

Towards Multimodal Error Management: Experimental Evaluation of User Strategies in Event of Faulty Application Behavior in Automotive Environments

Gregor McGlaun, Frank Althoff, Manfred Lang, and Gerhard Rigoll

Institute for Human-Machine Communication
Technical University of Munich
Arcisstr. 16, 80290 Munich, Germany
phone: +49 89 289-28541

{mcglaun, althoff, lang, rigoll}@ei.tum.de

ABSTRACT

In this work, we present the results of a study analyzing the reactions of subjects on simulated errors of a dedicated in-car interface for controlling infotainment and communication services. The test persons could operate the system, using different input modalities, such as natural or command speech as well as head and hand gestures, or classical tactile paradigms. In various situational contexts, we scrutinized the interaction patterns the test participants applied to overcome different operation tasks. Moreover, we evaluated individual user behavior concerning modality transitions and individual fallback strategies in case of system errors. Two different error types (Hidden System Errors and Apparent System Errors) were provoked. As a result, we found out that initially, i.e. with the system working properly, most users prefer tactile or speech interaction. In case of Hidden System Errors, mostly changes from speech to tactile interaction and vice versa occurred. Concerning Apparent System Errors, 87% of the subjects automatically interrupted or cancelled their input procedure. 73% of all test persons who continued interaction, when the reason for the faulty system behavior was gone, strictly kept the selected modality. Regarding the given input vocabulary, none of the subjects selected head or hand gesture input as the leading fallback modality.

Keywords: Error, management, user, study, behavior, human machine interaction, multimodal, automotive;

1. INTRODUCTION

Today's growing complexity of in-car infotainment and communication systems strongly implicates an enlargement of input modalities in cars. Multimodal interfaces (MI) offer a lot of advantages to the driver. Compared to

monomodal systems, MIs allow for shorter learning phases and a highly intuitive and individual interaction[1]. Prior Studies of Oviatt et al. showed that in purely speech-based systems, the recognition rate dropped by 20-50%, when input was provided during natural or spontaneous interaction, by different user groups (e.g., accented speakers, speech impaired persons, or children), or in noisy mobile environments[2].

Error-prone situations are very likely to occur during interaction with various applications in a car environment. If caused by heavy traffic noise, the signal-to-noise ratio gets drastically worse, e.g., speech recognition will probably no longer work properly. Hence, multimodal interfaces have great potential for a significant enhancement of error robustness. Oviatt et al. mention that in dedicated scenarios, up to 86% of all task critical errors can be avoided, if an alternative input modality is provided [3]. A special set of multimodal systems facilitates user interaction in a *synergistical* [4] way, i.e. the user can enter input temporally overlapping in different modalities. Besides the gain of efficiency, in case of *redundant* [4] input, recognition errors of a single modality could directly be avoided by *mutual disambiguation*[5]. For example, if a speech recognizer issues an *n*-best list with low confidence for the potential output candidates, additional visual information by lip-reading can result in correct recovery of the input. On the other hand, the user can at any time choose freely amongst the provided modalities, which allows for a highly natural and intuitive way of human-machine communication. In case the selected modality channel fails for some reason, it is necessary to have a comprehensive *error management* that assists the user in performing the desired interaction (e.g., offering so-called *fallback modalities* dependent on the context of the application and the system environment). One step in a targeted development of an effective error handler is to evaluate how the multimodal interaction behavior of the user changes in case of system errors.

2. THEORETICAL BACKGROUND

In the field of error theories, many researchers have contributed significant work.

Strictly following an absolute philosophical point of view, Festinger[6] has developed an approach of cognitive dissonance for describing user errors. In his model, human error is always an expression of certain habits that cannot automatically be used in specific situations and thus result in an error during the operation.

Rigby[7] differentiates between sporadic, accidental, and systematic errors. In his phenomenological approach, sporadic errors are singular events, and are often considered as outliers. Accidental errors have a high mean variation with regard to the intended target status, but in contrast to systematic errors, they do not show any clear tendency towards a special direction.

However, these two approaches can hardly be used in a practical application since they suffer from a significant drawback. As the flow of interactions is assumed to be controlled by the *system* exclusively, the *user* is not involved sufficiently.

The theoretical basis for modeling potential error-prone user interaction has been given by Reason[8]. Related to the skill-rule-knowledge framework of Rasmussen[9], he differentiates between errors on three different performance levels. User interactions at the *skill-based level* comprise operations which have already become routine by multiple execution. Characteristic errors are either execution failures (slips) or failures of memory (lapses). They imply a deviation (normally known in advance) from a well-trained routine. At the *rule-based level*, human performance is determined by stored rules (productions). Hence, error patterns are planning failures (mistakes) and typically related to the misclassification of situations. At the *knowledge-based level*, in novel situations, problems are solved by applying conscious analytical processes and stored knowledge. Significantly, errors arise from unpredictable changes in the environment one is not prepared for.

Interaction Errors

Based on the formal description and abstract classification of human errors discussed above, we will derive a practicable definition of an interaction error that additionally covers system failures and faults. In the following, we briefly list some prototypical error-prone situations in human-machine communication.

In the first case, the user gives a command, which is interpreted by the system in a certain context that does not match the primary intention of the user. In a second scenario, a given command is interpreted in the wrong way and executed. This can be done by both sides, the system and the user. If the mental model of the user (which is a combination of the task model and the system model) and the user model of the system differ to a certain degree, an issued command will be interpreted in the wrong context.

The significance of the error potential becomes higher, the later the existing divergence of the two models is detected.

Covering these individual cases, we can give the following definition of an interaction error: *An error in human machine communication is a consequent result, if the requirements and the intention of the acting part are not covered in a sufficient way by the reacting part.* Thereby, the acting part can be both the system and the user.

Evaluating errors that appear during human-machine interaction, it is very important to distinguish whether the user or the system actually caused an error. This work exclusively focuses on scenarios, in which the user as a (correctly) acting part faces a certain malfunction of the system. In this regard, we can identify two different error classes.

Hidden System Errors (HSE) are spontaneous errors that occur independently from any contextual conditions situation (e.g., sudden break down of a module). For the user, it is not apparent or comprehensible why the error happened. This class of errors is characterized by partial or total recognition failures in the used input modality.

Moreover, we evaluated so-called *Apparent System Errors* (ASE). Hence, the cause that leads to the error is clearly evident to the user (e.g., the user interrupts interaction with the system due to an incoming call on her or his mobile phone).

3. EXPERIMENTAL SETUP

For a dedicated analysis of user strategies induced by the errors of the classes listed above, we designed a non-field user study. In the following subsections, we will describe the boundary conditions, the basic method, and the schedule of the test.

Test platform

The study was conducted in the car laboratory of the institute, which is specially adapted to evaluate multimodal user interfaces in automotive environments. In a separate control room, a test supervisor monitors the run of the experiment. For simulating realistic test conditions, the laboratory is equipped with a simple driving simulator consisting of a specially prepared BMW limousine with a force-feedback steering wheel, gas and break pedals, as well as an automatic transmission. The test subjects have to use these devices to control a 3D driving task, which is projected on a white wall in front of the car. This allows for experiencing the driving scenario from a natural in-car perspective and a better anticipation of the driving course. The individual parameters of the simulation can totally be controlled by a dedicated run chart, e.g., the degree of the curves, day- or night sight conditions, speed regulations, obstacles, or passing cars. To carry out reproducible test runs, we have developed a special software suit[10] enabling a precisely timed management of



Figure 1: Screenshot of the test interface used in the study

various system parameters, semi-automatically announcing the operation tasks at specified points of time and logging all kind of transactions. This concept has successfully been applied in various experiments[11].

For permanently recording audio and video signals, the car is equipped with a microphone array and two cameras. Together with the driving data of the simulation, an effective analysis of the individual interaction style of the participant could be carried out.

Test System

The test vehicle is equipped with a multimedia interface for controlling an infotainment and communication application consisting of an MP3-player, a radio tuner, and an integrated telephone application. As can be seen in figure 1, the interface itself is organized in four separated horizontal areas. The top line is composed of four buttons representing the individual *main modes* of the application (mp3, radio, telephone, and control). Directly beneath this button line, as the central design element, the interface provides a list containing individual items that can vertically be scrolled through by the two buttons on the right. The area in the in the lower part contains context specific buttons varying from five buttons in MP3 mode, three in radio and control, and two in telephone mode. In dependence of the current application mode, the system provides the particular functionality by displaying the respective buttons. In addition, the interface contains a feedback line continuously informing the user of the current volume and status of the interface, e.g., indicating an incoming call connection or additional information regarding the tuned radio station or the MP3 track that is currently played.

All devices of the application have the commonly known functionalities of a standard CD player (like play, skip, stop, etc.). In the radio mode, the test participant can tune to 25 different previously stored stations. The telephone functions are limited to basic call handling (call, accept, deny, etc.) of 30 predefined address-book entries. Moreover, the volume of the audio signal can be controlled in a separate mode.

Using a key-word (“computer”) for initialization the system can be operated via natural speech (SPC). Furthermore, subjects can use head- (HEG) and hand-gestures (HAG). For interaction via HEG or HAG, there is no ini-



Figure 2: Overview over the Wizard-of-Oz principle

tialization paradigm, but subjects are told to make sure that their head or the hand is not outside the focus of the camera. For tactile interaction, there was a 10” touch-screen (TSC) located in the middle of the center console, as well as an keypad that was integrated in the armrest (AKC) of the test car. The AKC consists of a 2x4 button array, which is organized in direct analogy to the positioning of the buttons on the touch-screen. The buttons of the first row allow for controlling the main modes, the buttons of the second row change their functionality in dependence of the current system mode. The two turning knobs are used for adjusting the volume and for browsing in the list display. By pressing these knobs, subjects can mute the volume and select the current list item, respectively.

The test persons were given a set of six head and 15 hand gestures, as well as 30 speech commands that could be provided in natural speech expressions. Concerning the composition of the interaction vocabulary, six commands (e.g., “yes” and “no”) could be entered in any modality channel.

Test Methodology

The study was performed as a partial Wizard-of-Oz (WOO) test[12]. In this evaluation, the test supervisor (also referred to as “wizard”) simulates the recognizers for the semantic higher-level modalities (head- and hand-gestures, as well as the speech recognizer). The wizard interprets the user's intention and generates the appropriate system commands, which are sent back to the interface in the car to trigger the intended functionality (see figure 2). Haptical interactions via TSC and AKC are directly transcribed by the system, but for simulating error scenarios, the wizard can also interfere with this process.

The wizard is instructed to be extremely cooperative. In case of ambiguous user actions or actions that are similar or synonymous to the given vocabulary set, the test supervisor tries to interpret the interaction at best in the current system context. We chose the WOO technique, as

it allows for a deterministic system behavior and an arbitrarily adjustable recognition rate at any time.

As presented in former work[13], the driving simulation demands each test subject in a different way. For normalizing the cognitive load induced by the driving task, we have developed a dedicated baseline technique. This method rates the driving performance of the subject in a separate test run, and consequently adjusts the degree of difficulty of the driving task in the actual trial.

Test Procedure

At the beginning of the test procedure, there is a short training period in which the subjects get to learn the different ways of interaction with the system. Before the main part starts, we carry out the baseline analysis to make sure that each subject is exposed to the same cognitive load, as outlined above. The main part of the trial is split up into three parts, as follows:

Part I: Reference phase: In this phase, which contained 16 different operation tasks, the user could arbitrarily select and combine the modalities. Regarding the given vocabulary set mentioned above, in ten of 16 tasks, the respective functionality of the interface can be addressed via any modality channel. The driving task comprises a simple course (straight road, no obstacles). In this phase, we determine individual modality preferences and the quota of synergistic multimodal input.

Part II: HSE scenarios: This step of the test consists of 21 tasks. In turn, the user can freely choose and combine all modalities. In five scenarios, the system does not react on any kind of user input (e.g., “Call Mr. Miller,” but the dial command does not work). As a feature of the HSE scenarios, the actual reason why the system does not react in the current situation is not evident at all for the user. To get comparable conditions, the course of the driving task in the second part is identical with that of the reference phase.

Part III: ASE scenarios: In this trial part, which comprised 21 tasks, eight ASE scenarios are interspersed. These error situations are simulated by dazzling lights, noise (e.g., braking sounds or honks) or incoming telephone calls interrupting the current action of the user. Moreover, in eight tasks the test subject is forced to take a certain initial modality. Using a more complex driving task (obstacles on the road, keeping speed limits) than in part I, we increase the workload of the test participants.

4. RESULTS

In the study, 15 subjects (47 % female, 53 % male, average age 25.5 a) participated.

Regarding the left columns for each modality in figure 3, it can be seen that in the reference phase, tactile interaction followed by speech were the leading modalities.

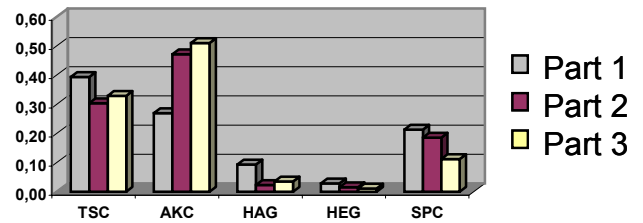


Figure 3: Modality distribution over all test parts (in %)

Very few used HAG (8%) or HEG (2%), respectively. 73% of the subjects stated that it was a new experience for them to operate a system via HAG or HEG and that it took getting used to. Despite massive system failures, the preference of the fallback modalities in part II and III of the test (middle and right column for each modality in figure 3) was nearly the same. SPC decreased, whereas AKC was even used more often than in the reference phase.

During the whole trial, we could only very sparsely observe synergistic multimodal input (8% of all interactions). Complementary input was mainly delivered sequentially or expressed in a single modality. When we asked test participants for the reason, they pointed out that while driving a car, they tried to keep an interaction as simple as practicable. Twelve of 15 subjects would rather execute two actions successively to reach a task goal, even if it eventually took more time.

In the HSE error scenarios, the subjects repeated a command 2.1 times on average, until they changed the modality. This is less than the subjects pointed out in a rating before the test (3.3 repetitions on average). Most retrials were done with tactile interaction via TSC (2.6 on average), AKC (2.4), and with speech (1.9). Contrary to the subjective data, the average number of command repetitions, using HAG (2.3) or HEG (2.2) was higher than AKC (1.7). If the system did not react for the third time, independently from the initial modality, the subjects used speech commands charged with emotions in combination with tactile interaction, i.e. hence, they performed redundant synergistic inputs. We could observe that towards the end of the test, the test persons showed a tendency to directly change the modality than to retry it in the current one. The subjects pointed out that they increasingly lost faith in the reliability of the modality and thus switched over to another one.

In the questionnaires, we also asked subjects to which fallback they would change, if they could no longer use their preferred modality. Concerning the situational context, we assumed a relaxed driving situation on an interstate. As a result, we got the transition matrix containing the averaged ratings (see table 1). Most test persons would prefer SPC, followed by tactile interaction. All participants dispreferred HEG and HAG. In the eyes of the subjects, some functionalities (e.g., a “random” or a “repeat” command) can hardly or only very laboriously be executed by gesturing. In the trial, 75% of the subjects

	TSC	SPC	HEG	HAG	AKC	median
TSC		2.93	2.33	2.53	2.53	2.58
SPC	1.13		1.93	1.60	1.67	1.58
HEG	2.60	2.86		3.00	3.00	2.86
HAG	3.33	3.40	3.15		3.36	3.31
AKC	2.27	2.60	2.60	2.07		2.38

Table 1: Transition matrix; first row: initial input modalities, first column: modalities the user tends to fall back to; rating was done, using a semantic differential scale without forced rating[12], where “1” stands for “definitely prefer”, and “6” means “definitely disprefer;”

switched from SPC to TSC. With AKC failing, only 33% of the subjects changed to TSC, whereas 47% switched over to SPC. In good agreement with the subjective ratings, HEG (0%) and HAG (6%) were hardly used as fall-back. Moreover, subjects tended to keep their modality as long as possible. In the ASE scenarios, 87% automatically interrupted the input, when an external event interfered with their action. In the experiment, 27% of the subjects forgot to finish the task they had begun. All of them pointed out that in such a case, they expect the system to remind them of the unfinished task in a way that they can proceed exactly from the point where they suspended. Those who continued interaction when the derangement was over, strictly kept the selected modality.

5. CONCLUSIONS AND FURTHER WORK

The study clearly proved that the situational context has to be considered in the error management. To be effective, the system must make a sensible taxonomy whether the current modality should be changed or the action can be retried in the initial modality.

In ongoing work, the findings are iteratively integrated into an error handling component of a multimodal in-car infotainment system[14]. The system is based on client server architecture, where information of the monomodal recognizers is processed via a late semantic fusion approach. To verify the usability of the error management component, extensive user studies are currently conducted with real recognizers for natural speech and gesture interaction[15,16,17].

6. ACKNOWLEDGMENTS

The work presented in this paper has been supported by the FERMUS-project[18], which is a cooperation between the BMW Group, SiemensVDO AG, Daimler-Chrysler AG, and the Institute of Human-Machine-Communication at the Technical University of Munich. FERMUS stands for error robust multimodal speech dialogues.

7. REFERENCES

- [1] S. Oviatt, “Taming Recognition Errors Within a Multimodal Interface,” in *Com. of the ACM*, 2000
- [2] S. Oviatt et al., “Designing the user interface for multimodal speech and gesture applications,” in *Human Computer Interaction, 2000, vol. 15, no. 4, pp. 263-322*, 2000
- [3] S. Oviatt et al. “Error resolution during human-computer error resolution,” in *Proc. of ICSLP '96*, 1996
- [4] L. Nigay et al., “A Generic Platform for Addressing the Multimodal Challenge,” in *Proc. of CHI '95*, 1995
- [5] S. Oviatt, “Mutual disambiguation of recognition errors in a multimodal architecture,” in *Proc. of CHI '99*, ACM Press, pp. 576-583, 1999
- [6] L. Festinger, “A Theory of Cognitive Dissonance,” *Stanford University Press*, 1957
- [7] L. Rigby, “The Nature of Human Error,” in *Annual technical conference transactions of the ASQC*, 1970
- [8] J. Reason, “Modeling the Basic Error Tendencies of Human Operators,” *Reliability Engineering and System Safety*, 22, pp. 137-153, 1988
- [9] J. Rasmussen, “Skills, Rules, Knowledge: Signals, Signs, and Symbols, and Other Distinctions in Human Performance Models,” in *Transac. of SMC '83*, SMC-13, pp. 257-267, 1983
- [10] B. Schuller et al., “Towards Automation of Usability Studies,” in *Proceedings of International Conference on Systems, Man and Cybernetics, SMC '02*, "Bridging the Digital Divide", Vol. 4, TP1N6, 2002
- [11] F. Althoff et al., “Experimental evaluation of user errors at the skill-based level in automotive environments,” in *Proc. of CHI '02*, 2002
- [12] J. Nielsen, “Usability Engineering,” Morgan Kaufmann Publishers, Inc., San Francisco, CA, 1999
- [13] G. McGlaun et al, “A new technique for adjusting distraction moments in multitasking non-field usability tests,” in *Proc. of CHI '02*, 2002
- [14] G. McGlaun et al, “A new approach for Integrating Multimodal Input via Late Semantic Fusion,” in *VDI/VDE-Proceedings 1678 of USEWARE 02*, Darmstadt, pp. , 2002
- [15] B. Schuller, “Towards intuitive speech interaction by the integration of emotional aspects,” in *Procs. of SMC '02*, "Bridging the Digital Divide", Vol. 6, WA2N1, 2002
- [16] M. Zobl et al., “Gesture-Based Control of In-Car Devices,” in *VDI/VDE-Proceedings 1678 of USEWARE 02*, Darmstadt, pp. 305-309, 2002
- [17] P. Morguet et al., “Comparison of Approaches to Continuous Hand Gesture Recognition for a Visual Dialog System,” in *Proceedings ICASSP 99, Vol. 6, pp. 3549-3552*, 1999
- [18] Project FERMUS: Error Robust Multimodal Speech Dialogs, website: www.fermus.de, 2002