

Towards Developing General Models of Usability with PARADISE

Marilyn Walker, Candace Kamm and Diane Litman

*AT&T Labs—Research
180 Park Avenue
Florham Park, NJ 07932-0971 USA*

(Received 3 November 1999)

Abstract

The design of methods for performance evaluation is a major open research issue in the area of spoken language dialogue systems. In this paper we present the PARADISE methodology for developing predictive models of spoken dialogue performance, and then show how to evaluate the predictive power and generalizability of such models. To illustrate our methodology, we develop a number of models for predicting system usability (as measured by user satisfaction), based on the application of PARADISE to experimental data from three different spoken dialogue systems. We then measure the extent to which our models generalize across different systems, different experimental conditions, and different user populations, by testing models trained on a subset of our corpus against a test set of dialogues. Our results show that our models generalize well across our three systems, and are thus a first approximation towards a general performance model of system usability.

1 Introduction

The development of methods for evaluating, comparing and predicting system performance is a major open research issue in the area of spoken language dialogue systems (SLDSs). In 1997, we proposed PARADISE (PARAdigm for DIalogue System Evaluation) as a general integrative framework for evaluating SLDSs (Walker et al., 1997b). PARADISE addressed three research goals:

- We wanted to be able to compare multiple systems or multiple versions of the same system doing the same domain tasks.
- We wanted to develop predictive models of the usability of a system as a function of a range of system properties.
- We wanted to be able to make generalizations across systems about which properties of a system impact usability, i.e. to figure out “what really matters to users”.

In previous work, we conducted a number of within-system comparisons, based on the application of PARADISE to experimental dialogues collected with three different systems: ANNIE, a SLDS for voice dialing and messaging (Kamm et al.,

1998); ELVIS, an SLDS for accessing email (Walker et al., 1998b; Walker et al., 1998a); and TOOT, an SLDS for accessing online train schedules (Litman et al., 1998; Litman and Pan, 1999). In this paper, we focus on the predictive power and generalizability of our models. We first develop a number of models for predicting system usability (as measured by user satisfaction). Then we measure the extent to which these models generalize across different systems, different experimental conditions, and different user populations, by testing models trained on a subset of our corpus against a test set of dialogues. Our experimental corpus consists of 544 ANNIE, ELVIS, and TOOT dialogues, totaling almost 42 hours of speech.

The structure of the paper is as follows. We first review the PARADISE evaluation framework in section 2. Then section 3 describes the ANNIE, ELVIS and TOOT systems and the experimental setup for collecting the corpus of dialogues. Section 4 presents results from developing predictive models of user satisfaction and testing these models on subsets of the corpus. We then discuss our results and suggest areas for future work.

2 PARADISE

The state of the art in evaluating SLDSs, prior to the proposal of the PARADISE framework, was to evaluate the SLDS in terms of a battery of both subjective and objective metrics. Some of these metrics focused on task completion or transaction success. Others were based on the performance of the SLDS's component technologies, such as speech recognizer performance (Sparck-Jones and Galliers, 1996; Hirschman et al., 1990; Ralston et al., 1995; Pallett, 1985). Subjective metrics included measures of user satisfaction (Shriberg et al., 1992), or ratings generated by dialogue experts as to how cooperative the system's utterances were (Bernsen et al., 1996).

The motivation for PARADISE is that problems arise when a battery of metrics is used. First, several different metrics may contradict one another. For example, Danieli and Gerbino compared two train timetable agents (Danieli and Gerbino, 1995), and found that one version of their SLDS had a higher transaction success rate and produced fewer inappropriate and repair utterances, but that the other version produced dialogues that were approximately half as long. However, they could not report whether the higher transaction success or the length of the dialogue was more critical to performance, or what tradeoffs between them might be acceptable.

Second, in order to make generalizations across different systems performing different tasks, it is important to know how multiple factors impact performance and how users' perceptions of system performance depend on the dialogue strategy and on tradeoffs among other factors like efficiency, usability and accuracy.

The PARADISE framework derives a combined performance metric for a dialogue system as a weighted linear combination of a task-based success measure and dialogue costs. In order to specify what factors should go into this combined performance metric, PARADISE posits a particular model of performance, illustrated in Figure 1. The model proposes that the system's primary objective is to maxi-

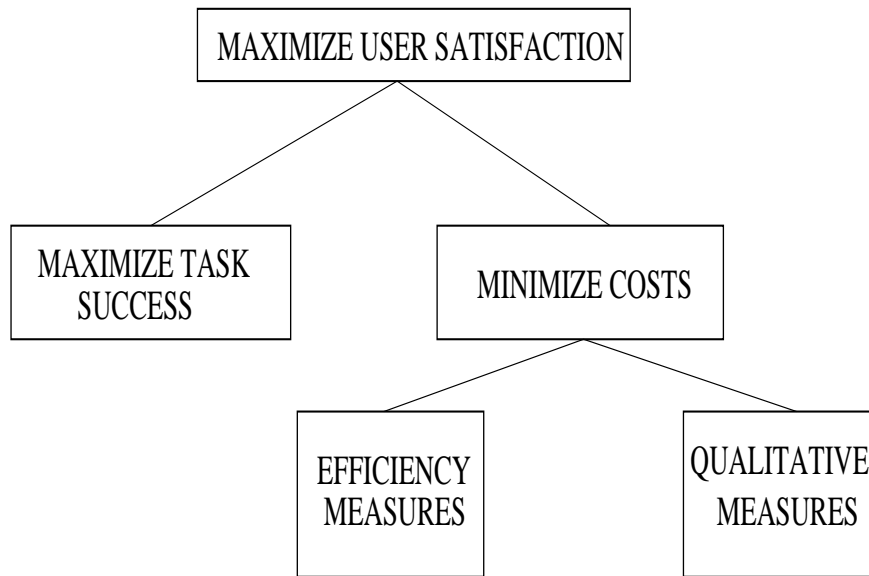


Fig. 1. PARADISE's structure of objectives for spoken dialogue performance.

mize user satisfaction. Task success and various costs that can be associated with the interaction are both contributors to user satisfaction. Dialogue costs are of two types: dialogue efficiency and quality. Efficiency costs are measures of the system's efficiency in helping the user complete the task, such as the number of utterances to completion of the dialogue. Dialogue quality costs are intended to capture other aspects of the system that may have strong effects on user's perception of the system, such as the number of times the user had to repeat an utterance in order to make the system understand the utterance.¹

Applying PARADISE to dialogue data requires that dialogue corpora be collected via controlled experiments during which users subjectively rate their satisfaction. In addition, the other components of the model, i.e. task success and interaction costs, must be either automatically logged by the system or be hand-labeled. The PARADISE model posits that a performance function can then be derived by applying multivariate linear regression with user satisfaction as the dependent variable and task success, dialogue quality, and dialogue efficiency measures as independent variables.

Modeling user satisfaction as a function of task success and dialogue cost metrics is intended to lead to predictive performance models for SLDSs, so that values for user satisfaction could be predicted on the basis of a number of simpler metrics that can be directly measured from the system logs, without the need for extensive experiments with users to assess user satisfaction. In order to make this predictive use of PARADISE a reality, the models that are derived from experiments with

¹ Developing qualitative cost measures that can be measured automatically is an active area of research (Walker and Hirschman, 1999).

one set of systems or user populations should be **generalizable** to other systems or other user populations. By applying PARADISE to three different systems, we can show what generalizations can be made across systems and user populations.

3 Experimental Design

3.1 *Experimental Setup*

All of our experiments applying PARADISE used a similar experimental setup. In each experiment, human subjects carried out a dialogue using one of the three dialogue systems. Each system was implemented using a general-purpose platform for phone-based spoken dialogue systems (Kamm et al., 1997). The platform consisted of a speech recognizer that supports barge-in so that the user can interrupt the agent when it is speaking. It also provided an audio server for both voice recordings and text-to-speech (TTS), an interface between the computer running the system and the telephone network, a module for application specific functions, and modules for specifying the application grammars and the dialogue manager.

The dialogues were obtained in controlled experiments. The ANNIE experiments (Kamm et al., 1998) were designed to evaluate (1) the effect of novice versus expert user populations, and (2) the impact of a short tutorial on the novice user population. The ELVIS experiments (Walker et al., 1997a; Walker et al., 1998b; Walker et al., 1998a) were designed (1) to evaluate the effect of different dialogue strategies for managing initiative and presentation of information; and (2) to conduct experiments applying reinforcement learning to dialogue so that the ELVIS system could learn from experience which dialogue strategies are optimal. The TOOT1 experiments (Litman et al., 1998) were designed to evaluate the effects of: (1) task difficulty and (2) cooperative versus literal response strategies. The TOOT2 experiments studied the effect of user-adapted interaction (Litman and Pan, 1999).

All of the experiments required users to complete a set of application tasks in conversations with a particular version of the system. Following PARADISE, we characterized these tasks in terms of particular items of information that the user must find out from the agent. Each set of tasks was also designed to be representative of typical tasks in that domain.

Experiments with ANNIE were designed to be comparable with ELVIS so the users in the ANNIE experiments performed the same tasks as those for ELVIS. Each user had to perform three tasks in sequence in three different conversations with the system. To examine the effect of a short tutorial, some users also performed an additional task as part of a tutorial interaction before starting the experimental tasks. The example email access task below was used for both ANNIE and ELVIS:

- You are working at home in the morning and plan to go directly to a meeting when you go into work. Kim said she would send you a message telling you where and when the meeting is. Find out **the Meeting Time** and **the Meeting Place**.

For TOOT, each user performed 4 tasks in sequence. Each task was represented by a scenario where the user had to find a train satisfying certain constraints,

by using the agent to retrieve and process online train schedules. A sample task scenario is as follows:

- Try to find a train going to **Boston** from **New York City** on **Saturday** at **6:00 pm**. If you cannot find an exact match, find the one with the **closest** departure time. Please write down the **exact departure time** of the train you found as well as the **total travel time**.

The experiments resulted in 108 dialogues with the ANNIE system, 268 dialogues with the ELVIS system, and 168 dialogues with the TOOT system for a total of 544 dialogues, consisting of almost 42 hours of speech.

3.2 Data Collection

Here, we discuss in detail the metrics that were logged and hand-labeled in experiments with all three systems. We have not previously made specific recommendations about what metrics to use, as we consider this an active area of research. However, any metric for task success, dialogue efficiency or dialogue quality previously proposed in the literature can be easily incorporated into the PARADISE framework.

In these experiments, we decided to use a combination of dialogue quality and efficiency measures, mainly focusing on those measures that could be automatically logged or computed. We used three different methods to collect the data: (1) All of the dialogues were recorded; (2) The dialogue manager logged the agent’s dialogue behavior and a number of other measures discussed below; (3) Users filled out web page forms after each task (task success and user satisfaction measures). Measures are summarized in Figure 2 and described in more detail below.

- **Dialogue Efficiency Metrics**
 - elapsed time, system turns, user turns
- **Dialogue Quality Metrics**
 - mean recognition score, timeouts, rejections, helps, cancels, bargeins (raw)
 - timeout%, rejection%, help%, cancel%, bargein% (normalized)
- **Task Success Metrics**
 - task completion as per survey
- **User Satisfaction**
 - the sum of TTS Performance, ASR Performance, Task Ease, Interaction Pace, User Expertise, System Response, Expected Behavior, Comparable Interface, Future Use.

Fig. 2. Metrics collected for spoken dialogues.

The **dialogue efficiency** metrics were calculated from the dialogue recordings and the system logs. The length of the recording was used to calculate the elapsed time in seconds (**ET**) from the beginning to the end of the interaction. Measures for the number of **System Turns**, and the number of **User Turns**, were calculated

on the basis of the system logging everything it said and everything it heard the user say.

The **dialogue quality** measures were derived from the recordings, the system logs and hand-labeling. A number of agent behaviors that affect the quality of the resulting dialogue were automatically logged. These included the number of timeout prompts (**timeouts**) played when the user didn't respond as quickly as expected, the number of recognizer rejections (**rejects**) where the system's confidence in its understanding was low and it said something like *I'm sorry I didn't understand you*. User behaviors that the system perceived that might affect the dialogue quality were also logged: these included the number of times the system played one of its context specific help messages because it believed that the user had said *Help* (**helps**), and the number of times the system reset the context and returned to an earlier state because it believed that the user had said *Cancel* (**cancel**s). The recordings were used to check whether users barged in on agent utterances, and these were labeled on a per-state basis (**bargeins**).

Another measure of dialogue quality was recognizer performance over the whole dialogue, calculated in terms of concept accuracy. The recording of the user's utterance was compared with the logged recognition result to calculate a concept accuracy measure for each utterance by hand. Concept accuracy is a measure of semantic understanding by the system, rather than word for word understanding. For example, the utterance *Read my messages from Kim* contains two concepts, the *read* function, and the *sender:kim* selection criterion. If the system understood only that the user said *Read*, then concept accuracy would be .5. Mean concept accuracy was then calculated over the whole dialogue and used, in conjunction with ASR rejections, to compute a Mean Recognition Score **MRS** for the dialogue.

Because our goal is to generate models that will generalize across systems and tasks, we also thought it important to introduce metrics that are likely to generalize. All of the efficiency metrics seemed unlikely to generalize since, e.g. the **elapsed time** to complete of a task depends on how difficult the task is. Other research suggested that the dialogue quality metrics were more likely to generalize (Litman et al., 1999), but we thought that the raw counts were likely to be task specific. Thus we normalized the dialogue quality metrics by dividing the raw counts by the total number of utterances in the dialogue. This resulted in the **timeout%**, **rejection%**, **help%**, **cancel%**, and **bargein%** metrics.

The web page forms are the basis for calculating Task Success and User Satisfaction measures. Users reported their perceptions as to whether they had completed the task (**Comp**).² They also had to provide objective evidence that they had in fact completed the task by filling in a form with the information that they had acquired from the agent.³

² *Yes, No* responses are converted to *1, 0*.

³ This supports an alternative way of calculating **Task Success** objectively by using the Kappa statistic to compare the information that the users filled in with a key for the task (Walker et al., 1997b). However some of our earlier results indicated that user's *perception* of task success was a better predictor of overall satisfaction, so here we simply use perceived task success as measured by **Comp**.

- Was the system easy to understand in this conversation? (**TTS Performance**)
- In this conversation, did the system understand what you said? (**ASR Performance**)
- In this conversation, was it easy to find the message you wanted? (**Task Ease**)
- Was the pace of interaction with the system appropriate in this conversation? (**Interaction Pace**)
- In this conversation, did you know what you could say at each point of the dialogue? (**User Expertise**)
- How often was the system sluggish and slow to reply to you in this conversation? (**System Response**)
- Did the system work the way you expected him to in this conversation? (**Expected Behavior**)
- In this conversation, how did the system’s voice interface compare to the touch-tone interface to voice mail? (**Comparable Interface**)
- From your current experience with using the system to get your email, do you think you’d use the system regularly to access your mail when you are away from your desk? (**Future Use**)

Fig. 3. User satisfaction survey for ELVIS and ANNIE experiments.

In order to calculate User Satisfaction, users were asked to evaluate the agent’s performance with a user satisfaction survey. The survey probed a number of different aspects of the users’ perceptions of their interaction with the SLDS in order to focus the user on the task of rating the system, as in (Shriberg et al., 1992; Jack et al., 1992; Love et al., 1994). A sample survey is in Figure 3. The surveys used for the three SLDSs were identical except that the **Comparable Interface** question was eliminated from the TOOT survey. The surveys were multiple choice and each survey response was mapped into the range of 1 to 5. Then the values for all the responses were summed, resulting in a **User Satisfaction** measure for each dialogue ranging from 8 to 40.

4 Results of Applying PARADISE

Table 1 summarizes the experimental corpus of 544 dialogues in terms of experimental conditions and the metrics that were collected. This section reports results for

- training models to predict user satisfaction via multivariate linear regression.
- testing these models across different systems, different experimental conditions, and different user populations to determine to what extent they generalize.

An overall summary of our results for testing and training predictive models is in Table 2. Each model discussed was trained by performing a stepwise multivariate linear regression on some set of dialogues with user satisfaction as the dependent variable and the other metrics discussed above and shown in Table 1 as the independent variables. For each model we describe the training set, list the significant factors of the model, indicate the performance of the trained model on the training

Table 1. Performance measure means per dialogue for different SLDSs.

Measure	SI ^a	MI ^b	Exps ^c	NovTut ^d	Novs ^e	Fixed ^f	Adapt ^g
Comp	.87	.80	1.0	.96	.73	.75	.87
User Turns	21.5	17.0	10.8	13.5	21.4	13.6	15.3
Sys Turns	24.2	21.2	11.7	15.4	25.5	14.0	15.7
ET ^h	339.14 s	296.18 s	156.5 s	195.6	280.3 s	238.6 s	234.1 s
MRS ⁱ	.88	.72	.80	.74	.67	.77	.78
TimeOuts	2.65	4.15	.64	1.28	2.4	.32	.25
TimeOut%	.11	.19	.04	.08	.10	.02	.02
Cancs	.34	.02	.17	.17	.61	1.02	.42
Canc%	.01	0	.01	.01	.03	.07	.02
Helps	.67	.92	.08	.75	3.3	.26	.05
Help%	.03	.05	.01	.05	.14	.02	.00
Bargelns	3.6	3.6	4.5	4.5	5.8	1.0	1.6
Bargeln%	.08	.09	.2	.17	.15	.08	.10
Rejects	.86	1.58	1.8	3.3	6.6	1.1	1.1
Reject%	.04	.08	.17	.22	.25	.07	.08
USAT ^j	28.9	25.0	32.8	30.8	23.0	28.3	31.0

^a ELVIS system initiative.

^b ELVIS mixed initiative.

^c ANNIE experts.

^d ANNIE novices who had a tutorial.

^e ANNIE novices.

^f TOOT1 and non-adaptable TOOT2.

^g Adaptable TOOT2.

^h Elapsed time.

ⁱ Mean recognition score.

^j User satisfaction.

set, describe the test set, and indicate the performance of the trained model on the test set.

We first trained models for each system on a random 90% of the dialogues for that system. Results for each system are given in the first three rows of Table 2, showing the degree to which the model is able to predict user satisfaction in the training set (Column labeled R^2 Training).⁴ We then tested whether a model trained on a random 90% of the corpus for a particular system could predict user satisfaction in the remaining 10% of the dialogues for that system. The results in the first three rows of Table 2 suggest that that the TOOT model generalizes most readily, with the model fits for ELVIS and ANNIE being somewhat worse.⁵

⁴ For the first 4 rows of the table, 10-fold cross-validation is used to estimate R^2 . Thus, the training and testing R^2 values represent means, and standard errors (SE) are shown. For the rest of the table, the R^2 values are obtained using a single training set and a separate held-out test set.

⁵ However, as we discuss in more detail below, the predictive capability of a model is dependent on the size of the training set and there are fewer ANNIE dialogues in the corpus.

Table 2. *Testing the predictive power of different models.*

Training Set ^a	R ² Training (SE) ^b	Test Set ^c	R ² Test (SE) ^d
ANNIE 90%	.50 (.009)	ANNIE 10%	.40 (.07)
ELVIS 90%	.39 (.003)	ELVIS 10%	.43 (.03)
TOOT 90%	.56 (.014)	TOOT 10%	.54 (.05)
ALL 90%	.47 (.004)	ALL 10%	.50 (.035)
ELVIS 90%	.42	TOOT	.55
ELVIS 90%	.42	ANNIE	.36
NOVICES	.47	ANNIE EXPERTS	.04
MRS ≤ .95	.46	MRS > .95	.23

^a The training set used to train the model.

^b The amount of variance in user satisfaction in the training set accounted for by the trained model. When 10 fold cross-validation is used to estimate R², the value reported is a mean and the standard error (SE) is shown.

^c A characterization of the test set.

^d The amount of variance in user satisfaction in the test set accounted for by the trained model.

We then examined models trained on the combined data for all three systems. The ALL row in Table 2 gives the results for training on a random 90% of the dialogues from the combined corpus and testing on the remaining 10% of the dialogues. As the table shows, the model is a good predictor of user satisfaction in the test set, accounting for 50% of the variance in R².

Next we examined whether models trained on a random 90% of one of the systems (ELVIS) could predict user satisfaction for users using one of the other systems (ANNIE, TOOT). The results for this test are shown in rows 5 and 6 of Table 2. Surprisingly, the ELVIS model accounts for 55% of the variance in user satisfaction for the TOOT1 and TOOT2 data sets, accounting for more of the variance in TOOT than it does on the ELVIS test set. Also the ELVIS model accounts for 36% of the variance in user satisfaction in the ANNIE data set.

We then turn to the question of whether models trained on one user population generalize to a different type of user population. We examined this in two ways.

First, we had a group of expert users who participated in the ANNIE experiments. These users had been using the ANNIE system for over 6 months on a regular basis. All of the other users who participated in our experiments were novice users who were not familiar with the limitations of spoken dialogue systems. Out of the 544 dialogues, 36 were dialogues with expert users and the remaining 508 dialogues were collected with novice users. We trained a model with the dialogues from the novice users and tested it against the dialogues with the experts. The results are shown in Table 2 in the NOVICES row. The results show that models trained on novice users **do not** easily generalize to an expert user population. The correlation of the predicted values to the actual values is only .19; the novice model only accounts for 4% of the variance in user satisfaction in the expert population.

Second, speech recognition performance varies a great deal from user to user, and we were interested in testing whether the models trained on ANNIE, ELVIS and

Table 3. *Factors that have the most predictive power in different models.*

Training Set ^a	Factors ^b
ANNIE 90%	help%, mrs, comp, bargein%
ELVIS 90%	mrs, comp, reject%, bargein%
TOOT 90%	mrs, comp, et, reject%
ALL 90%	mrs,comp,reject%, help%, bargein%, et
NOVICES	mrs, comp, reject%, bargein%, help%, et, timeout%
MRS \leq .95	mrs, comp, reject%, help%, et, bargein%, timeout%

^a The training set used to train the model.

^b Significant predictors in order of contribution.

TOOT would generalize to systems whose speech recognition performance overall was much better. In order to examine this question, we determined that 93 out of 544 of the dialogues had a Mean Recognition Score (MRS) better than .95 (GoodASR dialogues). We used the 451 dialogues where MRS was less than or equal to .95 to train a model for predicting user satisfaction and then tested it on the GoodASR dialogues. The results shown in the MRS \leq .95 row of Table 2 show that the model only accounts for 23% of the variance in R^2 in the test data. This is not surprising since the contribution of MRS to user satisfaction is large in the model trained on the rest of the corpus.

Table 3 also shows which factors were found to be significant predictors of user satisfaction in each of the trained models (Factors column), ordered by degree of contribution. Except for the ANNIE system model, recognizer performance (MRS) and task success (COMP) are always the largest contributors to user satisfaction, followed by the percentage of the dialogue turns that were rejections (REJECT%), where the system asked the user to repeat or paraphrase what they had just said. The models for the ANNIE experiments reflect the fact that the population of users was deliberately varied in these experiments. ANNIE experts had been using ANNIE for a period of about 6 months and many of them had considerable experience with other spoken dialogue systems. The ANNIE model reported in row 1 found that the percentage of requests for *help* was a significant predictor of user satisfaction, reflecting the fact that some users needed help significantly more than others, and these same users typically had much lower satisfaction with the system.

Finally, we examined how training set size affects performance by using cross-validation on the combined data from all three systems. We trained with 10 random samples of 27, 55, 110, 220, 330 and 490 dialogues and then tested on an unseen 10% of the dialogues. Figure 4 plots model fit (R^2) as a function of the amount of data used to train the model. The figure shows that the model fit is still improving with training size, which suggests that the predictive power of the models will increase when developed using larger corpora.

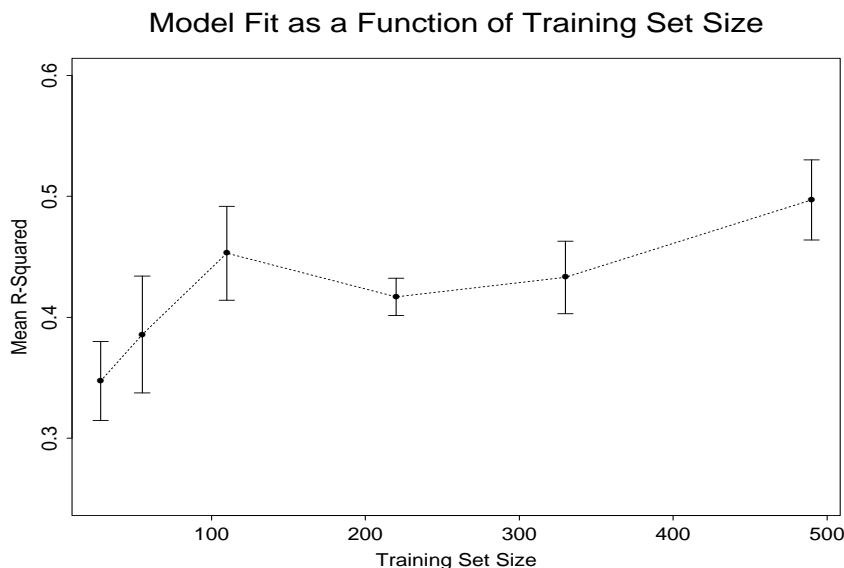


Fig. 4. Model fit as a function of training set size.

5 Discussion and Future Work

When we introduced the PARADISE framework for evaluating spoken language dialogue systems (SLDSs) we hoped to apply an empirical approach to developing general models of system usability. We proposed that user satisfaction could be used as a measure of system usability, and that it would be possible to develop models for predicting usability by learning which of range of other metrics were good predictors of user satisfaction.

This paper focused on developing general models for predicting user satisfaction based on experimental data collected with three different SLDS. We showed that models developed on 90% of the dialogues for each system does almost just as well predicting the variance in user satisfaction in the unseen 10% of the dialogues for that system as it does on the training set . We then showed that a general model, trained on 90% of the dialogues from all three systems, could also do as well on the test set as it did on the training set.

We also tested the extent to which a model trained on one system (ELVIS) can predict user satisfaction for the other two systems (ANNIE and TOOT). These tests showed that the ELVIS model could account for 55% of the variance in user satisfaction in TOOT and 36% of the variance in user satisfaction in ANNIE. The ELVIS model actually does better at predicting user satisfaction in TOOT than it does on on a random unseen 10% of ELVIS dialogues, and its predictive power on the ANNIE data is no worse than a model trained on 90% of the ANNIE dialogues. This suggests that there are general factors that predict user satisfaction that are shared by the models developed on each of the three systems, which is also supported by the factors as shown in Table 3.

We also examined the extent to which dialogues trained on one user population

generalize to dialogues collected with other user populations. We first tested models trained on dialogues with novices against dialogues with experts. The novice models did not predict user satisfaction in the expert population. This shows that it is important in future work to carry out longitudinal studies of users as they gain experience with the system, and that general models must be trained on a sample of dialogues including both expert and novice users. Finally, our results were also weaker when we tested models trained on dialogues with users who experienced average recognizer performance against dialogues with users who experienced excellent recognizer performance.

We believe that the results reported here make a contribution to the goal of developing general models of system usability, however there are several ways in which this work should be extended. First, we hope to repeat this analysis with a much larger corpus of dialogues collected via our participation in the DARPA COMMUNICATOR evaluations (Walker and Hirschman, 1999). Second, we hope to extend these techniques to dialogues collected in field studies as well as those collected in experiments such as we report here, by applying similar techniques to those reported in (Baggia et al., 1998). Finally, we believe that much research remains to be done to develop additional metrics for predicting user satisfaction that will generalize across systems.

References

- Baggia, P., Castagneri, G., and Danieli, M. (1998). Field trials of the italian arise train timetable system. In *Interactive Voice Technology for Telecommunications Applications, IVTTA*, pages 97–102.
- Bernsen, N. O., Dybkjaer, H., and Dybkjaer, L. (1996). Principles for the design of cooperative spoken human-machine dialogue. In *International Conference on Spoken Language Processing, ICSLP 96*, pages 729–732.
- Danieli, M. and Gerbino, E. (1995). Metrics for evaluating dialogue strategies in a spoken language system. In *Proceedings of the 1995 AAAI Spring Symposium on Empirical Methods in Discourse Interpretation and Generation*, pages 34–39.
- Hirschman, L., Dahl, D. A., McKay, D. P., Norton, L. M., and Linebarger, M. C. (1990). Beyond class A: A proposal for automatic evaluation of discourse. In *Proceedings of the Speech and Natural Language Workshop*, pages 109–113.
- Jack, M., Foster, J. C., and Stentiford, F. W. (1992). Intelligent dialogues in automated telephone services. In *International Conference on Spoken Language Processing, ICSLP*, pages 715 – 718.
- Kamm, C., Litman, D., and Walker, M. A. (1998). From novice to expert: The effect of tutorials on user expertise with spoken dialogue systems. In *Proceedings of the International Conference on Spoken Language Processing, ICSLP98*.
- Kamm, C., Narayanan, S., Dutton, D., and Ritenour, R. (1997). Evaluating spoken dialog systems for telecommunication services. In *5th European Conference on Speech Technology and Communication, EUROSPEECH 97*, pages 2203–2206.
- Litman, D. J. and Pan, S. (1999). Empirically Evaluating an Adaptable Spoken Dialogue System. In *Proceedings of the 7th International Conference on User Modeling*.
- Litman, D. J., Pan, S., and Walker, M. A. (1998). Evaluating Response Strategies in a Web-Based Spoken Dialogue Agent. In *Proceedings of ACL/COLING 98: 36th Annual Meeting of the Association of Computational Linguistics*, pages 780–787.
- Litman, D. J., Walker, M. A., and Kearns, M. J. (1999). Automatic detection of poor

- speech recognition at the dialogue level. In *Proceedings of the 37th Annual Meeting of the Association of Computational Linguistics, ACL99*, pages 309–316.
- Love, S., Dutton, R. T., Foster, J. C., Jack, M. A., and Stentiford, F. W. M. (1994). Identifying salient usability attributes for automated telephone services. In *International Conference on Spoken Language Processing, ICSLP*, pages 1307–1310.
- Pallett, D. S. (1985). Performance assessment of automatic speech recognizers. *J. Res. Natl. Bureau of Standards*, 90:371–387.
- Ralston, J. V., Pisoni, D. B., and Mullenix, J. W. (1995). Perception and comprehension of speech. In Syrdal, Bennet, and Greenspan, editors, *Applied Speech Technology*, pages 233–287. CRC Press.
- Shriberg, E., Wade, E., and Price, P. (1992). Human-machine problem solving using spoken language systems (SLS): Factors affecting performance and user satisfaction. In *Proceedings of the DARPA Speech and NL Workshop*, pages 49–54.
- Sparck-Jones, K. and Galliers, J. R. (1996). *Evaluating Natural Language Processing Systems*. Springer.
- Walker, M., , Fromer, J., Fabbrizio, G. D., Mestel, C., and Hindle, D. (1998a). What can I say: Evaluating a spoken language interface to email. In *Proceedings of the Conference on Computer Human Interaction (CHI 98)*.
- Walker, M., Hindle, D., Fromer, J., Fabbrizio, G. D., and Mestel, C. (1997a). Evaluating competing agent strategies for a voice email agent. In *Proceedings of the European Conference on Speech Communication and Technology, EUROSPEECH97*.
- Walker, M. and Hirschman, L. (September,1999). Darpa communicator evaluation proposal.
- Walker, M. A., Fromer, J. C., and Narayanan, S. (1998b). Learning optimal dialogue strategies: A case study of a spoken dialogue agent for email. In *Proceedings of the 36th Annual Meeting of the Association of Computational Linguistics, COLING/ACL 98*, pages 1345–1352.
- Walker, M. A., Litman, D., Kamm, C. A., and Abella, A. (1997b). PARADISE: A general framework for evaluating spoken dialogue agents. In *Proceedings of the 35th Annual Meeting of the Association of Computational Linguistics, ACL/EACL 97*, pages 271–280.