

International Telecommunication Union

**ITU-T**

TELECOMMUNICATION  
STANDARDIZATION SECTOR  
OF ITU

**Series P**  
**Supplement 24**  
(10/2005)

SERIES P: TELEPHONE TRANSMISSION QUALITY,  
TELEPHONE INSTALLATIONS, LOCAL LINE  
NETWORKS

---

**Parameters describing the interaction with  
spoken dialogue systems**

ITU-T P-series Recommendations – Supplement 24



## ITU-T P-SERIES RECOMMENDATIONS

### TELEPHONE TRANSMISSION QUALITY, TELEPHONE INSTALLATIONS, LOCAL LINE NETWORKS

Vocabulary and effects of transmission parameters on customer opinion of transmission quality	Series	P.10
Subscribers' lines and sets	Series	P.30
		P.300
Transmission standards	Series	P.40
Objective measuring apparatus	Series	P.50
		P.500
Objective electro-acoustical measurements	Series	P.60
Measurements related to speech loudness	Series	P.70
Methods for objective and subjective assessment of quality	Series	P.80
		P.800
Audiovisual quality in multimedia services	Series	P.900
Transmission performance and QoS aspects of IP end-points	Series	P.1000

*For further details, please refer to the list of ITU-T Recommendations.*

## **Supplement 24 to ITU-T P-series Recommendations**

### **Parameters describing the interaction with spoken dialogue systems**

#### **Summary**

This Supplement provides definitions for a set of parameters which can be extracted from services which rely on spoken dialogue systems. The parameters can be extracted from logged (test) user interactions with the service under consideration. They quantify the flow of the interaction, the behaviour of the user and the system, and the performance of the speech technology devices involved in the interaction. They provide useful information for system development, optimization and maintenance, and are complementary to subjective quality judgments collected according to ITU-T Rec. P.851.

#### **Source**

Supplement 24 to ITU-T P-series Recommendations was agreed on 21 October 2005 by ITU-T Study Group 12 (2005-2008).

#### **Keywords**

Assessment, automatic speech recognition, automatic speech understanding, dialogue management, interaction parameter, speech generation, speech technology, spoken dialogue system.

## FOREWORD

The International Telecommunication Union (ITU) is the United Nations specialized agency in the field of telecommunications. The ITU Telecommunication Standardization Sector (ITU-T) is a permanent organ of ITU. ITU-T is responsible for studying technical, operating and tariff questions and issuing Recommendations on them with a view to standardizing telecommunications on a worldwide basis.

The World Telecommunication Standardization Assembly (WTSA), which meets every four years, establishes the topics for study by the ITU-T study groups which, in turn, produce Recommendations on these topics.

The approval of ITU-T Recommendations is covered by the procedure laid down in WTSA Resolution 1.

In some areas of information technology which fall within ITU-T's purview, the necessary standards are prepared on a collaborative basis with ISO and IEC.

## NOTE

In this publication, the expression "Administration" is used for conciseness to indicate both a telecommunication administration and a recognized operating agency.

Compliance with this publication is voluntary. However, the publication may contain certain mandatory provisions (to ensure e.g. interoperability or applicability) and compliance with the publication is achieved when all of these mandatory provisions are met. The words "shall" or some other obligatory language such as "must" and the negative equivalents are used to express requirements. The use of such words does not suggest that compliance with the publication is required of any party.

## INTELLECTUAL PROPERTY RIGHTS

ITU draws attention to the possibility that the practice or implementation of this publication may involve the use of a claimed Intellectual Property Right. ITU takes no position concerning the evidence, validity or applicability of claimed Intellectual Property Rights, whether asserted by ITU members or others outside of the publication development process.

As of the date of approval of this publication, ITU had not received notice of intellectual property, protected by patents, which may be required to implement this publication. However, implementors are cautioned that this may not represent the latest information and are therefore strongly urged to consult the TSB patent database.

© ITU 2005

All rights reserved. No part of this publication may be reproduced, by any means whatsoever, without the prior written permission of ITU.

## CONTENTS

	<b>Page</b>
1 Scope .....	1
2 References.....	1
3 Definitions .....	1
4 Abbreviations.....	2
5 Introduction .....	2
6 Characteristics of interaction parameters.....	3
7 Review of interaction parameters .....	4
8 Interpretation of interaction parameter values.....	14
BIBLIOGRAPHY .....	16



## Supplement 24 to ITU-T P-series Recommendations

### Parameters describing the interaction with spoken dialogue systems

#### 1 Scope

This Supplement describes parameters providing information on the interaction with services which are based on spoken dialogue systems, as seen by the system developer and service operator. Spoken dialogue systems addressed by this Supplement enable a spoken language interaction with a human user via the telephone network on a turn-by-turn basis, and have automatic speech recognition, speech understanding, dialogue management, response generation, and speech output capabilities. They may provide access to information stored in a database, or allow different types of transactions to be performed.

The parameters defined here quantify the flow of the interaction, the behaviour of the user and the system, and the performance of the speech technology devices involved in the interaction. For extracting all parameters, the spoken dialogue system has to be accessible as a glass box; still, some parameters may also be extracted in a black-box approach, i.e., without access to the individual system components. The extraction can partially be performed automatically, and partially relies on a human expert transcribing and annotating interaction log files. The parameters address system performance from a system developer's point-of-view; thus, they provide complementary information to subjective evaluation experiments with spoken dialogue systems for which recommendations are given in ITU-T Rec. P.851. Further guidance on subjective evaluation methods in general and on the assessment of speech output devices, is available in ITU-T Recs P.800 and P.85, and in the Handbook on Telephonometry. The parameters listed in this Supplement do not specifically refer to possible degradations introduced by the transmission channel. These effects are an item for further study by ITU-T SG 12.

#### 2 References

- ITU-T Recommendation P.85 (1994), *A method for subjective performance assessment of the quality of speech voice output devices*.
- ITU-T Recommendation P.800 (1996), *Methods for subjective determination of transmission quality*.
- ITU-T Recommendation P.851 (2003), *Subjective quality evaluation of telephone services based on spoken dialogue systems*.
- ITU-T *Handbook on Telephonometry* (1992).

#### 3 Definitions

For definitions not listed here, please refer to ITU-T Rec. P.10.

**3.1 barge-in:** The ability of a human to speak over a system prompt or system output [10].

**3.2 dialogue:** A conversation or an exchange of information. As an evaluation unit: One of several possible paths through the dialogue structure.

**3.3 efficiency:** Measures of the accuracy and completeness of system tasks relative to the resources (e.g., time, human effort) used to achieve the specific system tasks.

**3.4 exchange:** A pair of contiguous and related turns, one spoken by each party in the dialogue [8].

**3.5 functionality:** Capability of the system to provide functions which meet stated and implied needs when the system is used under specific conditions.

**3.6 meta-communication:** The communication about communication, e.g., for resolving misunderstandings ("Did I understand you right?") or for reaching agreement on the use of the language.

**3.7 performance:** The ability of a unit to provide the function it has been designed for.

**3.8 speech technology:** The discipline concerned with the research and development of spoken language input and output systems, using contributions from the neighbouring disciplines of acoustics, electrical engineering, statistics, phonetics, natural language processing, and involving system requirements specification, design, implementation and evaluation, corpus and linguistic resource processing, and consumer oriented product evaluation [10].

**3.9 spoken dialogue system:** A computer system with which human users interact via spoken language on a turn-by-turn basis.

**3.10 task:** All the activities which a user must develop in order to attain a fixed objective in some domain.

**3.11 task-oriented dialogue:** A dialogue concerning a specific subject, aiming at an explicit goal (such as resolving a problem or obtaining specific information) [8].

**3.12 transaction:** The part of a dialogue devoted to a single high-level task (e.g., making a travel booking or checking a bank account balance). A transaction may be coextensive with a dialogue, or a dialogue may consist of more than one transaction [8].

**3.13 turn:** Utterance. A stretch of speech, spoken by one party in a dialogue, from when this party starts speaking until another party definitely takes over [1].

**3.14 utterance:** See turn.

## 4 Abbreviations

ASR	Automatic Speech Recognition
AVM	Attribute-Value Matrix
AVP	Attribute-Value Pair
DARPA	Defense Advanced Research Projects Agency
DP	Dynamic Programming
DTMF	Dual Tone Multiple Frequency
IVR	Interactive Voice Response
MOS	Mean Opinion Score
SDS	Spoken Dialogue System
WoZ	Wizard-of-Oz

## 5 Introduction

Spoken dialogue systems (SDSs), i.e., computer systems with which human users interact via spoken language on a turn-by-turn basis, may be part of modern telephone networks. They enable access to databases and transactions via the telephone, e.g., for obtaining train or airline timetable information, stock exchange rates, tourist information, or to perform bank account operations or make hotel reservations. In contrast to simple interactive voice response (IVR) systems with DTMF input, SDSs offer the full range of speech interaction capabilities, including the recognition of user speech, the assignment of meaning to the recognized words, the decision on how to continue the

dialogue, the formulation of a linguistic response, and the generation of spoken output to the user. In this way, a more-or-less "natural" spoken interaction between user and system is enabled.

In order to evaluate the quality of services which rely on SDSs from a user's perspective, ITU-T SG 12 set up ITU-T Rec. P.851 in 2003. This Recommendation describes methods for conducting subjective evaluation experiments in order to determine quality *from a user's point-of-view*, taking the SDS as a black box. With the help of experiments carried out according to ITU-T Rec. P.851, valuable information on quality, as it is seen by the user, may be obtained. However, it may be difficult to determine how the individual system components contribute to the overall quality experienced by the user, e.g., to determine which component needs improvement in case of interaction problems. Thus, the evaluation should be complemented by information which address the system performance *from a system designer and service operator's point-of-view*.

System-related information may be described in terms of so-called *interaction parameters*. Such parameters help to quantify the flow of the interaction, the behaviour of the user and the system, and the performance of the speech technology devices involved in the interaction. They address system performance from a system developer and service operator's point-of-view, and thus provide complementary information to subjective evaluation data. For extracting some of the parameters, the spoken dialogue system has to be accessible as a glass box; other parameters may also be extracted in a black-box approach, i.e., without an access to the individual system components.

This Supplement provides a collection of interaction parameters which have been used for evaluating SDSs in the past 15 years. The listed parameters are related to the overall communication of information between user and system, the meta-communication in case of misunderstandings, the cooperativity of the system, the task which can be carried out with the help of the system, and the system's speech input capabilities. No parametric description is yet available for speech output quality (e.g., with respect to synthesized speech quality). The collection is based on theoretical work which is described in [17].

Not all of the interaction parameters are in a direct relationship to the perceived quality of SDS-based services. In fact, correlations between individual parameters and users' quality judgments are generally quite moderate. Still, it will be advantageous to dispose of a large set of parameters describing the interaction between user and system, in this way, capturing most of the information which is potentially relevant for perceived quality from a system designer's perspective. Such parameters provide useful information for system development, optimization, and maintenance.

The parameters having once being defined and applied in evaluation experiments at different test sites, may facilitate an estimation of their impact on perceived quality for a wide range of systems and services. In this way, it may become possible to develop algorithms for predicting quality on the basis of interaction parameters. Work in this direction is still under way within ITU-T SG 12 and elsewhere.

## **6 Characteristics of interaction parameters**

Interaction parameters can be extracted when real or test users interact with the service. The extraction can be performed partly instrumentally and partly with the help of log files which have to be transcribed and annotated by a human expert. Simple parameters, like the duration of the interaction or of single utterances, can usually be measured fully instrumentally, with appropriate algorithms. On the other hand, human transcription and annotation is necessary when not only the surface form (speech signals) is addressed, but also the contents and meaning of system or user utterances (e.g., to determine the accuracy of a word or concept).

SDSs are of such high complexity that a description of system behaviour and a comparison between systems or system versions needs to be based on a multitude of different parameters [24]. As a consequence, both (instrumental and expert-based) ways of collecting interaction parameters should

be followed in order to get as much information as possible. Based on the collected information, spoken dialogue services can be optimized and maintained very efficiently.

Because interaction parameters are based on data which has been collected in an interaction between user and system, they are influenced by the characteristics of the system, of the user, and of the interaction between both. These influences can usually not be separated, because the user's behaviour is strongly influenced by that of the system (e.g., the questions asked by the system), and vice versa (e.g., the vocabulary and speaking style of the user influences the system's accuracy of recognition and understanding). Consequently, interaction parameters strongly reflect the characteristics of the user group they have been collected with.

Interaction parameters are either determined in a laboratory test setting under controlled conditions, or in a field test. In the latter case, it may not be possible to extract all parameters, because not all necessary information can be gathered. For example, if the success of a task-oriented interaction (e.g., collection of a train timetable) is to be determined, then it is necessary to know about the exact aims of the user. Such information can only be collected in a laboratory setting, e.g., in the way it is described in ITU-T Rec. P.851. In case the fully integrated system is not yet available, it is possible to collect parameters from a so-called "Wizard-of-Oz" (WoZ) simulation, where a human experimenter replaces missing parts of the system under test. The characteristics of such a simulation have to be taken into account when interpreting the obtained parameters.

Interaction parameters can be calculated on a word level, on a sentence or utterance level, or on the level of a full interaction or dialogue. In case of word or utterance level parameters, average values are often calculated for each dialogue. The parameters collected with a specific group of users may be analysed with respect to the impact of the system (version), the user group, and the experimental setting (scenarios, test environment, etc.), using standard statistical methods. A characterization of these influences can be found in ITU-T Rec. P.851.

## **7 Review of interaction parameters**

Based on a broad literature survey, parameters were identified which have been used in different assessment and evaluation experiments during the past 15 years. The respective literature can be found in [2][3][4][6][7][8][9][11][12][14][16][21][22][23][24][25][26][27][28][30][31][32], and the parameters have been summarized in [17]. The parameters can broadly be classified as follows:

- Dialogue- and communication-related parameters;
- Meta-communication-related parameters;
- Cooperativity-related parameters;
- Task-related parameters;
- Speech-input-related parameters.

These categories will be briefly discussed in the following clauses. For each category, the respective parameters will be listed, together with a definition, the interaction level addressed by the parameter (word, utterance or dialogue), as well as the measurement method (instrumental or based on expert annotation).

## 7.1 Dialogue- and communication-related parameters

Parameters which refer to the overall dialogue and to the communication of information give a very rough indication of how the interaction takes place. They do not specify the communicative function of each individual utterance in detail. Parameters belonging to this category are listed in Table 1, and include duration-related parameters (overall dialogue duration, duration of system and user turns, system and user response delay), and word- and turn-related parameters (average number of system and user turns, average number of words per system and per user turn, number of system and user questions).

Two parameters which have been proposed in [11] are worth noting: The *query density* gives an indication of how efficiently a user can provide new information to a system, and the *concept efficiency* describes how efficiently the system can absorb this information from the user. These parameters also refer to the system's language understanding capability, but they have been included in this clause because they result from the system's interaction capabilities as a whole, and not purely from the language understanding capabilities.

All parameters in this category are of a global character and refer to the dialogue as a whole, although they are partly calculated on an utterance level. Global parameters are sometimes problematic, because individual differences in cognitive skills may be large in relation to the system-originated differences, and because subjects might learn strategies for task solution which have a significant impact on global parameters.

**Table 1 – Dialogue- and communication-related interaction parameters**

Abbr.	Name	Definition	Int. level	Meas. meth.
<i>DD</i>	dialogue duration	Overall duration of a dialogue in [ms], see e.g., [8][6][12][21].	dial.	instr.
<i>STD</i>	system turn duration	Average duration of a system turn, from the system starting speaking to the system stopping speaking, in [ms]. A turn is an utterance, i.e., a stretch of speech spoken by one party in the dialogue. [8]	utter.	instr.
<i>UTD</i>	user turn duration	Average duration of a user turn, from the user starting speaking to the user stopping speaking, in [ms]. [8]	utter.	instr.
<i>SRD</i>	system response delay	Average delay of a system response, from the user stopping speaking to the system starting speaking, in [ms]. [22]	utter.	instr.
<i>URD</i>	user response delay	Average delay of a user response, from the system stopping speaking to the user starting speaking, in [ms]. [22]	utter.	instr.
<i># turns</i>	number of turns	Overall number of turns uttered in a dialogue. [30]	dial.	instr./ expert.
<i># system turns</i>	number of system turns	Overall number of system turns uttered in a dialogue. [30]	dial.	instr./ expert.
<i># user turns</i>	number of user turns	Overall number of user turns uttered in a dialogue. [30]	dial.	instr./ expert.
<i>WPST</i>	words per system turn	Average number of words per system turn in a dialogue. [6]	utter.	instr./ expert.
<i>WPUT</i>	words per user turn	Average number of words per user turn in a dialogue. [6]	utter.	instr./ expert.

**Table 1 – Dialogue- and communication-related interaction parameters**

Abbr.	Name	Definition	Int. level	Meas. meth.
# <i>system questions</i>	number of system questions	Overall number of questions from the system per dialogue.	dial.	expert.
# <i>user questions</i>	number of user questions	Overall number of questions from the user per dialogue. [12][21]	dial.	expert.
<i>QD</i>	query density	<p>Average number of new concepts (slots, see 7.4) introduced per user query. Being <math>n_d</math> the number of dialogues, <math>n_q(i)</math> the total number of user queries in the <math>i^{\text{th}}</math> dialogue, and <math>n_u(i)</math> the number of unique concepts correctly "understood" by the system in the <math>i^{\text{th}}</math> dialogue, then</p> $QD = \frac{1}{n_d} \sum_{i=1}^{n_d} \frac{n_u(i)}{n_q(i)}$ <p>A concept is not counted to <math>n_u(i)</math> if the system already understood it in one of the previous utterances. [11]</p>	set of dial.	expert.
<i>CE</i>	concept efficiency	<p>Average number of turns which are necessary for each concept to be "understood" by the system. Being <math>n_d</math> the number of dialogues, <math>n_u(i)</math> the number of unique concepts correctly "understood" by the system in the <math>i^{\text{th}}</math> dialogue, and <math>n_c(i)</math> the total number of concepts in the <math>i^{\text{th}}</math> dialogue, then</p> $CE = \frac{1}{n_d} \sum_{i=1}^{N_d} \frac{n_u(i)}{n_c(i)}$ <p>A concept is counted whenever it was uttered by the user and was not already understood by the system. [11]</p>	set of dial.	expert.

## 7.2 Meta-communication-related parameters

Meta-communication, i.e., the communication about communication, is particularly important for the spoken interaction with systems which have limited recognition, understanding and reasoning capabilities. In this case, correction and clarification utterances or even sub-dialogues are needed to recover from misunderstandings.

The parameters belonging to this group quantify the number of system and user utterances which are part of meta-communication. Most of the parameters are calculated as the absolute number of utterances in a dialogue which relate to a specific interaction problem, and are then averaged over a set of dialogues. They include the number of help requests from the user, of time-out prompts from the system, of user utterances rejected by the system in cases where no semantic content could be extracted (ASR rejections), of diagnostic system error messages, of barge-in attempts from the user, and of user attempts to cancel a previous action.

The ability of the system (and of the user) to recover from interaction problems can be described in two ways: Either explicitly by the correction rate, i.e., the percentage of all (system or user) turns which are primarily concerned with rectifying an interaction problem, or implicitly with the *implicit*

*recovery* parameter, which quantifies the capacity of the system to regain utterances which have partially failed to be recognized or understood.

In contrast to the global measures, most meta-communication-related parameters describe the function of system and user utterances in the communication process. Thus, most parameters have to be determined with the help of an annotating expert. The parameters are listed in Table 2.

**Table 2 – Meta-communication-related interaction parameters**

<b>Abbr.</b>	<b>Name</b>	<b>Definition</b>	<b>Int. level</b>	<b>Meas. meth.</b>
<i># help request</i>	number of help requests from the user	Overall number of user help requests in a dialogue. A user help request is labelled by the annotation expert if the user explicitly asks for help. This request may be formulated as a question (e.g., "What are the available options?") or as a statement ("Give me the available options!"). [30]	utter.	expert.
<i># system help</i>	number of diagnostic system help messages	Overall number of help messages generated by the system in a dialogue. A help message is a system utterance which informs the user about available options at a certain point in the dialogue.	utter.	instr./expert.
<i># time-out</i>	number of time-out prompts	Overall number of time-out prompts, due to no response from the user, in a dialogue. [30]	utter.	instr.
<i># ASR rejection</i>	number of ASR rejections	Overall number of ASR rejections in a dialogue. An ASR rejection is defined as a system prompt indicating that the system was unable to "hear" or to "understand" the user, i.e., that the system was unable to extract any meaning from a user utterance. [30]	utter.	instr.
<i># system error</i>	number of diagnostic system error messages	Overall number of diagnostic error messages from the system in a dialogue. A diagnostic error message is defined as a system utterance in which the system indicates that it is unable to perform a certain task or to provide a certain information. [22]	utter.	instr./expert.
<i># barge-in</i>	number of user barge-in attempts	Overall number of user barge-in attempts in a dialogue. A user barge-in attempt is counted when the user intentionally addresses the system while the system is still speaking. In this definition, user utterances which are not intended to influence the course of the dialogue (laughing, expressions of anger or politeness) are not counted as barge-ins. [30]	utter.	expert.
<i># cancel</i>	number of user cancel attempts	Overall number of user cancel attempts in a dialogue. A user turn is classified as a cancel attempt if the user tries to restart the dialogue from the beginning, or if he/she explicitly wants to step one or several levels backwards in the dialogue hierarchy. [16][23]	utter.	expert.

**Table 2 – Meta-communication-related interaction parameters**

Abbr.	Name	Definition	Int. level	Meas. meth.
<i>SCT, SCR</i>	number of system correction turns, system correction rate	Overall number (SCT) or percentage (SCR) of all system turns in a dialogue which are primarily concerned with rectifying a "trouble", thus not contributing new propositional content and interrupting the dialogue flow. A "trouble" may be caused by speech recognition or understanding errors, or by illogical, contradictory, or undefined user utterances. In the case where the user does not give an answer to a system question, the corresponding system answer is labelled as a system correction turn, except when the user asks for information or action which is not supported by the current system functionality. [8][24][9][7]	utter.	expert.
<i>UCT, UCR</i>	number of user correction turns, user correction rate	Overall number (UCT) or percentage (UCR) of all user turns in a dialogue which are primarily concerned with rectifying a "trouble", thus not contributing new propositional content and interrupting the dialogue flow (see SCT, SCR). [8][24][9][7]	utter.	expert.
<i>IR</i>	implicit recovery	Capacity of the system to recover from user utterances for which the speech recognition or understanding process partly failed. Determined by labelling the partially parsed utterances (see definition of PA:PA in 7.5) as to whether the system response was "appropriate" or not: $IR = \frac{\# \text{utterances with appropriate system answer}}{PA:PA}$ For the definition of "appropriateness" see 7.3. [7]	utter.	expert.

### 7.3 Cooperativity-related parameters

Cooperativity has been identified as a key aspect for a successful interaction with a spoken dialogue system [1]. Unfortunately, it is difficult to quantify whether a system behaves cooperatively or not. Several of the dialogue- and meta-communication-related parameters somehow relate to system cooperativity, but they do not attempt to quantify this aspect.

Direct measures of cooperativity are the contextual appropriateness parameters introduced by Simpson and Fraser [24]. Each system utterance has to be judged by a number of experts as to whether it violates one or more of Grice's maxims for cooperativity, see [13]:

- *Quantity* of information: Make your contribution as informative as required (for the current purpose of the exchange); do not make your contribution more informative than is required.
- *Quality*: Try to make your contribution one that is true; do not say what you believe to be false; do not say that for which you lack adequate evidence.
- *Relation*: Be relevant.
- *Manner*: Be perspicuous; avoid obscurity of expression; avoid ambiguity; be brief (avoid unnecessary prolixity); be orderly.

These principles have been stated more precisely by Bernsen and Dybkjær [1] with respect to spoken dialogue systems.

The utterances are classified into the categories of appropriate (not violating Grice's maxims), inappropriate (violating one or more maxims), appropriate/inappropriate (the experts cannot reach agreement in their classification), incomprehensible (the content of the utterance cannot be discerned in the dialogue context), or total failure (no linguistic response from the system). It has to be noted that the classification is not always straightforward, and that interpretation principles may be necessary.

**Table 3 – Cooperativity-related interaction parameters**

Abbr.	Name	Definition	Int. level	Meas. meth.
<i>CA:AP</i> , <i>CA:IA</i> , <i>CA:TF</i> , <i>CA:IC</i> ,  <i>%CA:AP</i> , <i>%CA:IA</i> , <i>%CA:TF</i> , <i>%CA:IC</i>	contextual appropriateness	Overall number or percentage of system utterances which are judged to be appropriate in their immediate dialogue context. Determined by labelling utterances according to whether they violate one or more of Grice's maxims for cooperativity: <ul style="list-style-type: none"> <li>• <i>CA:AP</i>: Appropriate, not violating Grice's maxims, not unexpectedly conspicuous or marked in some way.</li> <li>• <i>CA:IA</i>: Inappropriate, violating one or more of Grice's maxims.</li> <li>• <i>CA:TF</i>: Total failure, no linguistic response.</li> <li>• <i>CA:IC</i>: Incomprehensible, content cannot be discerned by the annotation expert.</li> </ul> For more details see [24][8][9]; the classification is similar to the one adopted in [14].	utter.	expert.

#### 7.4 Task-related parameters

Current state-of-the-art services enable task-orientated interactions between system and user, and task success is a key issue for the usefulness of a service. Task success may best be determined in a laboratory situation where explicit tasks are given to the test subjects, see ITU-T Rec. P.851. However, realistic measures of task success have to take into account potential deviations from the scenario by the user, either because he/she did not pay attention to the instructions given in the scenario, because of his/her inattentiveness to the system utterances, or because the task was unresolvable and had to be modified in the course of the dialogue.

Modification of the experimental task is considered in most definitions of task success which are reported in the literature. Success may be reached by simply providing the right answer to the constraints set in the instructions, by constraint relaxation from the system or from the user (or both), or by spotting that no solution exists for the defined task. Task failure may be tentatively attributed to the system's or to the user's behaviour, the latter however being influenced by that of the system.

A different approach to determine task success is the  $\kappa$  coefficient. It assumes a speech-understanding approach which is based on attributes (concepts, slots) for which allowed values have to be assigned in the course of the dialogue between system and user. The pairs of attributes and assigned values are called attribute-value pairs (AVPs). A set of all available attributes, together with the values assigned by the task (a so-called attribute-value matrix (AVM)), completely describes a task which can be carried out with the help of the system. In order to determine the  $\kappa$  coefficient, a confusion matrix  $M(i,j)$  is set up for the attributes in the key (scenario definition) and in the reported solution (log file of the dialogue). Then, the agreement between key and solution  $P(A)$  and the chance agreement  $P(E)$  can be calculated from this matrix, see Table 4.  $M(i,j)$  can be

calculated for individual dialogues, or for a set of dialogues which belong to a specific system or system configuration.

The  $\kappa$  coefficient relies on the availability of a simple task coding scheme, namely in terms of an AVM. However, some tasks cannot be characterized as easily. In that case, more elaborate approaches to task success are needed, approaches which usually depend on the type of task under consideration.

**Table 4 – Task-related interaction parameters**

Abbr.	Name	Definition	Int. level	Meas. meth.
<i>TS</i>	task success	<p>Label of task success according to whether the user has reached his/her goal by the end of a dialogue, provided that this goal could be reached with the help of the system. The labels indicate whether the goal was reached or not, and the assumed source of problems:</p> <ul style="list-style-type: none"> <li>• <i>TS:S</i>: Succeeded (task for which solutions exist)</li> <li>• <i>TS:SCs</i>: Succeeded with constraint relaxation by the system</li> <li>• <i>TS:SCu</i>: Succeeded with constraint relaxation by the user</li> <li>• <i>TS:SCsCu</i>: Succeeded with constraint relaxation both from the system and from the user</li> <li>• <i>TS:SN</i>: Succeeded in spotting that no solution exists</li> <li>• <i>TS:F<sub>s</sub></i>: Failed because of the system's behaviour, due to system inadequacies</li> <li>• <i>TS:F<sub>u</sub></i>: Failed because of the user's behaviour, due to non-cooperative user behaviour</li> </ul> <p>See also [8][7][24].</p>	dial.	expert.
$\kappa$	kappa coefficient	<p>Percentage of task completion according to the kappa statistics. Determined on the basis of the correctness of the result AVM reached at the end of a dialogue with respect to the scenario (key) AVM. A confusion matrix <math>M(i,j)</math> is set up for the attributes in the result and in the key, with <math>T</math> the number of counts in <math>M</math>, and <math>t_i</math> the sum of counts in column <math>i</math> of <math>M</math>. Then</p> $\kappa = \frac{P(A) - P(E)}{1 - P(E)}$ <p>with <math>P(A)</math> the proportion of times that the AVM of the actual dialogue and the key agree, <math>P(A) = \sum_{i=1}^n \frac{M(i,i)}{T}</math>. <math>P(E)</math> can be estimated from the proportion of times that they are expected to agree by chance,</p> $P(E) = \sum_{i=1}^n \left(\frac{t_i}{T}\right)^2.$ <p>[31][4]</p>	dial. or set of dial.	expert.

## 7.5 Speech-input-related parameters

The speech input capability of a spoken dialogue system is determined by its capability to recognize words and utterances, and to extract the meaning from the recognized string (so-called "speech understanding"). For automatic speech recognition, two approaches have to be distinguished: Word recognizers are able to extract single words from the user's speech, when spoken in isolation (isolated word recognition) or continuously (keyword spotting). On the other hand, continuous speech recognizers are able to recognize whole sentences or utterances. Speech understanding is often performed on the basis of attribute-value pairs, see clause 7.4. The parameters described in the following paragraph address both speech recognition and speech understanding.

Continuous speech recognizers generally provide a word string hypothesis as an output. In order to judge whether the string correctly represents what has been said, a reference transcription has to be provided by the transcribing expert. For each utterance, hypothesized and reference string are first aligned on a word level, using a Dynamic Programming (DP) matching algorithm [19] [20]. On the basis of the alignment, the number of correctly determined words  $c_w$ , of substitutions  $s_w$ , of insertions  $i_w$ , and of deletions  $d_w$  is counted. These counts can be related to the total number of words in the reference  $n_w$ , resulting in two alternative measures of recognition performance, the word error rate *WER* and the word accuracy *WA*, see Table 5.

Complementary performance measures can be defined on the sentence level, in terms of a sentence accuracy, *SA*, or a sentence error rate, *SER*, see Table 5. In general, *SA* is lower than *WA*, because a single misrecognized word in a sentence impacts the *SA* parameter. It may however become higher than the word accuracy, especially when many single-word sentences are correctly recognized. The fact that *SER* and *SA* penalize a whole utterance when a single misrecognized word occurs has been pointed out by Strik et al. [26] [27]; the problem can be circumvented with the parameters *NES* and *WES*, see Table 5. When utterances are not separated into sentences, all sentence-related metrics can also be calculated on an utterance instead of a sentence level.

Isolated word recognizers provide an output hypothesis for each input word or utterance. Input and output words can be directly compared, and similar performance measures, as in the continuous recognition case, can be defined, omitting the insertions. Instead of the insertions, the number of false alarms in a time period can be counted, see van Leeuwen and Steeneken [28]. *WA* and *WER* can also be determined for keywords only when the recognizer operates in a keyword-spotting mode.

For speech understanding assessment, two common approaches have to be distinguished. The first one is based on the classification of system answers to user questions into categories of correctly answered, partially correctly answered, incorrectly answered, or failed answers. The individual answer categories can be combined into measures which have been used in the US DARPA program, see Table 5. The second way is to classify the system's parsing capabilities, either in terms of correctly parsed utterances, or of correctly identified AVPs. On the basis of the identified AVPs, global measures such as the concept accuracy, *CA*, the concept error rate, *CER*, or the understanding accuracy, *UA*, can be calculated. All parameters are listed in Table 5.

**Table 5 – Speech-input-related interaction parameters**

Abbr.	Name	Definition	Int. level	Meas. meth.
<i>WER, WA</i>	word error rate, word accuracy	<p>Percentage of words which have been correctly recognized, based on the orthographic form of the hypothesized and the (transcribed) reference utterance, and an alignment carried out with the help of the "scilite" algorithm, see [18]. Designating <math>n_w</math> the overall number of words from all user utterances of a dialogue, and <math>s_w</math>, <math>d_w</math> and <math>i_w</math> the number of substituted, deleted and inserted words, respectively, then the word error rate and word accuracy can be determined as follows:</p> $WER = \frac{s_w + i_w + d_w}{n_w}$ $WA = 1 - \frac{s_w + i_w + d_w}{n_w} = 1 - WER$ <p>See [24]; details on how these parameters can be calculated in case of isolated word recognition are given in [28].</p>	word	instr./expert
<i>SER, SA</i>	sentence error rate, sentence accuracy	<p>Percentage of entire sentences which have been correctly identified. Denoting <math>n_s</math> the total number of sentences, and <math>s_s</math>, <math>i_s</math> and <math>d_s</math> the number of substituted, inserted and deleted sentences, respectively, then:</p> $SER = \frac{s_s + i_s + d_s}{n_s}$ $SA = 1 - \frac{s_s + i_s + d_s}{n_s} = 1 - SER$ <p>[24]</p>	utter.	instr./expert.
<i>NES</i>	number of errors per sentence	<p>Average number of recognition errors in a sentence. Being <math>s_w(k)</math>, <math>i_w(k)</math> and <math>d_w(k)</math> the number of substituted, inserted and deleted words in sentence <math>k</math>, then</p> $NES(k) = s_w(k) + i_w(k) + d_w(k)$ <p>The average <i>NES</i> can be calculated as follows:</p> $NES = \frac{\sum_{k=1}^{\# \text{ user turns}} NES(k)}{\# \text{ user turns}} = \frac{WER \cdot \# \text{ user words}}{\# \text{ user turns}}$ <p>[26]</p>	utter.	instr./expert.

**Table 5 – Speech-input-related interaction parameters**

Abbr.	Name	Definition	Int. level	Meas. meth.
<i>WES</i>	word error per sentence	<p>Related to <i>NES</i>, but normalized to the number of words in sentence <i>k</i>, <math>w(k)</math>:</p> $WES(k) = \frac{NES(k)}{w(k)}$ <p>The average <i>WES</i> can be calculated as follows:</p> $WES = \frac{\sum_{k=1}^{\#user\ turns} WES(k)}{\#user\ turns}$ <p>[26]</p>	word	instr./expert.
<i>AN:CO</i> , <i>AN:IN</i> , <i>AN:PA</i> , <i>AN:FA</i> , <i>%AN:CO</i> , <i>%AN:IN</i> , <i>%AN:PA</i> , <i>%AN:FA</i>	number or percentage of correct/incorrect/partially correct/failed system answers	<p>Overall number or percentage of questions from the user which are:</p> <ul style="list-style-type: none"> <li>• correctly (<i>AN:CO</i>);</li> <li>• incorrectly (<i>AN:IC</i>);</li> <li>• partially correctly (<i>AN:PA</i>);</li> <li>• not at all (<i>AN:FA</i>).</li> </ul> <p>answered by the system, per dialogue, see [21][12][14].</p>	utter.	expert
<i>DARPA<sub>s</sub></i> , <i>DARPA<sub>me</sub></i>	DARPA score, DARPA modified error	<p>Measures according to the DARPA speech understanding initiative, modified by Skowronek [25] [17] to account for partially correct answers:</p> $DARPA_s = \frac{AN : CO - AN : IC}{\#user\ questions}$ $DARPA_{me} = \frac{AN : FA + 2 \cdot (AN : IC + AN : PA)}{\#user\ questions}$ <p>[21][12][25]</p>	utter.	expert.
<i>PA:CO</i> , <i>PA:PA</i> , <i>PA:IC</i> , <i>%PA:CO</i> , <i>%PA:PA</i> , <i>%PA:IC</i>	number of correctly/partially correctly/incorrectly parsed user utterances	<p>Evaluation of the number of concepts (attribute-value pairs, AVPs) in an utterance which have been extracted by the system:</p> <ul style="list-style-type: none"> <li>• <i>PA:CO</i>: All concepts of a user utterance have been correctly understood by the system.</li> <li>• <i>PA:PA</i>: Not all but at least one concept of a user utterance has been correctly understood by the system.</li> <li>• <i>PA:IC</i>: No concept of a user utterance has been correctly understood by the system.</li> </ul> <p>Expressed as the overall number or percentage of user utterances in a dialogue which have been parsed correctly/ partially correctly/ incorrectly. [7]</p>	utter.	expert.

**Table 5 – Speech-input-related interaction parameters**

Abbr.	Name	Definition	Int. level	Meas. meth.
<i>CA, CER</i>	concept accuracy, concept error rate	Percentage of correctly understood semantic units, per dialogue. Concepts are defined as attribute-value pairs (AVPs), with $n_{AVP}$ the total number of AVPs, and $s_{AVP}$ , $i_{AVP}$ and $d_{AVP}$ the number of substituted, inserted and deleted AVPs. The concept accuracy and the concept error rate can then be determined as follows:  $CA = 1 - \frac{s_{AVP} + i_{AVP} + d_{AVP}}{n_{AVP}}$ $CER = \frac{s_{AVP} + i_{AVP} + d_{AVP}}{n_{AVP}}$ [9][24][3][2]	utter.	expert.
<i>UA</i>	understanding accuracy	Percentage of user utterances in which all semantic units (AVPs) have been correctly extracted:  $UA = \frac{PA:CO}{\# \text{ user turns}}$ [32]	utter.	expert.

## 7.6 Further parameters

The majority of interaction parameters listed in the tables describe the behaviour of the system, which is obvious because it is the system and service quality which is of interest. In addition to these, user-related parameters can be defined. They are specific to the test user group, but may nevertheless be closely related to quality features perceived by the user.

When separating the quality of an SDS-based service into quality aspects, in the way which is indicated in 5.3/P.851, it can be observed that several aspects of quality are not addressed by interaction parameters. No parameters directly relate to usability, user satisfaction, acceptability, or speech output quality. So far, only very few approaches have been made which address the quality of speech output (be it concatenated or synthesized) in a parametric way. Instrumental measures related to speech intelligibility are defined e.g., in IEC 60268-16 [15], but they have not been designed for a telephone environment. Concatenation cost measures have been proposed which can be calculated from the input text and the speech database of a concatenative synthesis system [5]. Although they sometimes show high correlations to MOS scores obtained in auditory experiments, such measures are very specific to the speech synthesizer and its concatenation corpus.

## 8 Interpretation of interaction parameter values

Although interaction parameters, as the ones defined in this Supplement, are important for system design, optimization and maintenance, they are not directly linked to the quality which is perceived by the human user. Consequently, the collection of interaction parameters should be complemented by a collection of user judgements on different quality aspects, as described in ITU-T Rec. P.851. Only in this way can valid information on the quality of services, which are based on spoken dialogue systems, be obtained.

An interpretation of interaction parameter values may be based on experimental findings which are, however, often specific to the considered system or service. As an example, an increased number of time-out prompts may indicate that the user does not know what to say at specific points in a

dialogue, or that he/she is confused about system actions [29]. Increasing barge-in attempts may simply reflect that the user learned that it is possible to interrupt the system. In contrast, a reduced number may equally indicate that the user does not know what to say to the system. Lengthy user utterances may result from a large amount of initiative attributed to the user. A decrease of meta-communication-related parameter values (especially of user-initiated meta-communication) can be expected to increase system robustness, dialogue smoothness, and communication efficiency [1].

## BIBLIOGRAPHY

- [1] BERNSEN, N.O., DYBKJÆR, H., DYBKJÆR, L.: *Designing interactive speech systems: From first ideas to user testing*, Springer, DE-Berlin, 1998.
- [2] BILLI, R., CASTAGNERI, G., DANIELI, M.: Field trial evaluations of two different information inquiry systems, *Proc. 3rd IEEE Workshop on Interactive Voice Technology for Telecommunications Applications (IVTTA'96)*, US-Basking Ridge NJ, pp. 129-134, 1996.
- [3] BOROS, M., ECKERT, W., GALLWITZ, F., GORZ, G., HANRIEDER, G., NIEMANN, H.: Towards understanding spontaneous speech: Word accuracy vs. concept accuracy, *Proc. 4th Int. Conf. on Spoken Language Processing (ICSLP'96)*, IEEE, US-Piscataway NJ, 2, pp. 1009-1012, 1996.
- [4] CARLETTA, J.: Assessing agreement of classification tasks: The kappa statistics, *Computational Linguistics*, Vol. 22(2), pp. 249-254, 1996.
- [5] CHU, M., PENG, H.: An objective measure for estimating MOS of synthesized speech, *Proc. 7th Europ. Conf. on Speech Communication and Technology (Eurospeech 2001 – Scandinavia)*, DK-Aalborg, 3, pp. 2087-2090, 2001.
- [6] COOKSON, S.: Final evaluation of VODIS – Voice operated data inquiry system, *Proc. of Speech'88, 7th FASE Symposium*, UK-Edinburgh, 4, pp. 1311-1320, 1988.
- [7] DANIELI, M., GERBINO, E.: Metrics for evaluating dialogue strategies in a spoken language system, *Empirical Methods in Discourse Interpretation and Generation. Papers from the 1995 AAAI Symposium*, US-Stanford CA, AAAI Press, US-Menlo Park CA, pp. 34-39, 1995.
- [8] FRASER, N.: Assessment of interactive systems, *Handbook on Standards and Resources for Spoken Language Systems* (D. Gibbon, R. Moore and R. Winski, eds.), Mouton de Gruyter, DE-Berlin, pp. 564-615, 1997.
- [9] GERBINO, E., BAGGIA, P., CIARAMELLA, A., RULLENT, C.: Test and evaluation of a spoken dialogue system, *Proc. Int. Conf. Acoustics Speech and Signal Processing (ICASSP'93)*, IEEE, US-Piscataway NJ, 2, pp. 135-138, 1993.
- [10] GIBBON, D., MOORE, R., WINSKI, R., Eds.: *Handbook on Standards and Resources for Spoken Language Systems*. Mouton de Gruyter, DE-Berlin, 2000.
- [11] GLASS, J., POLIFRONI, J., SENEFF, S., ZUE, V.: Data collection and performance evaluation of spoken dialogue systems: The MIT experience, *Proc. 6th Int. Conf. on Spoken Language Processing (ICSLP 2000)*, CN-Beijing, 4, pp. 1-4., 2000.
- [12] GOODINE, D., HIRSCHMAN, L., POLIFRONI, J., SENEFF, S., ZUE, V.: Evaluating interactive spoken language systems, *Proc. 2nd Int. Conf. on Spoken Language Processing (ICSLP'92)*, CA-Banff, 1, pp. 201-204, 1992.
- [13] GRICE, H.P.: Logic and conversation, *Syntax and Semantics, Vol. 3: Speech Acts* (P. Cole and J.L. Morgan, eds.), Academic Press, US-New York NY, pp. 41-58, 1975.
- [14] HIRSCHMAN, L., PAO, C.: The cost of errors in a spoken language system, *Proc. 3rd Europ. Conf. on Speech Communication and Technology (Eurospeech'93)*, DE-Berlin, 2, pp. 1419-1422, 1993.
- [15] IEC 60268-16 (2003), *Sound system equipment – Part 16: Objective rating of speech intelligibility by speech transmission index*. International Electrotechnical Commission, CH-Geneva.

- [16] KAMM, C.A., LITMAN, D.J., WALKER, M.A.: From novice to expert: The effect of tutorials on user expertise with spoken dialogue systems, *Proc. 5th Int. Conf. on Spoken Language Processing (ICSLP'98)*, AU-Sydney, 4, pp. 1211-1214, 1998.
- [17] MÖLLER, S.: *Quality of telephone-based spoken dialogue systems*. Springer, US-New York NY, 2005.
- [18] NIST Speech Recognition Scoring Toolkit, *Speech recognition scoring toolkit*, National Institute of Standards and technology, <http://www.nist.gov/speech/tools>, US-Gaithersburg MD, 2001.
- [19] PICONE, J., DODDINGTON, G.R., PALLETT, D.S.: Phone-mediated word alignment for speech recognition evaluation, *IEEE Trans. Acoustics Speech and Signal Processing*, Vol. 38(3), pp. 559-562, 1990.
- [20] PICONE, J., GOUDIE-MARSHALL, K.M., DODDINGTON, G.R., FISHER, W.: Automatic text alignment for speech system evaluation, *IEEE Trans. Acoustics Speech and Signal Processing*, Vol. 34(4), pp. 780-784, 1986.
- [21] POLIFRONI, J., HIRSCHMAN, L., SENEFF, S., ZUE, V.: Experiments in evaluating interactive spoken language systems, *Proc. DARPA Speech and Natural Language Workshop*, US-Harriman CA, pp. 28-33, 1992.
- [22] PRICE, P.J., HIRSCHMAN, L., SHRIBERG, E., WADE, E.: Subject-based evaluation measures for interactive spoken language systems, *Proc. DARPA Speech and Natural Language Workshop*, US-Harriman CA, pp. 34-39, 1992.
- [23] SAN-SEGUNDO, R., MONTERO, J.M., COLÁS, J., GUTIÉRREZ, J., RAMOS, J.M., PARDO, J.M.: Methodology for dialogue design in telephone-based spoken dialogue systems: A Spanish train information system, *Proc. 7th Europ. Conf. on Speech Communication and Technology (Eurospeech 2001 – Scandinavia)*, DK-Aalborg, 3, pp. 2165-2168, 2001.
- [24] SIMPSON, A., FRASER, N.M.: Black box and glass box evaluation of the SUNDIAL system, *Proc. 3rd Europ. Conf. on Speech Communication and Technology (Eurospeech'93)*, DE-Berlin, 2, pp. 1423-1426, 1993.
- [25] SKOWRONEK, J.: *Entwicklung von Modellierungsansätzen zur Vorhersage der Dienstqualität bei der Interaktion mit einem natürlichsprachlichen Dialogsystem*. Diploma thesis (unpublished), Institut für Kommunikationsakustik, Ruhr-Universität, DE-Bochum, 2002.
- [26] STRIK, H., CUCCHIARINI, C., KESSENS, J.M.: Comparing the performance of two CSRs: How to determine the significance level of the differences, *Proc. 7th Europ. Conf. on Speech Communication and Technology (Eurospeech 2001 – Scandinavia)*, DK-Aalborg, 3, pp. 2091-2094, 2001.
- [27] STRIK, H., CUCCHIARINI, C., KESSENS, J.M.: Comparing the recognition performance of CSRs: In search of an adequate metric and statistical significance test, *Proc. 6th Int. Conf. on Spoken Language Processing (ICSLP 2000)*, CN-Beijing, 4, pp. 740-743, 2000.
- [28] VAN LEEUWEN, D., STEENEKEN, H.: Assessment of recognition systems, *Handbook on Standards and Resources for Spoken Language Systems* (D. Gibbon, R. Moore and R. Winski, eds.), Mouton de Gruyter, DE-Berlin, pp. 381-407, 1997.
- [29] WALKER, M.A., FROMER, J., DI FABBRIZIO, G., MESTEL, C., HINDLE, D.: What can I say?: Evaluating a spoken language interface to email, *Human Factors in Computing Systems. CHI'98 Conference Proc.*, US-Los Angeles CA, ACM, US-New York NY, pp. 582-589, 1998.

- [30] WALKER, M.A., LITMAN, D.J., KAMM, C.A., ABELLA, A.: Evaluating spoken dialogue agents with PARADISE: Two case studies, *Computer Speech and Language*, Vol. 12(3), pp. 317-347, 1998.
- [31] WALKER, M.A., LITMAN, D.J., KAMM, C.A., ABELLA, A.: PARADISE: A framework for evaluating spoken dialogue agents, *Proc. of the 35th Ann. Meeting of the Assoc. for Computational Linguistics*, ES-Madrid, pp. 271-280, 1997.
- [32] ZUE, V., SENEFF, S., GLASS, J.R., POLIFRONI, J., PAO, C., HAZEN, T.J., HETHERINGTON, L.: JUPITER: A telephone-based conversational interface for weather information, *IEEE Trans. Speech and Audio Processing*, Vol. 8(1), pp. 85-96, 2000.



## SERIES OF ITU-T RECOMMENDATIONS

Series A	Organization of the work of ITU-T
Series D	General tariff principles
Series E	Overall network operation, telephone service, service operation and human factors
Series F	Non-telephone telecommunication services
Series G	Transmission systems and media, digital systems and networks
Series H	Audiovisual and multimedia systems
Series I	Integrated services digital network
Series J	Cable networks and transmission of television, sound programme and other multimedia signals
Series K	Protection against interference
Series L	Construction, installation and protection of cables and other elements of outside plant
Series M	Telecommunication management, including TMN and network maintenance
Series N	Maintenance: international sound programme and television transmission circuits
Series O	Specifications of measuring equipment
<b>Series P</b>	<b>Telephone transmission quality, telephone installations, local line networks</b>
Series Q	Switching and signalling
Series R	Telegraph transmission
Series S	Telegraph services terminal equipment
Series T	Terminals for telematic services
Series U	Telegraph switching
Series V	Data communication over the telephone network
Series X	Data networks, open system communications and security
Series Y	Global information infrastructure, Internet protocol aspects and next-generation networks
Series Z	Languages and general software aspects for telecommunication systems