# Creating conversational interfaces for children

Shrikanth Narayanan[1], *Member, IEEE*, and Alexandros Potamianos[2], *Member, IEEE*

[1]AT&T Labs–Research, 180 Park Ave, P.O. Box 971, Florham Park, NJ 07932-0971, U.S.A.

[2]600 Mountain Ave., Murray Hill, NJ 07974-0636, U.S.A.

`shri@research.att.com`, `potam@research.bell-labs.com`

### Abstract

Creating conversational interfaces for children is challenging in several respects. These include acoustic modeling for automatic speech recognition (ASR), language and dialog modeling, and multimodal-multimedia user interface design. First, issues in ASR of children speech are introduced by an analysis of developmental changes in the spectral and temporal characteristics of the speech signal using data obtained from 456 children, ages 5-18 years. Acoustic modeling adaptation and vocal tract normalization algorithms that yielded state-of-the-art ASR performance on children speech are described. Second, an experiment designed to better understand how children interact with machines using spoken language is described. Realistic conversational multimedia interaction data were obtained from 160 children who played a voice-activated computer game in a Wizard of Oz (WoZ) scenario. Results of using these data in developing novel language and dialog models as well as in a unified maximum likelihood framework for acoustic decoding in ASR and semantic classification for spoken language understanding are described. Leveraging the lessons learned from the WoZ study and a concurrent user experience evaluation, a multimedia personal agent prototype for children was designed. Details of the architecture and application details are described. Informal evaluation by children was very positive especially for the animated agent and the speech interface.

### Keywords

Automatic Speech Recognition for Children, Spoken Dialogue Systems, Multimodal Systems, Multimedia Interfaces, Human Computer Interactions, User Experience, Computer Games, Personal Agents.

## Introduction

Recent advances in speech and multimedia technology have spurred worldwide deployment of several prototype and commercial applications that provide natural spoken dialogue with machines [33], [8], [11], [4]. Such efforts are, however, primarily targeted toward the adult user population.

While the state of the art in speech technology is still not perfect for the adult population, the task of building spoken dialogue applications for children poses even greater challenges. Children speak and interact with computers differently from adults. Several aspects of these differences can be identified such as in the acoustic and linguistic characteristics of speech, dialog interaction strategies, problem solving skills and user preferences. Further, the sources of these differences reflect physiological and anatomical changes associated with development of articulators and the effects of socio-economic factors during a child's growth. Some of these research challenges will be considered in this paper.

The CHildren's Interactive Multimedia Project (acronym: CHIMP) at AT&T aimed at providing essential guidelines for engineering successful multimodal-input multimedia-output applications for children with an emphasis on the spoken dialog interface. Factors that motivated this study include: (1) Children form a crucial segment of customer population for interactive multimedia systems, and (2) Children are eager and quick to embrace, and use, new technologies. There are several statistics supporting these two facts. More than 60% of children (4-11 years) use a PC at home compared to 40% for the total U.S. population [15]. Games were found to take up 40% of a preteen's time on the computer, out of a total 6.5-12.5 hours/week, while the rest of the time was devoted to school work [1]. Another study found that out of a total population of about 44 million 2-12 year old children in the U.S., about 25% were projected to be online by the year 2000; the projection increased to 45% for the year 2002 [9]. About 67% of children reported using the internet to gather information, 65% to play games, 49% to do chats, 48% to do creative activities, and 46% to download 'stuff' [17]. Similar statistics are available for teenagers. Over 90% of teenagers, across economic boundaries, use computers and more than 70% use the internet [16]. Over 98% of them credit technology for making a positive difference in their lives and 92% believe technology will improve education and job opportunities. Over 71% of them want to *talk* to their computers: speech recognition was the number one high-tech product kids would like to see developed. A third motivating factor for this study was: (3) The lack of speech-technology resources for creating voice-enabled applications for children. The resources designed for the adult population are not directly usable for children users. For instance, the age-dependent acoustic and linguistic variability in children's speech [13] makes automatic speech recognition (ASR) for children more difficult compared to adults, hence requiring special algorithms to be designed for providing satisfactory levels of ASR performance. For example,

while analyzing ASR performance of the live usage dataobtained from their Jupiter system, Zue et al [33, Fig. 9] found that the in-vocabulary word error rate for children was almost twice that for adult users.

The idea of providing voice interfaces for children's applications is not a new one, however the scope of the systems that have been developed thus far has been relatively limited. Examples of spoken dialog system prototypes for children include word games for pre-schoolers [30], aids for reading [14] and pronunciation tutoring [28]. There is also an increasing number of commercial products being brought to the market – toys and computer games for children – that have limited speech recognition capabilities (small vocabulary, typically isolated or key word recognition). But due to the inherent poor automatic speech recognition and understanding performance, speech has rarely been used as the primary interaction modality in these applications.

Recently, there has also been increasing interest in the design of multimodal interfaces that combine speech with a variety of other input modalities such as text, touch, mouse clicks, hand-writing, and gestures [29], [31], [6]. Results of these investigations suggest that the use of multiple modalities, rather than a single modality, leads to more efficient and natural interaction and en-hances the overall user experience (for example, [3]). Multimodality is attractive in the creation of conversational interfaces for children in the sense of both overcoming inherent limitations in speech technology and exploiting the ubiquitous availability and/or familiarity with conventional modalities such as the computer mouse, keyboard, joy stick and pen. There are several open research issues, that need to be addressed, including multimodal input integration and interpre-tation, multimodal dialog design, multimedia output presentation and performance evaluation. Realistic case studies and prototype designs are crucial to further our understanding of multi-modal interactions. The design of a multimodal prototype application for children that will be described in this paper represents an effort in this direction.

Building conversational interfaces for children is a challenging problem and needs to be carried out in several stages. The first step is to establish a proof of concept for the use of speech as a viable means for children to interact with a machine, both in terms of feasibility and usability. Second, data from children need to be collected for quantifying the variability present in their speech and to train and test models for automatic speech recognition (ASR) and spoken language understanding (SLU). This is necessary for ensuring *acceptable* levels of ASR and SLU performance across all ages and environments. Finally, the intuition and results obtained from

such data analyses and modeling [21] can be used to create prototype systems.

The rest of the paper is organized as follows. In Section I, acoustic variability issues in children speech and the problem of ASR for children are addressed. The collection, analysis, and modeling of acoustic, spoken language and dialog data obtained from a Wizard of Oz (WoZ) experiment are described in Section II. Results from a user experience evaluation that investigated subjective impressions of children regarding various aspects of the conversational interface are also presented. In Section III, a design of a conversational multimodal system that leverages the results of the WoZ experiments is presented. The "*Agent CHIMP*" prototype combines speech, keyboard and mouse input modalities and uses text, graphics, speech and animation for output presentation. The application is controlled by animated agents. An informal user evaluation of the prototype and future directions are provided.

## I. ASR FOR CHILDREN

Many present day ASR systems, including the ones considered in this paper, use a hidden Markov model (HMM) based pattern recognition wherein statistical models constructed from speech-data samples (training) are used to discern patterns in new unseen speech samples (testing) [26]. Acoustic variability in children's speech, which renders pattern classification difficult, is identified as a major hurdle in building high performance ASR applications for children (Sec. I-A). In Section I-B the effect of age on the ASR performance of children speech is described. This is followed by a description of how speaker normalization was used to reduce variability and increase the resolution between pattern (phone) classes: A speaker normalization procedure that combines spectral shaping and frequency warping [25] was implemented that resulted in recognition error rate reduction of up to 45%.

### A. Acoustic characteristics of children's speech

Investigations of children speech have shown systematic age-dependent variation in the acoustic correlates of speech such as formants, pitch and duration [5], [10], [7], [13]. As a part of the AT&T CHIMP project, changes in the temporal and spectral parameters of children's speech were investigated using speech data (23454 utterances) obtained from 436 children ages between 5 and 18 years and 56 adults [13]. Results showed a systematic decrease in the values of the mean and variance of the acoustic correlates such as formants, pitch and duration with age, reaching adult ranges around 13 or 14 years. The intra-speaker variability was larger for young children,

especially for those under 10 years (Fig. 2). Formant frequency values exhibited linear scaling with age, especially for males (Fig. 1).

There are several implications of these observed age-dependent trends on the ASR of children's speech. The increased spectral and temporal variability in formant values results in greater overlap among phonemic classes for children than for adult speakers, thus rendering the pattern classification problem inherently more difficult. Further, the range of values for most acoustic parameters is much larger for children than for adults. For example, five-year old children have formant values up to 50% higher than male adults [13]. The combination of a large acoustic parameter range and increased acoustic variability seriously degrades ASR performance, as shown in the next section.

Additionally, there are some fundamental issues in processing children's speech. Spectral feature extraction, the typical front-end signal processing step in ASR, is more difficult for children's speech because the fundamental frequency and the formant bandwidths are of comparable magnitude. Moreover, for a given signal bandwidth (say a 4 kHz telephony band), there are fewer formants in the spectra of children's speech compared to those of adults. Thus, the sparse sampling of the spectrum (due to high F0 values) and relatively fewer formants in the given bandwidth (due to high formant values) in children's speech pose fundamental limitations on the amount of phoneme-dependent information available at the ASR front-end.

*B. Baseline ASR performance: Children vs. Adults*

Baseline ASR performance was evaluated as a function of speaker's age for two tasks: (1) Connected digit recognition and (2) Command and control phrase recognition. The acoustic models for these experiments were trained from speech utterances collected over the public switched telephone network from both adult and children speakers. Details of the training and testing databases are provided in Table I. A mixture of 6 Gaussians was used to model each state of the context-dependent digit units. Separate phone HMMs for adult and children speakers were trained from the corpora DgtI, DgtII, SubwI and SubwII ("CHLD"), respectively (see Table I for corpus details). A mixture of 16 Gaussians was used to model each state of the the 40 context-independent (subword) English phone units.

In Fig. 3(a), word recognition accuracy for a connected digit recognition task (corpus DgtTest) is plotted as a function of age for two model training conditions: models trained from adult speakers (corpus DgtI), labeled "Adult HMM", and from children speakers (corpus DgtII), la-

beled "Child HMM". For both matched and (especially for) mismatched training and testing conditions, the recognition performance decreases substantially for young children. Performance reaches adult levels at approximately thirteen or fourteen years of age. This agrees with the observation in [13] that by the age of fourteen *both* the mean and standard deviation of most acoustic parameters reach adult levels.

Overall recognition performance for children speakers was *up to four times worse* than for adults depending on the speaker's age. For mismatched training and testing conditions ("Adult HMM"), word error rate is approximately two to three times higher than for matched conditions. The major reasons for performance degradation in younger speakers are acoustic mismatch between the training and testing data, increased acoustic variability and the large range of acoustic parameters. Speaker normalization and model adaptation were used to reduce the mismatch and variability. These procedures are summarized in the following section.

## C. Linear Frequency Warping and Model Adaptation

The frequency warping approach to speaker normalization compensates mostly for inter-speaker vocal tract length variability by linear warping of the frequency axis by a factor $\alpha$ [12]. Frequency warping is implemented in the mel-frequency filterbank front-end by linear scaling of the spacing and bandwidth of the filters. For each utterance, the optimal warping factor $\hat{\alpha}$ is selected from a discrete ensemble of possible values so that the likelihood of the warped utterance is maximized with respect to a given HMM and a given transcription. Let $X^\alpha$ denote the sequence of cepstrum observation vectors warped by a linear frequency warping function. If $\lambda$ denotes the parameters of the HMM model, then the optimal warping factor is defined as

$$\hat{\alpha} = \arg \max_\alpha P(X^\alpha | \alpha, \lambda, H) \tag{1}$$

where $H$ is a decoded string obtained from an initial recognition pass. The selected observation vector sequence $X^{\hat{\alpha}}$ is decoded in a second recognition pass to obtain the recognized string.

There is a large class of maximum likelihood based model adaptation procedures that can be described as parametric transformations of the HMM model or the observation sequence. For these procedures, we let $\lambda_\gamma = h_\gamma(\lambda)$ denote the model obtained by a parametric linear transformation $h_\gamma()$. The optimal parameters of the linear transformation $\hat{\gamma}$ and the frequency warping $\hat{\alpha}$ can be simultaneously estimated. The maximum likelihood criterion can be used to select the appropriate model and also optimize the parameters of the speaker normalization and

model adaptation algorithms as follows

$$\{\hat{\alpha}, \hat{\gamma}\} = \arg\max_{\{\alpha, \gamma\}} P(X^\alpha | \alpha, \gamma, H) \ . \tag{2}$$

The potential of this class of procedures was investigated in the context of speaker adaptation from single utterances. In our case, $h_\gamma()$ is a simple linear bias applied to the means of the model distributions or the observation sequence [25], and $\lambda^n$, $n = 1, .., N$ is a family of age-group dependent acoustic models. The results of speaker normalization and model adaptation applied to the connected digits and command and control recognition tasks are described next.

**Experimental Results:**

In Fig. 3(a), digit recognition accuracies before and after speaker normalization are shown (test corpus DgtTest) for HMMs trained from both adult (DgtI corpus) and children (DgtII corpus) speaker populations. The allowed range of formant frequency scaling was from –20% to +12% and a total of 17 warping factors were examined during frequency warping. The error rate reduction due to speaker normalization was up to 50%, and was greater for young speakers under twelve years of age and when mismatched models trained from adult speakers ("Adult HMM") were used (dotted vs. dashed line in Fig. 3(a)). After speaker normalization, the recognition accuracy for children speakers over 9 years of age was comparable to that of adults. The summary of the cumulative results for all ages is given in Table II. In addition, the performance of an HMM trained from data (equally) mixed from the adult and children corpora DgtI and DgtII is shown (labeled "Cld+Adlt HMM"). Overall, digit error rate reduction by 25-45% was achieved using speaker normalization.

In Fig. 3(b), (c), word recognition accuracy is shown as a function of age for the command and control task are shown under similar training and testing conditions. CommI and CommII consist of 10 possible phrases (16 word vocabulary) and 50 phrases (68 word vocabulary), respectively. Similar to the digit recognition task, speaker normalization helped significantly to bridge the gap in performance between the models trained from adult and from children speaker populations. However, recognition accuracy levels obtained for adult speakers were still not obtainable for the younger age group (6-9 years), suggesting that normalization strategies more sophisticated than simple linear frequency warping may be needed.

## II. Creating conversational systems for children: WoZ experiments

To investigate how children converse with interactive systems and to collect speech, dialog and user experience data in a realistic spoken language application environment, a Wizard of Oz (WoZ) experiment was designed. Increased acoustic and linguistic variability are typical of spontaneous speech, and the WoZ experiment was aimed to provide valuable data toward evaluating ASR and SLU performance of spontaneous child-machine interactions in realistic scenarios. Note that the ASR performance described in Section I was based on read speech obtained in a relatively controlled set up.

About 160 children, ages 8-14 years, participated in the study by playing an interactive computer game using voice commands, or keyboard and mouse control [21]. The software selected for this WoZ experiment was the popular computer game "Where in the U.S.A. is Carmen Sandiego?" (WITUICS) by Brøderbund Software (now, The Learning Company). WITUICS is an interactive detective game for children ages eight years and older. There were several reasons why this computer game was chosen for the study. Overall, the game was rich in dialog subtasks including navigation and multiple queries, database entry, and database search. Further, the fact that (during a substantial part of the game) the child conversed with cartoon characters on the screen made the dialog more natural and human-like. As a result spontaneous speech could be elicited.[1] The structure of the game was not changed (no adaptation to voice inputs). The only modification was the addition of four generic text-to-speech synthesized dialog error control and clarification messages: (1) I can't do that now, what else would you like me to do? (2) Can you spell that for me? (3) What was that? (4) I don't know how to do that. What else would you like me to do?

### A. Game Description

To successfully complete the game, i.e., arrest the appropriate suspect, two subtasks had to be completed, namely, (i) determining the physical characteristics of the suspect and completing a profile sketch to enable an arrest warrant, and (ii) tracking and apprehending the suspect (by traveling through at least five of the fifty states in U.S.A every game). The player could talk to various characters appearing on the game screen seeking clues about the suspect's trail and

---

[1]The children players were not informed of the existence of a wizard and an observation room. Further, for approximately half of the experimental runs the player was alone in the game room without a moderator present.

physical appearance. To help interpret the clues thus obtained, the player can use aids such as geographical databases that can be queried using single or multiple word searches. A game was deemed successful when the player traveled to the correct location, and identified the suspect correctly (using the profile information) from among several cartoon characters on the screen.

## B. Experimental Setup

The Wizard of Oz (WoZ) experimental setup is shown in Fig. 4. The player sat in front of a slave monitor wearing headphones, i.e., watching and listening to the audio-visual output piped from the wizard's computer. In the observation room, the wizard controlled the experiment by providing the appropriate output in response to the user's input. Since the audio-channel of the game was not intercepted, the pre-defined dialog error-control and clarification messages were played through a separate audio channel connected to a loudspeaker placed next to the slave monitor. High-quality audio recordings of the player's voice commands are collected using a close-talking head-mounted microphone (Sennheiser) and a far-field microphone (??). The audio output from the game was also recorded for reference. A video recording of the "picture-in-picture" image of the player and the game screen complete with the (mixed) audio from player and computer was also obtained. Neither the loudspeaker nor the video camera were found to be intrusive by any of our subjects.

## C. Experiments and Population Statistics

Although a variety of experiments were conducted using voice(V), keyboard and mouse(K+M), or voice, keyboard and mouse (V+K+M) inputs to control the game, the primary focus in this paper will be on the voice modality interactions.

Prior to each experiment, the game and the voice interface were explained to the child by a moderator. The players were not informed of the existence of a wizard. The wizard followed a set of pre-defined rules while playing the game while an assistant helped with managing the data collection during the experiment. Data from a total of 160 children and 7 adults were collected. Most players played two games (23% played one game and 3% played three games). The total number of games played (using voice with no recognition errors) per age group and gender are shown in Table III. A total of about 50000 utterances were collected. After the completion of the experiment, the moderator interviewed the subject to gauge the user's perception regarding the game and the interface.

*D. Subjective user evaluation*

All the children who took part in the WoZ experiment participated in an exit interview wherein subjective impressions about the game and the interface were obtained. Sample questions that they were asked included: (1) What did you like about using voice activation? (2) What did you like/dislike about the game? (3) Would you like to use voice with keyboard and mouse? The participants were also asked to rate on a scale from 1 to 5 (5 being the highest) the following: voice interface, game, use of headset, TTS-generated error messages, and the use of multimodal inputs.

The children gave very high ratings to the speech interface (93% rated the interface 4 or 5). The game also received high marks but somewhat lower than the interface (only 81% rated the game 4 or 5). The speech interface ratings degraded only slightly when 5% misrecognitions and 5% rejections were randomly introduced into the game by the wizard. It is interesting to compare the relation between the number of games won and the ratings. Losing a game had a significant negative effect on the rating of the game. However, there was no significant effect of the game outcome on the children's rating of the voice input. The 11-12 year olds gave the highest ratings for the voice input. Gender effects were negligible.

Other results showed that the dislike for TTS generated error messages and for spelling (for the purpose of "ASR ambiguity resolution") decreased with age. The enjoyment of the use of a headset microphone roughly correlated with the enjoyment of the game. Finally, about two-thirds of the children preferred having a multimodal interface to a voice-only interface. In summary, user experience results were promising for the inclusion of voice as one of the interaction modalities in the design of interactive applications for children.

*E. Dialog data analysis*

In this section, analysis of "dialog" i.e., user-system interaction data is presented for the voice and keyboard-mouse modalities. Speech utterances were manually assigned to dialog states according to the game actions they triggered [21]. Dialog states were defined to roughly correspond to one (or a group of similar) actions taken by the wizard in response to a user input. For example, the dialog state "Talk2Him" incorporated user queries asking for a cartoon character's attention, while states "WhereDid" and "TellMeAbout" corresponded to queries about the suspect's whereabouts and physical characteristics, respectively. Spoken utterances were assigned

to predefined dialog states by the wizard's assistant while the game was being played and were later verified by a group of human labelers. A sample interaction illustrating dialog state tagging is given in Table II-E. A total of about thirty dialog states were identified for this application.

The flow through the task was characterized by a sequence of dialog state transitions. The game dialog flow primarily consisted of navigation/query, database search and database entry subdialogs. Fig. 5 shows the dialog flow diagram for the navigation/query subdialog. The total number of times a state is visited (in parenthesis) and the total number of state transitions (arrow labels) are shown for all games played by children players (total of 290 games). Such graphs provided useful information about problem-solving and dialog strategies of children. For example, it can be deduced from Fig. 5 that only 30% of the time a game character was asked to provide information about both the location and the physical characteristics of a suspect, i.e., most children preferred to concentrate on a single task per turn. Similar remarks can be made for the frequency of skipping states e.g., executing a database query (state "Find") without first opening the database (state "Database") in the database search subdialog, or the frequency of superfluous greetings (e.g., "Goodbye").

Age and speaker dependencies in dialog state transitions were analyzed. Key observations regarding spoken interactions included: (1) queries seeking multiple attributes were far less common than those seeking a single attribute (2) frequency of skipping states in the canonical game structure was low (3) frequency of superfluous commands such as "goodbye" was relatively high. There were no noticeable differences in the dialog patterns of male and female children. However, the dialog patterns of older children (11-14 years) were different from those of younger ones (8-10 years). The older children tended to complete the game faster, did fewer database lookups, used more advanced dialog patterns, and had fewer out-of-domain utterances (about half the number as the younger group).

*F. Dialog strategies: Keyboard and Mouse vs. Voice*

A total of 12 children players alternated on using voice and keyboard-mouse (K+M) to control the game. The dialog/action flow and underlying task solving strategies were very similar for both modalities. The total number of commands was roughly the same for the navigation/query and database entry subtasks. However, children took fewer turns (almost 50%) using keyboard and mouse than voice to carry out the relatively high-perplexity database search and retrieval tasks. This suggests that for the database search task voice is not the most efficient modality

(with the current interface). A final observation is that when using K+M superfluous greetings at the navigation/query menu (dialog state: "Goodbye") were reduced by a factor of three compared to using voice. This reinforces the belief that *although voice might not be the most efficient modality, it is a more natural modality.*

## G. Linguistic and Acoustic Analysis

In this section, inter- and intra-speaker linguistic variability for groups of utterances that are semantically equivalent (i.e., those that trigger the same game action) are investigated. Specifically, the average Levenshtein string distance was computed among all strings belonging to the same dialog state and speaker category, and compared with average string distance among all speakers. In addition, the frequency of occurrence of disfluencies and filled pauses were measured for each age group. Finally, average word length of utterances, average utterance duration and speaking rate were measured.

### G.1 Linguistic variability

Linguistic variability for semantically equivalent sentences was measured for "simple" dialog states (corresponding to a single unambiguous game action) in the navigate/query and database entry subtasks. This subset of the data contained 22422 utterances. All sentences collected from speaker $n$ that belonged to the "simple" dialog state $k$ were deemed elements of class $C_{k,n}$. The intra-speaker linguistic variability for dialog state $k$ was then defined as $(1/\sum_n L_{k,n}) \sum_n \sum_{i,j} d(S_i, S_j)$, where $d$ is the Levenshtein word-string distance (with 0.75 penalty for word insertion/deletions and 1 for substitutions), $S_i, S_j \in C_{k,n}$, and $L_{k,n}$ is the total number of words in $C_{k,n} \times C_{k,n}$. Similarly, inter-speaker linguistic variability was defined as $(1/L_k) \sum_{i,j} d(S_i, S_j)$, where $S_i, S_j \in C_k = \bigcup_n C_{k,n}$. Table V shows the linguistic variability for various dialog states. Overall, *inter-speaker variability is almost twice as high as intra-speaker variability.* This suggests that there is potential gain from building speaker-specific language models or from performing speaker adaptation on the language models. Note also that both the inter- and intra-speaker variability in Table V varies a lot among dialog states. Finally, intra-speaker linguistic variability was computed for the 8-10 and 11-14 age groups, and between male and female speakers. Overall, female speakers displayed higher intra-speaker variability by about 10% than male speakers but this trend was dialog state-dependent. Similarly, about 10% increase in linguistic variability was found in the 11-14 age group vs. the 8-10 age group.

G.2 Extraneous-speech modeling

Any speech utterance that triggered no valid game response or action was defined to be ex-
traneous i.e., out of domain. Statistical modeling of (sequences of) dialog states that precede
extraneous speech events is important for designing robust dialog systems for children. In the
WITUICS data, extraneous speech utterances corresponded to approximately 5% of all utter-
ances spoken for the 8-10 year-olds (3.7% for all subjects), ranging from 0% to 25% among
individual subjects (7% variance). Most extraneous speech utterances fell in one of the following
categories: (i) those expressing excitement (disappointment) when vital (useless) information
is provided by the game or success (failure) was achieved in one of the game stages, (ii) those
requesting game-strategy information, interpretation of game output or approval by other people
in the room (an adult moderator or other children were present in the game room for about half
of games played), and (iii) interacting with characters on the screen irrelevant to game goals and
objectives. Overall, the extraneous speech utterances were found to be highly speaker-dependent,
age-dependent, and to be preceded by a small subset of dialog states. Results imply that mod-
eling extraneous-speech at the dialog level can significantly contribute to successful utterance
verification strategies.

Disfluencies and hesitations in the speech data were analyzed as a function of age and gender.
Mispronounciations, false-starts, (excessive) breath noise and filled pauses (e.g., um, uh) were
manually labeled for a subset of the data (22422 utterances). About 2% of the labeled utterances
contained false-starts and 2% contained (obvious) mispronounciations. Breathing and filled
pauses were found in 4% and 8% of the utterances, respectively. While no gender dependency
was found for any of the disfluency measures, there was a distinct age dependency. The frequency
of of mispronounciations was almost twice as high for the younger (8-10 years) age group than for
the older group (11-14 years). Breathing noises occurred 60% more often for younger children.
Surprisingly, this trend was reversed for filled pauses which occurred almost twice as often for
the 11-14 age group. Although disfluencies and hesitation phenomena occur more frequently in
children than in adults, our experience showed that ASR performance does not suffer significantly
due to these effects, hence requiring no special acoustic modeling strategies.

Finally, small differences in duration and average string length were found between the young
and old age groups. No gender or age bias was found in the average utterance length (in words).
The average sentence duration was about 10% longer for younger children. As a result, the

speaking rate for the 11-14 year-olds was about 10% higher than for the younger group which is in agreement with [13].

*H. ASR and SLU Performance – WITUICS task*

Baseline ASR performance and the effects of speaker normalization and model adaptation for connected digit and command/control phrase recognition task using read speech from children were described in Sec. I-B. In this section, ASR and spoken language understanding (SLU) performance for the conversational WITUICS task will be presented.

For interactive task-oriented applications such as command and control, unlike dictation applications, it is not always necessary to recognize and understand every spoken word. The SLU problem for the WITUICS task was defined as identifying the next dialog state (ACTION classification) and the values of any associated attributes (ATTRIBUTE recognition), given an utterance transcription. For example, SLU for the utterance "I'd like to go to Indiana" will result in (ACTION: travel, ATTRIBUTE: Indiana) while SLU for the utterance "I'd like to travel" will produce (ACTION: travel, ATTRIBUTE: null). In [24], [27], a unified maximum likelihood probabilistic framework for performing both ASR and SLU was proposed and applied to the WITUICS task. The approach and the results are summarized below.

The joint likelihood maximization for acoustic decoding and SLU can be given as

$$\max_{S_t,W_t} P(S_t,W_t|O_t,S_1..S_{t-1}) = \max_{S_t,W_t} P(O_t,|W_t,S_t,S_1..S_{t-1})P(S_t,W_t|S_1..S_{t-1})/P(O_t|S_1..S_{t-1}) \quad (3)$$

$$\equiv \max_{S_t,W_t} P(O_t|W_t)P(W_t|S_1..S_t)P(S_t|S_1..S_{t-1}) \quad (4)$$

where $S_t$ is the dialog state, $W_t$ is the transcribed user input and $O_t$ is the acoustic observation sequence at dialog turn $t$. Since is $S_t$ is not known at decoding, $P(W_t|S_1..S_t)$ is approximated by $P(W_t|S_1..S_{t-1})$. For computation, the problem is decomposed into acoustic-language decoding and understanding from transcription i.e., the posterior probability is maximized first with respect to $W_t$ and then with respect to $S_t$.

$$\hat{W}_t = \arg\max_{W_t} \underbrace{P(O_t|W_t)}_{Acoustic}\underbrace{P(W_t|S_1..S_{t-1})}_{Language} \quad (5)$$

$$\hat{S}_t = \arg\max_{S_t} \underbrace{P(\hat{W}_t|S_1..S_t)}_{Understanding}\underbrace{P(S_t|S_1..S_{t-1})}_{Dialog} \quad (6)$$

Both the language model and the dialog model components in Eqns. (5) and (6) were specified by N-gram automatons.

To evaluate ASR and SLU performance, the WITUICS data was partitioned into a training set with 6039 utterances (59 speakers, 102 interactions) and a test set with 2050 utterances (21 speakers, 37 interactions). The data included orthographic transcriptions of the spoken utterances and manually-assigned dialog state tags. The acoustic models were context-independent phone models, with 16 mixture Gaussians and the language models were dialog state-dependent [27] word trigrams. Bigrams for both the dialog model and the understanding models provided the best SLU classification results. Classification accuracy from true transcriptions (obtained manually) was 94.4% while combined the ASR and SLU yielded 86.3% correct classification. Overall, attribute recognition was comparable to the overall ASR word accuracy level of 78%. In summary, 5%-20% error rate reduction came from language adaptation, and 15%-25% from dialog modeling. It was also shown in [24], that further improvements in SLU can be obtained by utilizing acoustic confidence scores in the understanding model (10%). The results indicated feasibility of building viable spoken dialog systems for children using these ASR and SLU technologies.

## III. Building a Prototype

In this Section, a case study of designing a conversational multimodal prototype application for children is presented. The design exercise leverages the lessons learned from the WoZ study described in Sec. II. A personal communications assistant (providing telephony, web access, email) and a computer game application for children were used as a vehicle to achieve the following goals: (1) Define a general conversational multimodal system architecture. (2) Provide a means for investigating merging of multimodal inputs (keyboard text, mouse clicks, voice) and multimedia presentation strategies. (3) Demonstrate the concept of agent and sub-agent embodiments that handle different modules and functionalities within an application. (4) Demonstrate the role of intelligence and personality of the user interface through spontaneous conversation, audio, animation (gestures), and graphics.

While the user interface design focused on children, the system architecture itself was generic. Since the chosen application for prototyping was different from the one in the WoZ study, language was unavailable. Hence, corpus-driven language and understanding modeling could not be applied at the initial stage of the application creation.

*A. System building blocks*

Figure 6 shows the main functional building blocks of the system from a user's perspective. The central part of the system is the controller. The user interface enables interactions using voice, typed text, mouse clicks or combinations there of. Output to the user is presented through audio, graphics, animation and textual modalities. The speech and language processing unit, comprising the ASR and SLU components, enables spoken language interactions. The dialog manager also communicates with information resources such as databases. Further details of the various modules are given in the following sections.

The prototype system consisted of the following components: input/output (I/O) event handler, dialog manager, graphical user interface (GUI), spoken language understanding (SLU), speech recognizer (ASR), speech synthesizer (TTS), animator and database. The speech recognizer used children-specific acoustic models that were built using the data obtained from the WoZ study described in Sec. II. Since language data was unavailable at the time of the creation of this prototype application domain, application-specific finite state grammars were hand crafted to boot-strap the language models. ASR was performed using the AT&T Watson speech recognition engine.

The dialog manager defines the strategies and actions to be taken based on the user's input and decides what to present to the user. The dialog manager used in the CHIMP prototype was based on AMICA (AT&T's Mixed Initiative Conversational Architecture) which provided a library of dialogue actions and a means for specifying dialog strategies either through a JAVA-based GUI or through a high-level scripting language called DMD [19]. The 'template' is a data structure that AMICA uses to maintain the dialog state information. The dialog manager can communicate with external modules such as ASR and databases through TCP/IP socket connections.

A semantic representation corresponding to the spoken utterance was derived by performing a lexical analysis (CHRONUS tools, [18]) followed by a rule transduction on the resulting lexical lattice. Due to the lack of in-domain data, rules for deriving semantic representation were boot-strapped by hand-crafted rules in a finite state machine representation. The result of the semantic analysis was represented in a template form. The template structure that provides an 'attribute-value' mapping of concepts for use by the dialog manager was derived using a template generator. For simplicity and consistency sake, all actions underlying mouse button clicks and touch on the GUI were also mapped to equivalent, semantically unambiguous, prototype natural

language expressions and are henceforth handled by the semantic analyzer in a mode-independent way similar to spoken or typed inputs.

The GUI consisted of five main areas (Figure 9): graphics/animation area, text area, buttons area, command line area and user command echo area. The GUI design aimed to provide a consistent look and feel across various applications. The personality and appearance of the animated agents provided orientation for the user. Function buttons provided an alternate means of accomplishing several key commands that could also be achieved through voice or typed inputs. A history of user inputs was maintained and any previous input could be easily repeated by highlighting and clicking on the desired entry.

In summary, the input event handler synchronizes the (asynchronous) input from the user (speech or keyboard or mouse events). All inputs are transmitted to the understanding system by the way of the dialog manager which in turn returns a template with the semantic representation. Based on this semantic representation, the dialog manager decides on the next action to take which results in a output template which contains commands for the multimedia presentation. The output template is parsed by the output event handler resulting in a multimedia presentation of the system's response (text, graphics, animation, speech). Further details on the communication between the various modules is given in the next section.

*B. Architecture*

An architectural diagram is shown in Fig. 7, where the focus is on the communication between the controller and various servers. The controller is an asynchronous I/O event handler and comprises two separate event loops for the input and output modalities, respectively. The program flow of the controller event loops is straightforward: (i) input events are sent to the dialog manager in the form of text strings (input loop), (ii) templates received from the dialog manager are parsed for multimodal output fields and passed on to the appropriate output modules (text-to-speech synthesizer, animation/movie player or the graphical user interface). During processing of requests by the dialog manager all input events are queued in a stack. If multiple events have been queued up in the input, only the latest event is sent to the dialog manager for processing. The output event handler has the additional functionality of being able to surrender control or simply start up external applications (e.g., pop up an e-mail reader or a web browser). Finally, the controller handles multimodal barge-in events (speaking- or typing-over voice prompts or animation sequences) from the input clients by informing the text-to-speech and movie player

modules to stop playback of speech and/or animation sequences.

The dialog manager processes incoming strings by following a control language specification. Strings are interpreted by calling the understanding module and domain information is retrieved from a SQL backend database server. State information and state history are encapsulated in a template form. Based on the current input and the dialog manager specification, the state of the dialog manager gets updated and generates appropriate response to the user.

The speech input is processed by the ASR client which communicates with the AT&T ASR WATSON server through a wireline protocol. The commands and audio flows from the client to the server while the results and notifications flow the other way. The interaction between the ASR client and server is through polling or notification modes. Input commands from other modalities are queued in a stack and processed in a similar fashion.

An important principle followed in the design of this prototype is that all the dialog state information is maintained internally in the template that is generated and updated by the dialog manager. In an initial implementation, the state information was transmitted to the I/O handler which in turn decided the form of the multimedia presentation. This strategy was found to be inefficient in the sense that there was duplication of program control logic at the dialog manager and the I/O handler. In a subsequent implementation, the output generation module was folded into the dialog manager: only specific commands for speech synthesis, animation and text display were transmitted to the I/O handler using a predefined protocol. A detailed discussion of the multimodal architecture can be found in [6].

## C. Application details

**Tasks:** The prototype consisted of two distinct tasks as shown in Figure 8: a communications agent application (information retrieval from a personal directory, placing phone-calls, accessing the Internet, sending email) and a computer game (spelling bee). Both applications were controlled by a conversational animated agent. The agent embodiment chosen for the CHIMP prototype was a cartoon chimpanzee character. The personality and appearance of the agent, however, were different across the two applications. The user could freely switch back and forth between the applications at any point during the interaction. The system also has 'go to sleep' and 'wake up' features by means of which the personal agent's attention can be controlled. The personality, including speaking style, and the appearance of the agent for the communications task, Agent Chimp, was designed to be courteous and sophisticated personal agent. Agent Spell,

who handled the spelling game, on the other had a caricature of a studious personality.

The spelling game provided a richer and more challenging application domain than the somewhat simple command and control nature of the personal assistant task. One of the design objectives was to explore how to design interactive educational tutors for children and provide a test bed for investigating automatic dialog strategy adaptation during the course of an interaction. Agent Spell operated in an "intelligent" mode and provided appropriate feedback and guidance to the child user depending on how the game was progressing. The game progress was conveyed through a score meter on the GUI that kept track of the number of attempts per word and overall scores with appropriate accompanying audio prompts. The child could take initiative and choose to play a spelling game at a difficulty level of choice (easy/medium/hard) and in a subject area of choice.

**Dialog Features:** The dialog manager design supported both mixed-initiative (i.e., both the agent and the user can take initiative) and system-initiative strategies. Ambiguity resolution (for example, in the case of the retrieval of multiple records) and error control (using confidence measures provided by utterance verification) were implemented as a part of the dialog strategy. Help messages were available through GUI and through audio prompts. Animation provided a key presentation modality particularly in providing both an engaging interface and useful orientation about the current dialog state. The agent personality transformation when the user switched between one task to another (game to communications agent or vice versa) provided task level orientation. This was accomplished by a animation sequence that consisted of the old agent personality disappearing behind a dropping curtain and the new personality emerging when the curtain lifts back. The following meta dialog presentation features were included: pointing for information presentation (particularly, to indicate the desired selection from among a list), shrugging for conveying retrieval or action failure, and nodding for error conditions or confusions. Sleep and wake agent modes were animated with vanishing and re-appearing agent animation sequences: the agent in sleep mode was represented by a transparent ghost image on the screen which transforms back to the normal image when awakened by appropriate attention command. There were several idle sequences implemented which would automatically put the agent in sleep mode should there be no user action within a specified time out period.

Informal evaluation of the prototype by children users was very positive, especially for the animated agent and the speech interface. Children enjoyed the naturalness and flexibility of

the interface and communicating with the animated agent. Their main criticism was the limited range of interaction one could have with the animated agent (there were only about ten animated sequences) and the lack of understanding of out-of-domain user requests. Formal evaluation of the prototype remains to be done. The prototype can also serve as a test bed for investigating synchronization and merging of input modalities and multimedia presentation.

## IV. Summary

Voice interaction as a component of the multimedia experience fabric is an input modality greatly desired by the children users. The addition of conversational capability to children's multimedia applications contributes to more natural user interactions and improved user experience. The experiments and results reported in this paper show that it is feasible to build conversational systems for children. The inherent variability in children's speech makes ASR difficult. Speaker normalization and model adaptation were used to improve speech recognition performance. A WoZ experiment in a gaming environment provided data for creating novel language models and understanding strategies for dialog systems. Lessons learned from this study and its concurrent user experience evaluation were leveraged in the design of a prototype multimodal-multimedia application for children. Since the system was designed for children users, heavy emphasis was placed on the interface design. Indeed it was found that using animated sequences to communicate information and adding 'personality' to the interface significantly improved the user experience. In addition, the flexible choice of input modality (any of speech, natural language, commands or buttons) made the application easy to use even for novice users. In addition to the user interface, the prototype served as a test bed for creating a general multimodal system architecture. The main design principle of our system was a modular architecture, where the controller communicates with 'stateless' servers via text messages (all state information resides in the template). Seamless integration of all input modalities for our applications is achieved by translating all inputs into text strings that are in turn handled by the spoken language understanding system (or equivalently directly generating an equivalent semantic representation corresponding to certain input events). Other features of our system not yet implemented include customizable application content and customizable agent personality. Overall, the prototype represents a successful first effort in building a multimodal system for children with an emphasis on conversational speech. We hope that data from such prototypes will help further conversational human-machine interaction technology.

## REFERENCES

[1] The American Learning Household Survey, "A study of household demand for Educational Programming, Software and Technology," *Conducted by FIND/SVP,: The Emerging Technology Research Group*, Sep. 1998.

[2] D. C. Burnett and M. Fanty, "Rapid unsupervised adaptation to children's speech on a connected-digit task," in *Proc. ICSLP*, Oct. 1996.

[3] P. Cohen, M. Johnston, D. McGee. and S. Oviatt, J. Clow, and J. Smith, "The efficiency of multimodal interaction: A case study," in *Proc. ICSLP 98*, (Sydney, Australia), 1998.

[4] B. Buntschuh, C. Kamm, G. DiFabbrizio, A. Abella, M. Mohri, S. Narayanan, I. Zeljkovic, J. Wright, S. Marcus, R. D. Sharp, R. Duncan, and J. Wilpon, "VPQ: A Spoken Language Interface to Large Scale Directory Information," in *Proc. ICSLP*, Sydney, Australia, pp. 2863–2867, 1998.

[5] S. Eguchi. and I. J. Hirsh, "Development of speech sounds in children," *Acta. Otolaryng.*, Suppl. 257, 1969.

[6] G. DiFabrizzio, P. Ruscitti, S. Narayanan, and C. Kamm, "Extending Computer Telephony and IP Telephony Standards for Voice-Enabled Services in a Multi-modal User Interface Environment," in *Proceedings of Interactive Dialogue in Multi-modal systems*, Kloster Irsee, Germany, pp. 9–12, June 1999.

[7] U. Goldstein, U. G. "An articulatory model for the vocal tracts of growing children," Ph.D. Thesis, MIT, Cambridge, MA, 1980.

[8] A. Gorin, G. Riccardi, and J. Wright, "How may I help you?," *Speech Commun.*, vol. 23, pp. 113–127, 1997.

[9] Jupiter Communications Research, *cited in CNN* , Aug 11, 1998.

[10] R. D. Kent, "Anatomical and neuromuscular maturation of the speech mechanism: Evidence from acoustic studies," *J. Speech Hear. Res.*, vol. 19, pp. 421–447, 1976.

[11] L. Lamel *et al*, "The LIMSI RailTel system: Field trial of a telephone service for rail travel information," *Speech Commun.*, vol. 23, pp. 67–82, 1997.

[12] L. Lee and R. C. Rose, "Speaker normalization using efficient frequency warping procedures," in *Proc. ICASSP*, pp. 353–356, May 1996.

[13] S. Lee, A. Potamianos, and S. Narayanan, "Acoustics of children's speech: Developmental changes of temporal and spectral parameters," *J. Acoust. Soc. Am.*, vol. 105, pp. 1455–1468, Mar. 1999.

[14] J. Mostow, A. G. Hauptmann, and S. F. Roth, "Demonstration of a reading coach that listens," *Proceedings of the ACM Symposium on User Interface Software and Technology*, pp. 77–78, 1995.

[15] MTV Network, "Leisure Time Study ," *cited in CNN* , Aug 11, 1998.

[16] Newsweek, "Teenagers and Technology," p. 86, April 28, 1997.

[17] Nickelodeon, "Online Attitude and Usage Study," *cited in CNN* , Aug 11, 1998.

[18] R. Pieraccini and E. Levin, "A spontaneous-speech understanding system for database query applications ," in *ESCA Workshop on Spoken Dialogue Systems - Theories and Applications*, 1995.

[19] R. Pieraccini, E. Levin, and W. Eckert, "AMICA: the AT&T Mixed Initiative Conversational Architecture," in *Proc. EuroSpeech*, (Rhodes, Greece), Sept. 1997.

[20] A. Potamianos et al, "Design principles and tools for multimodal dialog systems," in *Proc. ESCA Workshop Interact. Dialog. Multi-Modal Syst.*, (Kloster Irsee, Germany), June 1999.

[21] A. Potamianos and S. Narayanan, "Spoken dialog systems for children," in *Proc. ICASSP*, (Seattle, WA), pp. 197–200, May 1998.

[22] A. Potamianos, S. Narayanan, and S. Lee, "Automatic speech recognition for children," in *Proc. EuroSpeech*, (Rhodes, Greece), pp. 2371–2374, Sept. 1997.

[23] A. Potamianos and S. Narayanan, "Automatic speech recognition and semantic understanding for children," in preparation.

[24] A. Potamianos, G. Riccardi, and S. Narayanan, "Categorical understanding using statistical N-gram models," in *Proc. EuroSpeech*, (Budapest, Hungary), Sept. 1999.

[25] A. Potamianos and R. C. Rose, "On combining frequency warping and spectral shaping in HMM-based speech recognition," in *Proc. ICASSP*, Apr. 1997.

[26] L. Rabiner and B. -H. Juang, *Fundamentals of Speech Recognition.* Englewood Cliffs, NJ: Prentice Hall, 1993.

[27] G. Riccardi, A. Potamianos, and S. Narayanan, "Language model adaptation for spoken language systems," in *Proc. ICSLP 98*, (Sydney, Australia), pp. 2327–2330, 1998.

[28] M. Russell, B. Brown, A. Skilling, R. Series, J. Wallace, B. Bonham, and P. Barker, "Applications of automatic speech recognition to speech and language development in young children," in *Proc. ICSLP*, Philadelphia, PA, Oct. 1996.

[29] R. Sharma, V. Pavlovic, and T. Huang, "Toward multimodal human computer interface," *Proceedings of the IEEE*, vol. 86, no. 5, pp. 853–869, 1998.

[30] E. F. Strommen and F. S. Frome, "Talking back to big bird: Preschool users and a simple speech recognition system," *Educational Technology Research and Development*, vol. 41, pp. 5–16, 1993.

[31] T. Takezawa and T. Morimoto, "A multimodal-input multimedia-output guidance system: MMGS," in *Proc. ICSLP 98*, (Sydney, Australia), 1998.

[32] J. G. Wilpon and C. N. Jacobsen, "A study of automatic speech recognition for children and the elderly," in *Proc. ICASSP*, pp. 349–352, May 1996.

[33] V. Zue, S. Seneff, J. Glass, J. Polifroni, C. Pao, T. Hazen, and L. Hetherington, "Jupiter: A telephone-based conversational interface for weather information," *IEEE Trans. Speech Audio Processing*, vol. 8, Jan. 2000, pp. 85-96.
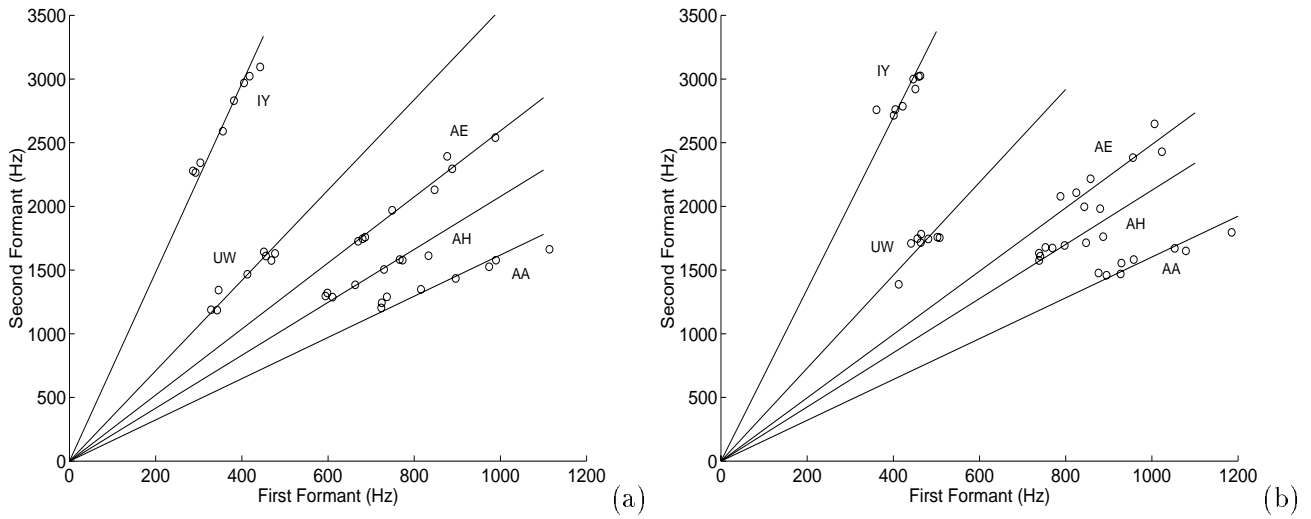
Fig. 1. Vowel F1-F2 scaling as a function of age for (a) Male speakers and (b) Female speakers (seven data points per vowel for ages 5-6, 7-8, ..., 17-18).
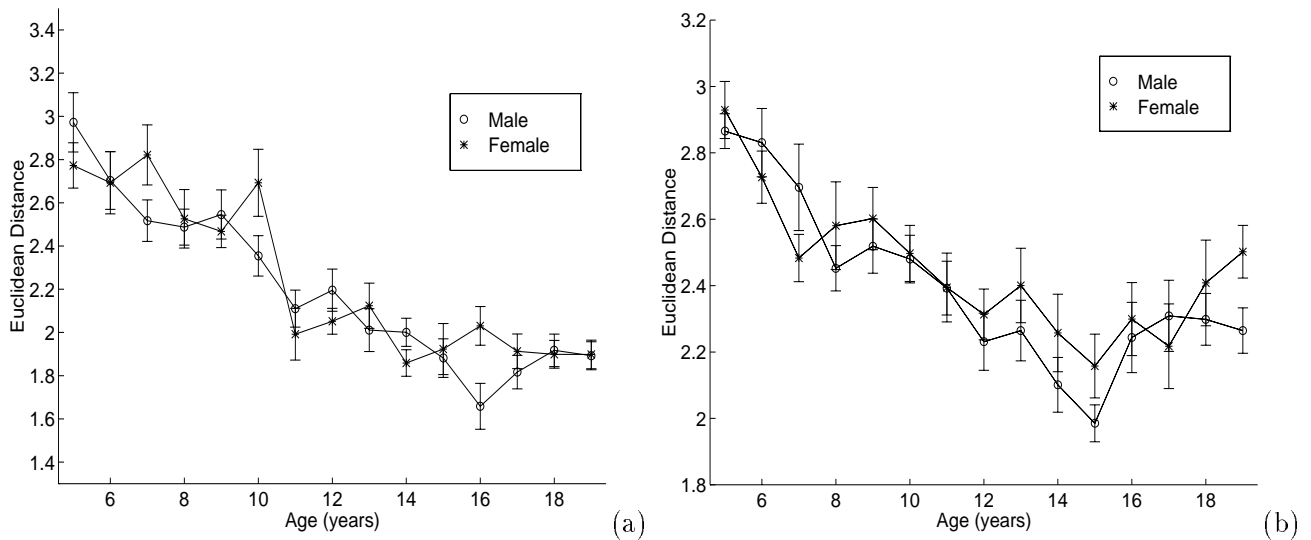


Fig. 2. Intra-speaker variability as a function of age: (a) Mean cepstral distance between the two repetitions of the same vowels and (b) Mean cepstral distance between the first- and second-half segments within the same vowel realization.

| Name | Speaker Population | Content | No. of speakers | No. of strings |
|---|---|---|---|---|
| DgtI | Adults | digits | 3026 | 4781 |
| DgtII | 10-17 yrs. | digits | 1234 | 5767 |
| SubwI | Adults | phrases | 242 | 12144 |
| SubwII | 10-17 yrs. | phrases | 1234 | 14267 |
| DgtTest | 6-17 yrs. | digits | 501 | 2656 |
| CommTest | 6-17 yrs. | commands | 501 | 3554 |

TABLE I

Training and testing databases.

| Model | Baseline | Norm. | Improv. |
|---|---|---|---|
| Adult HMM | 15.9% | 8.7% | +45% |
| Children HMM | 6.7% | 4.9% | +25% |
| Cld+Adlt HMM | 7.6% | 5.6% | +25% |

TABLE II

Digit error rate for children speakers before and after speaker normalization.

| | Age | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Gnd | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 8-14 | >21 |
| F | 18 | 23 | 32 | 24 | 10 | 8 | 4 | 119 | 5 |
| M | 21 | 51 | 16 | 23 | 21 | 25 | 14 | 171 | 8 |

TABLE III

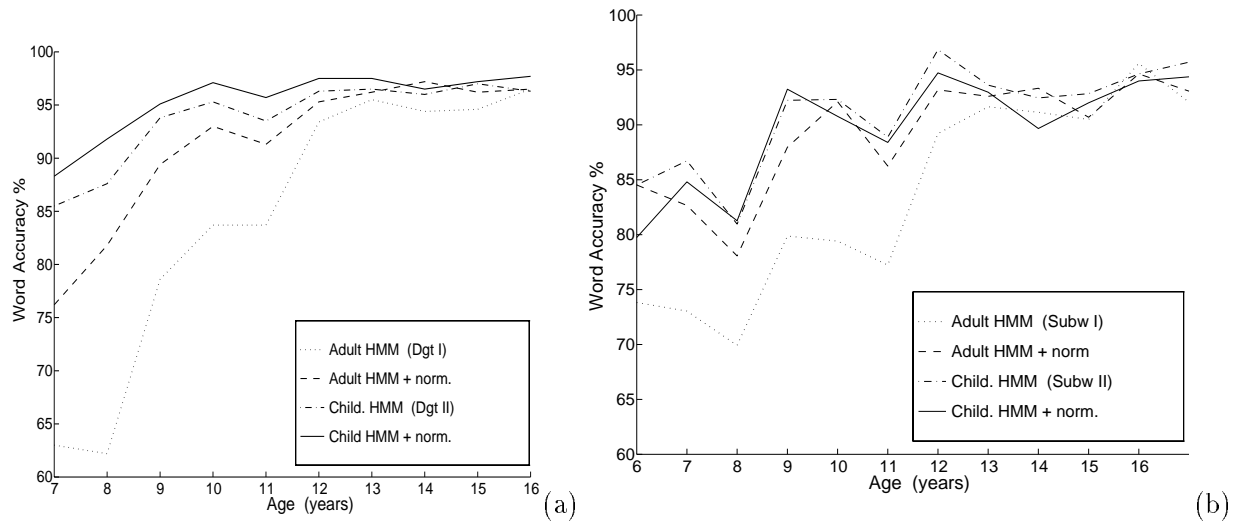Number of games per player's age and gender.

Fig. 3. Word accuracy (%) vs. speaker's age using HMMs trained from children ("Child. HMM") or adult ("Adult HMM") speaker population with ("norm.") and without speaker normalization for: (a) connected digit task, (b) command and control task.
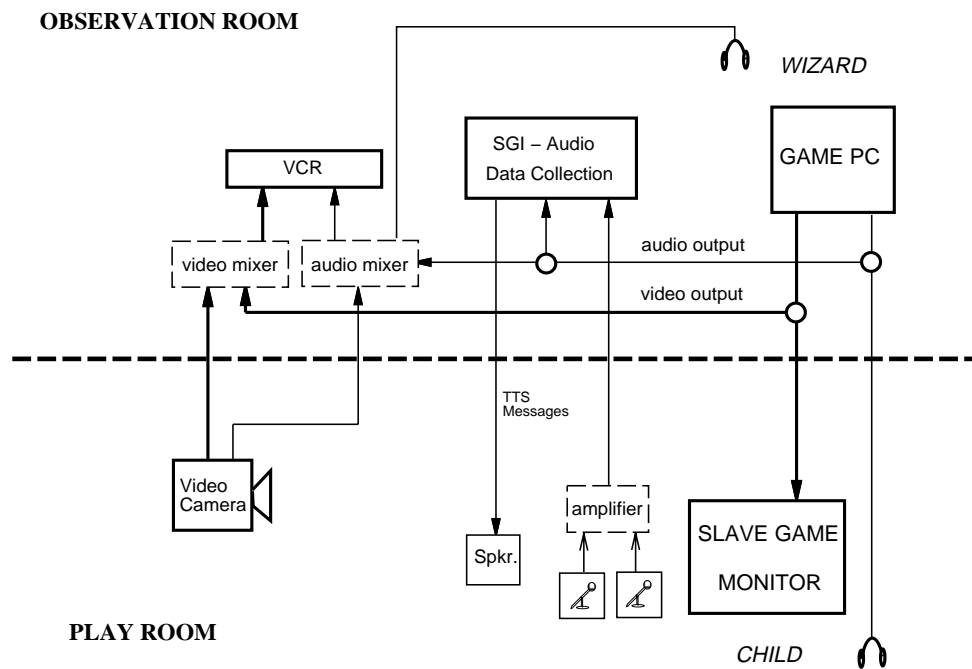


Fig. 4. The experimental WoZ setup.

| User input/System output | Dialogue State |
|---|---|
| User: Tell me about the suspect? <br><br> System: She is neither long- nor short-legged | $S_{t-3}$: TellmeAbout |
| U: Her <u>height</u> is <u>average</u> <br><br> S: ... [updating suspect's drawing] | $S_{t-2}$: EnterFeature |
| U: Where did the suspect go? <br><br> S: She is picking peonies in Bloomington | $S_{t-1}$: WhereDid |
| U: Go to <u>Indiana</u> <br><br> S: ... [travel theme] | $S_t$: GoToState |

TABLE IV

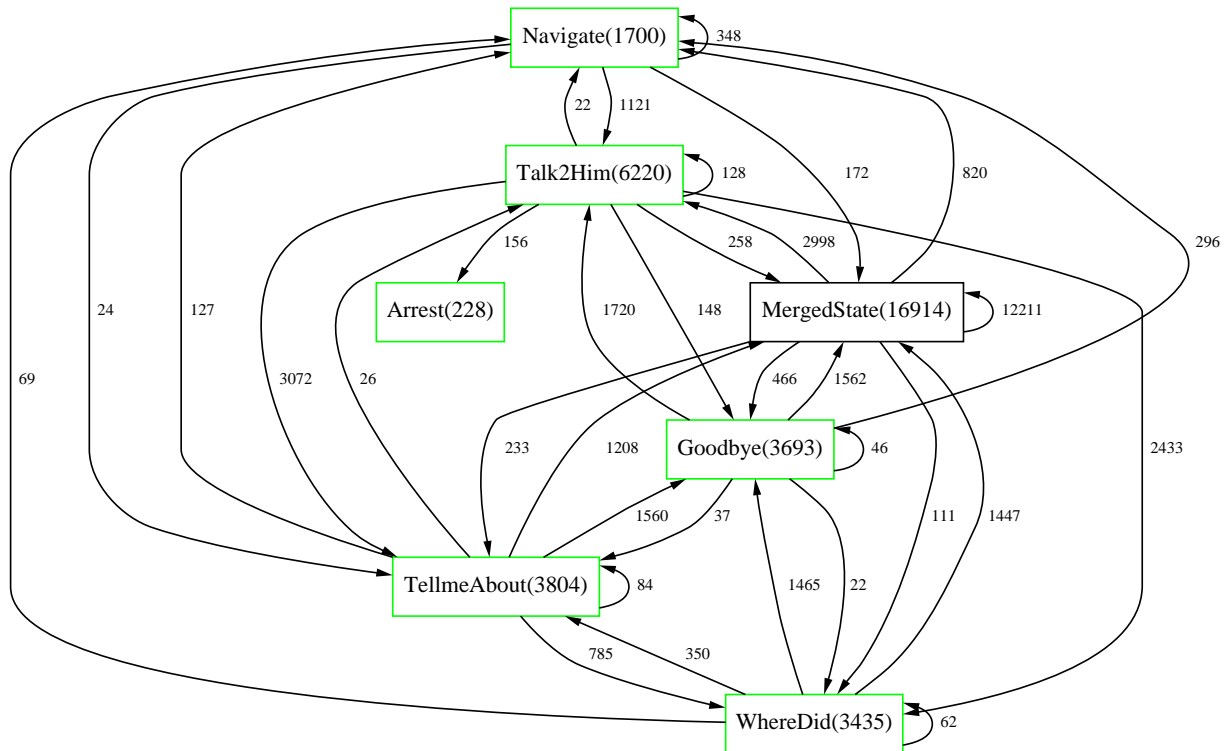Transcript of a sample interaction along with dialog state tags.



Fig. 5. Dialog state and state transition diagram (with counts) for all children players for the navigation/query subdialog ("MergedState" denotes combination of all dialog states not shown in plot).

|  | Dialog State | | | | |
|---|---|---|---|---|---|
| Variability | 1 | 2 | 3 | 4 | 5 |
| Intra-speaker | 0.43 | 0.26 | 0.22 | 0.32 | 0.40 |
| Inter-speaker | 1.05 | 0.48 | 0.40 | 0.54 | 0.72 |

TABLE V

INTER- AND INTRA-SPEAKER LINGUISTIC VARIABILITY FOR DIALOG STATES "TALK2HIM" (1), "WHEREDID" (2), "TELLMEABOUT (3)", "GOODBYE" (4), "OPENCLUEBOOK" (5).
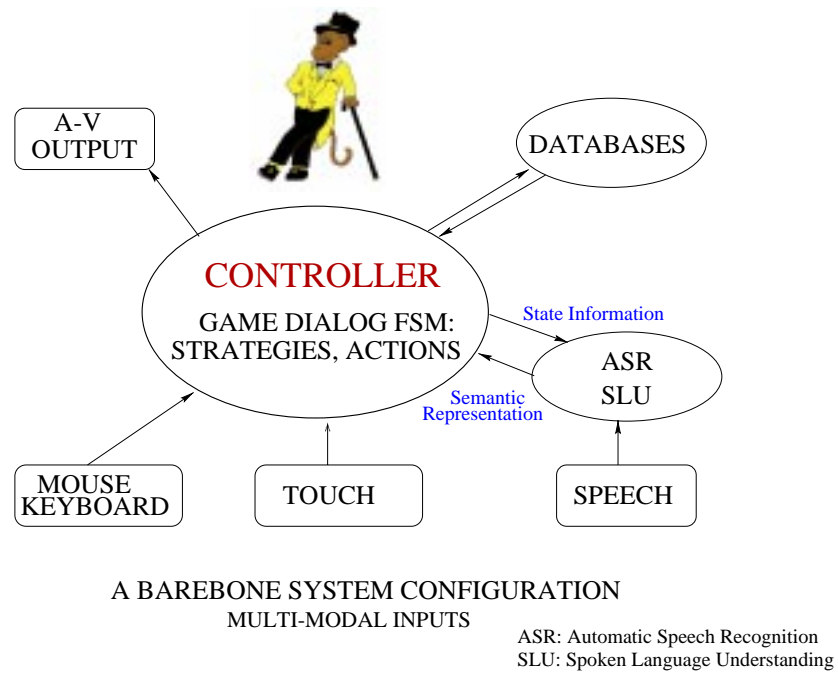


Fig. 6. Functional block diagram of the CHIMP prototype.

speech input → ASR CLIENT → WATSON ASR SERVER

keyboard and mouse input

**GUI**

**EVENT HANDLER**

sync

GRAPHICAL INPUT PROC.

INPUT EVENT HANDLER

AMICA DIALOG SERVER

user echo

blocking

GRAPHICAL OUTPUT GEN.

template

OUPUT EVENT HANDLER

template

spawn and wait

movie

text images buttons

APPLICATION (call,email,web)

text

TTS SERVER

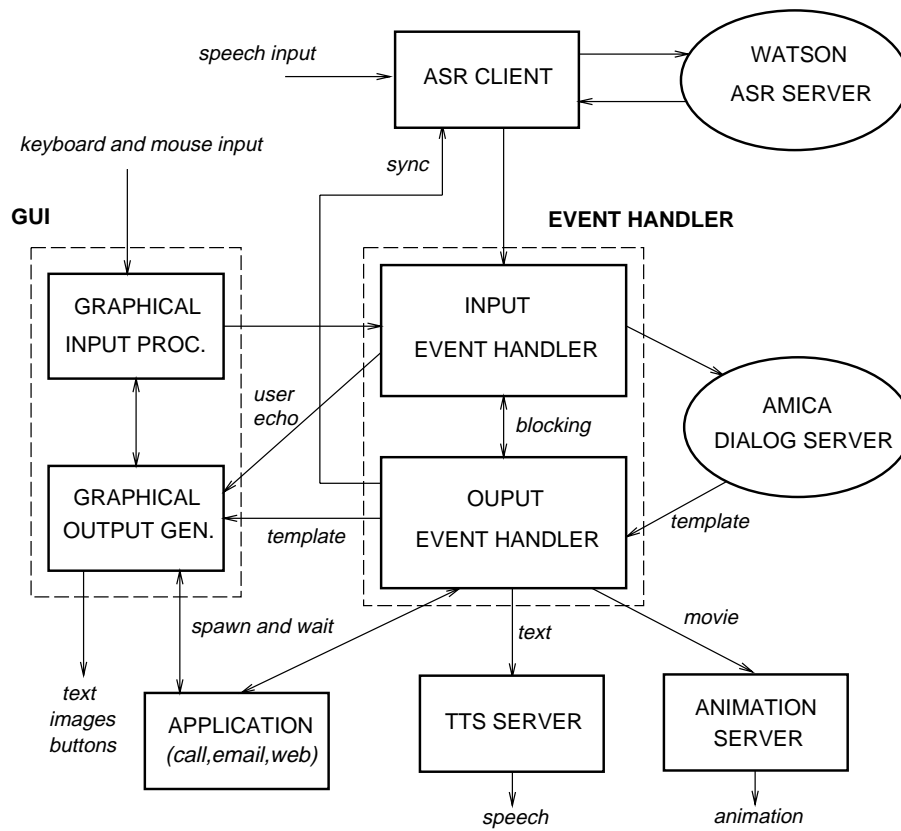ANIMATION SERVER

speech

animation
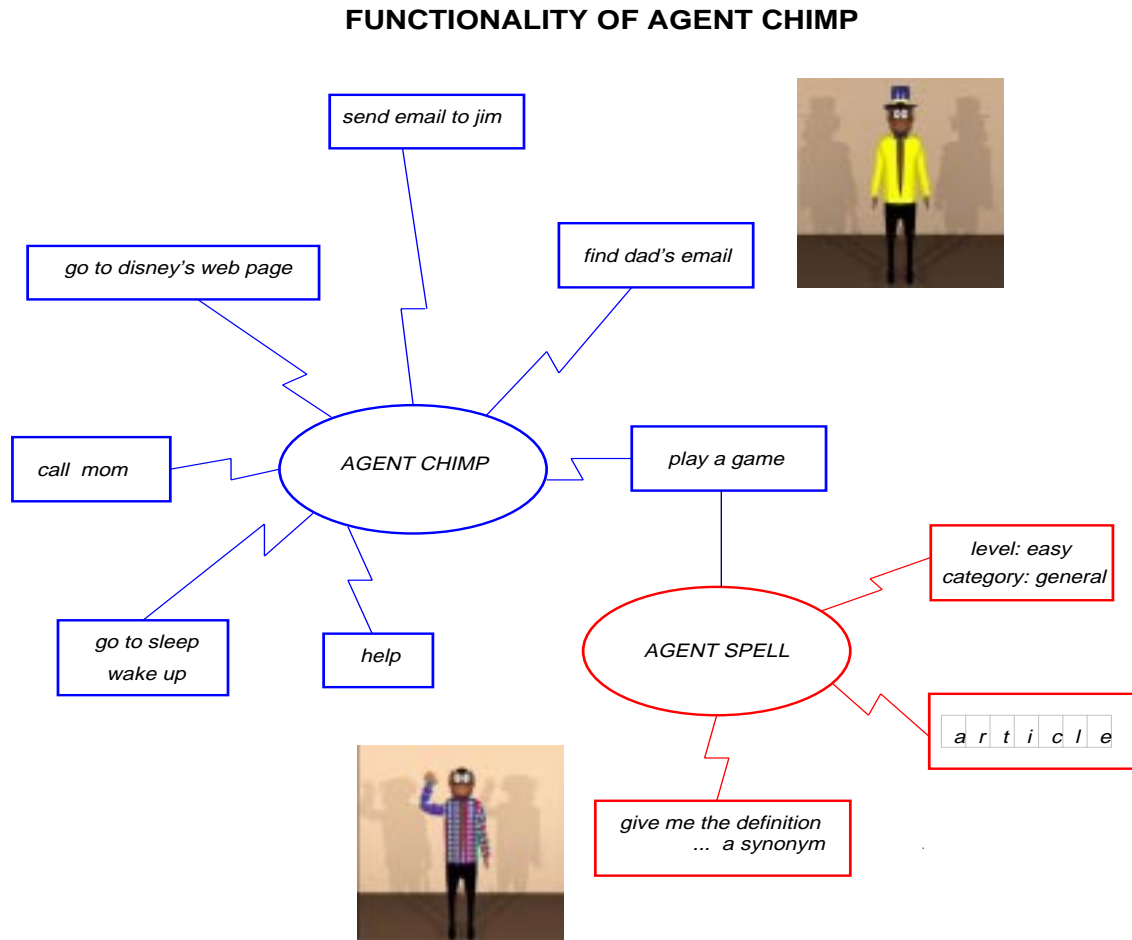
Fig. 7. Architectural diagram of CHIMP prototype.

Fig. 8. Overview of some of the application features: Agent Chimp handling the personal agent task and Agent Spell handling the spelling game.

Fig. 9. Screen shot of the prototype graphical user interface.