech signal  
 687 (e.g., Mel-cepstrum), the bandwidth can be further  
 688 saved by transmitting only these features. Our recent  
 689 work on noise robustness has been concentrated on the  
 690 Aurora2 and 3 tasks [8, 15], an effort to standardize  
 691 a DSR front-end that addresses the issues surrounding  
 692 robustness to noise.

693 In DSR applications, it is easier to update software  
 694 on the server because one cannot assume that the client  
 695 is always running the latest version of the algorithm.  
 696 With this consideration in mind, while designing noise-  
 697 robust algorithms, we strive to make the algorithms  
 698 front-end agnostic. That is, the algorithms should make  
 699 no assumptions on the structure and processing of the  
 700 front end and merely try to undo whatever acoustic  
 701 corruption that has been shown during training. This  
 702 consideration also favors noise-robust approaches in  
 703 the feature rather than in the model domain.

704 We have developed several high-performance  
 705 speech feature enhancement algorithms on the Au-  
 706 rora2 and 3 tasks and on other Microsoft internal tasks  
 707 with much larger vocabularies. One most effective

algorithm is called SPLICE, short for Stereo-based 708  
 Piecewise Linear Compensation for Environments 709  
 [3–5]. In a DSR system, the SPLICE may be applied 710  
 either within the front end on the client device, or on 711  
 the server, or on both with collaboration. Certainly a 712  
 server side implementation has some advantages as 713  
 computational complexity and memory requirements 714  
 become less of an issue and continuing improvements 715  
 can be made to benefit even devices already deployed 716  
 in the field. Another useful property of SPLICE in 717  
 the server implementation is that new noise conditions 718  
 can be added as they are identified by the server. This 719  
 can make SPLICE quickly adapt to any new acoustic 720  
 environment with minimum additional resource. 721

## 6. Summary and Discussion 722

Recent progress in signal processing and speech recog- 723  
 nition technologies has created a promising direction 724  
 for speech-centric multimodal HCI research. These 725  
 HCI modalities include speech, vision (e.g., gesture), 726  
 pen, mouse, keyboard, screen display, and other GUI 727  
 elements. The speech-centric perspective for HCI ad- 728  
 vocated in this paper is based on the recognition that 729  
 speech is a necessary modality to enable a pervasive 730  
 and consistent user interaction with computers across 731  
 a full range of devices—large or small, fixed or mo- 732  
 bile, and that speech has the potential to provide a 733  
 natural user interaction model. However, the ambigu- 734  
 ity of spoken language, the memory burden of using 735  
 speech as output modality on the user, and the lim- 736  
 itations of current speech technology have prevented 737  
 speech from becoming the choice of mainstream inter- 738  
 face. Multimodality is capable of dramatically enhanc- 739  
 ing the usability of speech interface because GUI and 740  
 speech have complementary strengths. Multimodal ac- 741  
 cess will enable users to interact with an application in 742  
 a variety of ways—including input with speech, key- 743  
 board, mouse and/or pen, and output with graphical 744  
 display, plain text, motion video, audio, and/or synthe- 745  
 sized speech. 746

Two prototype systems, MiPad and MapPointS, de- 747  
 veloped at Microsoft Research take the speech-centric 748  
 perspective in their design. They fully exploit the effi- 749  
 ciency of the speech input, while using other modalities 750  
 to enhance the interaction and to overcome imperfec- 751  
 tion of the speech recognition technology. This paper 752  
 provides a detailed account for the design of the Map- 753  
 PointS system. The system adds the “Speech” modality 754

Deng and Yu

755 into the existing Microsoft product of MapPoint, which  
 756 provides a comprehensive location-based database  
 757 such as maps, routes, driving directions, and proxim-  
 758 ity searches. MapPoint also provides an extensive set  
 759 of mapping-related content, such as business listings,  
 760 points-of-interest, and other types of data that can be  
 761 used within applications. In particular, it is equipped  
 762 with highly accurate address finding and geo-coding  
 763 capabilities in North America, and contains finely  
 764 tuned driving direction algorithms using blended in-  
 765 formation from best-in-class data sources covering 6.7  
 766 million miles of roads in the United States. Loaded with  
 767 the speech functionality, the value of MapPointS to the  
 768 users is the quick, convenient, and accurate location-  
 769 based information when they plan a long-distance trip,  
 770 want to find their way around an unfamiliar town or try  
 771 to find the closest post office, bank, gas station, or ATM  
 772 in any town in North America. The MapPointS system  
 773 has implemented a subset of the desired functionalities  
 774 provided by MapPoint, limited mainly by the complex-  
 775 ity of the grammar (used for semantic parsing),  
 776 which defines what kind of queries the users can make  
 777 verbally, possibly in conjunction with the other input  
 778 modalities such as the mouse click and keyboard input.

779 We in this paper provided an overview of the Map-  
 780 PointS system architecture and its major functional  
 781 components. We also presented several key software  
 782 design engineering principles and considerations in de-  
 783 veloping MapPointS. One useful lesson we learned in  
 784 developing MapPointS is the importance of user or  
 785 environmental modeling, where the user-specific in-  
 786 formation and the user's interaction history with the  
 787 system are exploited to beneficially adapt the LM. The  
 788 drastically reduced perplexity of the LM not only im-  
 789 proves speech recognition performance, but more sig-  
 790 nificantly enhances semantic parsing (understanding)  
 791 which acts on all types of input modalities, speech or  
 792 otherwise. Some quantitative results we presented in  
 793 Table 4 substantiated this conclusion.

794 Our current work is to apply the lessons learned  
 795 from the MapPointS case study, user modeling in  
 796 particular, as presented in detail in this paper to  
 797 other speech-centric HCI tasks. For the extension of  
 798 the prototype MapPointS system, we perceive the  
 799 following future work:

- 800 • Port the system into mobile devices such as Pocket  
801 PC.
- 802 • Incorporate GPS information into the existing Map-  
803 PointS functionality.

- Include new system functionalities such as direct 804  
address searching through speech. 805
- Improve the dialog system component in order to 806  
provide the speech response (instead of only the 807  
GUI response as is now), and to resolve confusability 808  
using speech interaction. 809

## Acknowledgments 810

This invited paper is an expanded version of the 811  
 keynote speech given by the first author at the IEEE 812  
 Fifth Workshop on Multimedia Signal Processing in 813  
 St. Thomas, US Virgin Islands on December 11, 2002. 814  
 The authors wish to thank the organizing committee of 815  
 the 2002 IEEE Workshop on Multimedia Signal Pro- 816  
 cessing for the invitation to present the content of this 817  
 paper at the workshop, and for the invitation to write 818  
 this full-length paper. They also thank the members of 819  
 the entire speech technology group at Microsoft Re- 820  
 search in Redmond, who contributed many different 821  
 components of the system described in this paper. 822

## References 823

1. L. Comerford, D. Frank, P. Gopalakrishnan, R. Gopinath, and 824  
J. Sediwy, "The IBM Personal Speech Assistant," in *Proc. Int.* 825  
*Conf. on Acoustics, Speech and Signal Processing*, vol. 1, Salt 826  
Lake City, UT, May 2001. 827
2. L. Deng, K. Wang, A. Acero, H. Hon, J. Droppo, C. Boulis, 828  
Y. Wang, D. Jacoby, M. Mahajan, C. Chelba, and X.D. 829  
Huang, "Distributed Speech Processing in MiPad's Multimodal 830  
User Interface," *IEEE Transactions on Speech and Audio* 831  
*Processing*, vol. 10, no. 8, 2002, pp. 605–619. 832
3. L. Deng, A. Acero, M. Plumpe, and X.D. Huang, "Large- 833  
Vocabulary Speech Recognition Under Adverse Acoustic 834  
Environments," in *Proc-ICSLP-2000*, Beijing, China, Oct. 835  
2000, vol. 3, p. 806–809. 836
4. L. Deng, A. Acero, L. Jiang, J. Droppo, and XD Huang. 837  
"High-Performance Robust Speech Recognition Using Stereo 838  
Training Data," in *Proc. ICASSP-2000*, vol. 1, Salt Lake City, 839  
Utah, April 2001, pp. 301–304. 840
5. J. Droppo, L. Deng, and A. Acero, "Evaluation of SPLICE on 841  
the Aurora2 and 3 Databases," in *Proc. ICSLP-2002*, Denver, 842  
CO., 2002. 843
6. S. Dusan, G. Gadbois, and J. Flanagan, "Multimodal Interac- 844  
tion on PDA's Integrating Speech and Pen Inputs," in *Proc.* 845  
*Eurospeech-2003*, Geneva, Switzerland, Sept. 2003. 846
7. R. Hamburgden, D. Wallach, M. Viredaz, L. Brakmo, C. 847  
Waldspurger, J. Bartlett, T. Mann, and K. Farkas, "Itsy: Stretch 848  
the Bounds of Mobile Computing," in *IEEE Computer*, April 849  
2001, pp. 28–36. 850
8. H.G. Hirsch and D. Pearce, "The AURORA Experimental 851  
Framework for the Performance Evaluations of Speech 852  
Recognition Systems Under Noisy Conditions," in *Proc. ISCA* 853

Speech-Centric Perspective for Human-Computer Interface

**854** ITRW-ASR2000 Automatic Speech Recognition: Challenges for  
**855** the Next Millennium, Paris, France, Sept. 2000.  
**856** 9. M. Johnston et al., "MATCH: An architecture for Multimodal  
**857** Dialogue Systems," in *Proc. ACL-2002*, Philadelphia, PA, 2002.  
**858** 10. J. Lai, (Ed.), "Conversational Interfaces: Special section,"  
**859** Communications of the ACM, vol. 43, no. 9, 2000.  
**860** 11. A. Nefian, L. Liang, X. Pi, L. Xiaoxiang, C. Mao, and K.  
**861** Murphy, "Dynamic Bayesian Networks for Audio-Visual  
**862** Speech Recognition," EURASIP Journal of Applied Signal  
**863** Processing, vol. 11, 2002, pp. 1–15.  
**864** 12. C. Neti, G. Iyengar, G. Potamianos, A. Senior, and B. Maison.  
**865** "Perceptual Interfaces for Information Interaction: Joint Pro-  
**866** cessing of Audio and Visual Information for Human-Computer  
**867** Interaction," in *Proc. ICSLP-2000*, Beijing, vol. 1, 2000,  
**868** pp. 11–14.  
**869** 13. S. Oviatt, "Breaking the Robustness Barrier: Recent Progress on  
**870** the Design of Robust Multimodal Systems," in M. Zelkowitz,  
 Advances in Computers (Ed.), Academic Press, vol. 56, 2002, **871**  
 pp. 305–341. **872**  
 14. R. Rose, S. Parthasarathy, B. Gajic, A. Rosenberg, and S. **873**  
 Narayanan, "On the Implementation of ASR Algorithm for **874**  
 Hand-held Wireless Mobile Devices," in *Proc. ICASSP-2001*, **875**  
 Salt Lake City, Utah, vol. I, 2001. **876**  
 15. J. Segura, A. Torre, M. Benitez, and A. Peinado, "Model-Based **877**  
 Compensation of the Additive Noise for Continuous Speech **878**  
 Recognition: Experiments using the AURORA2 Database **879**  
 and Tasks," *Proc. Eurospeech-2001*, Aalborg, Denmark, **880**  
 2001. **881**  
 16. R. Shama, V. Pavlovic, and T. Huang, "Toward Multimodal **882**  
 Human-Computer Interface," in *Proc. IEEE*, vol. 86, no. 5, **883**  
 1998, pp. 853–869. **884**  
 17. K. Wang, "Implementation of a Multimodal Dialog System **885**  
 Using Extended Markup Language," in *Proc. ICSLP-2000*, **886**  
 Beijing, China, 2000. **887**



Au: Pls.  
 provide  
 Bios and  
 Photo.

UNCORRECTED PROOF