



STANDARDISED SPEECH INTERFACES - KEY FOR OBJECTIVE EVALUATION OF RECOGNITION ACCURACY

Christel Müller and Karsten Schröder

Deutsche Telekom Berkomp GmbH, Research Section Speech Systems, Computer Telephony Systems
10589 Berlin, Goslarer Ufer 35, Germany
c.mueller@berkomp.de, k.schroeder@berkomp.de

ABSTRACT

Today existing speech recognisers differ not only in their technology from single word to natural language understanding, from hardware based and software-only solutions but also in the recognition accuracy under real-time conditions in the whole variety of environments and networks. Standardised APIs for a defined state diagram of recognition process are the key requirement for integrating different technologies and their evaluation. First time evaluation tests based on standardised interfaces and optimised vocabularies were carried out in an open system environment at Deutsche Telekom's Research Lab. The approach of a unique lexicon development system enhances the whole process of vocabulary generation. Optimising the phoneme transcription is one of the main future issues at the lab.

1. INTRODUCTION

Existing application development tools differ significantly in the integration of ASR recognition technologies, in the handling of vocabularies and grammars for the same dialogue state. The main conditions for getting independent and comparable results of recognition accuracy are on a unique set of speech data one hand and a unique set of ASR functions at the resource layer on the other hand. Furthermore the development of vocabulary, grammar rules and parameter settings have to correspond for one selected test application. Therefore a recogniser evaluation system configuration for unified interfaces of different ASR technologies was set up and several tests for key word spotter were carried out at the Deutsche Telekom research lab. Additionally the differences in phonetic transcriptions of vocabularies led to first conclusions for an optimised deployment of ASR technology of German Language. Finally the characteristics of speech data were described for this research approach of testing recognition accuracy under different conditions in telecommunication networks. It was important to keep the balance between specific features of speech interfaces and restrictions of voice processing technology for the test application, since both

factors influence and enhance each other in the process of technological advance.

2. STANDARDISED SPEECH INTERFACES

The ASR resource management layer has the task to implement unified ASR interfaces and accordingly to achieve vendor independence. This includes the recognition functions as well as the vocabulary, grammar and parameter settings for hardware and software based solutions. The standardisation of ASR testing was carried out according to the S.100 and S.300 Interfaces, which were created by the members of ASR Task Group of Enterprise Computer Telephony Forum. These interfaces correspond to the existing SAPI and JSRAPI standards in terms of main functionality.

2.1. Subset of Functionality for ASR Testing

The independence of basic recognition functions and a selected set of parameters [1] are realised for three selected technologies:

- ◆ HMM-Keyword Spotter RA, hardware based, multilingual
- ◆ HMM-Keyword Spotter RB, hardware based, multilingual
- ◆ HMM-Keyword Spotter RC, software based, multilingual. The handling parameters not only includes the initial settings but also the modification of their returned values and a repetitive setting:

```
„Recognize(Group,nNumResults,RTC,  
OptArgs,TranInfo,Mode)“  
„RetrieveRecognition(Group,ResultType,  
Results,TranInfo,Mode)“.
```

The 'Talk over' functionality was not taken into account for the different experiments.

2.2. Implementation into an open System Architecture

ASR evaluation in an open environment is a completely new approach to obtain true evaluation results under real conditions. An open system is characterised by: variable choice of ASR technology provider, multiple applications running at this environment, standardised interfaces as well as administration and minimal effort of resource implementation.

Such an open system architecture consists of :

- ◆ User Interface Layer
- ◆ Network Layer
- ◆ Resource Layer
- ◆ Resource Management Layer and
- ◆ Application Layer,

where the resource layer contains the different ASR technologies and ensures the communication via a standardised speech interface.

3. DESIGN OF THE ASR EVALUATION SYSTEM

The approach covered ‘black box’ and ‘glass box’ aspects as well. Thus, the term evaluation system was used according to the definitions of assessment and evaluation given in [2].

The fully automated system for the speech recogniser assessment was built of the following components (s. figure 1):

- ◆ Test application A, B
- ◆ Service Creation Environment with a graphical user interface,
- ◆ Computer telephony middleware for resource management of ASR and telephony and media processing,
- ◆ S.300 drivers for ASR and media control
- ◆ media control unit
- ◆ ASR software and hardware
- ◆ Primary Rate Interface (PRI) running in network mode
- ◆ Primary Rate Interface running in terminal mode
- ◆ PRI back-to-back cable.

Using the back-to-back dial mode via the digital telephone interfaces the timeslots on the SCBus for the speech recognition and for the playback of the speech samples were connected.

A special feature of the CT-software allows to start a play action and a recognition action in parallel within a single application. Thus, there was no need to synchronize two sites as shown in [3] and [4].

The speech coding format was PCM A-Law, 64 kBit/s. No conversions took place on the signal way in the system.

All recognisers got the identical digital speech input over the same interface (SCBus).

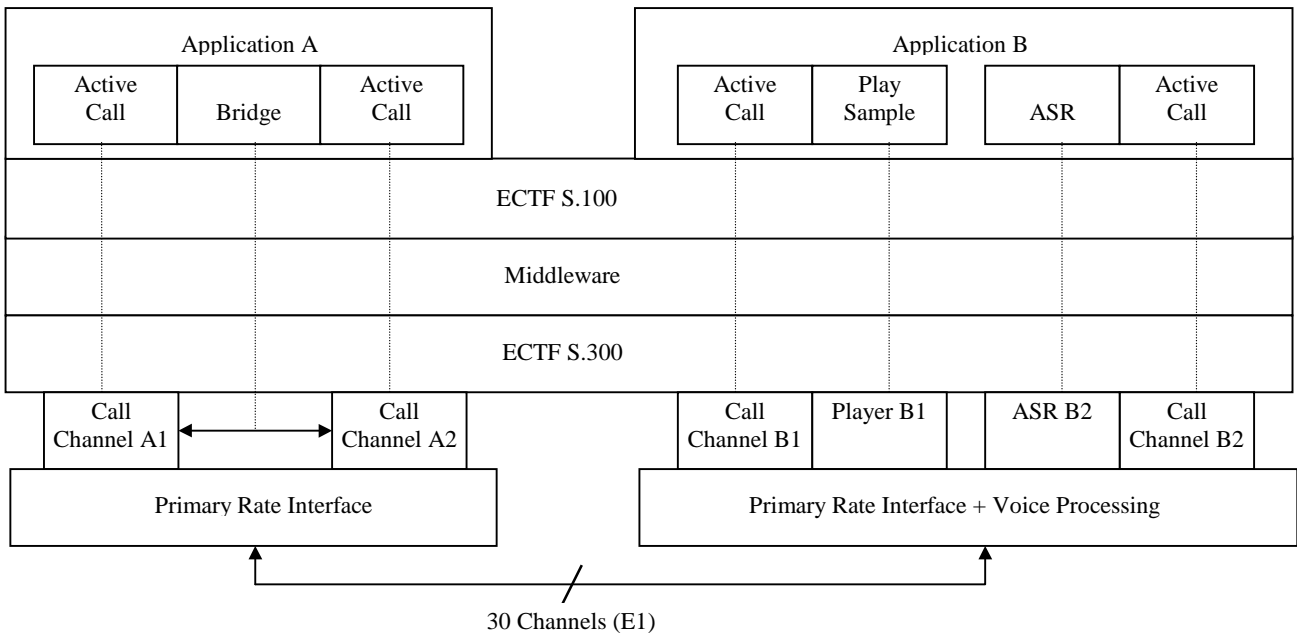


Figure 1: System Architecture of ASR Evaluation

4. EXAMPLE OF A COMPARATIVE ASR EVALUATION TASK

All three technologies were tested under the same environment with the same application.

Two test dialogues were started. After receiving and answering the call of dialogue 1 the second application played the test set of words and recognized the speech patterns in parallel.

4.1. Characteristics of Speech Data for objective ASR evaluation

Deutsche Telekom created a pool of speakers which covered all German dialects. The following information were available:

- ◆ gender,
- ◆ age,
- ◆ primary and secondary language
- ◆ primary and secondary dialect.

Speech samples read from word lists were collected for large telecommunication applications over the phone.

The features of every session were registered concerning:

- a) the telephones characteristics:
 - ◆ analogue standard and wireless,
 - ◆ ISDN standard and wireless / DECT,
 - ◆ analogue mobile / TACS,
 - ◆ digital mobile / GSM,
 - ◆ public telephone,
- b) the exact name of the telephone and speakers' environments:
 - ◆ car,
 - ◆ home,
 - ◆ office,
 - ◆ landscape,
 - ◆ street,
 - ◆ telephone booth,
 - ◆ public building and other.

Every speech sample was recorded into one separate speech file, 64 kbit/s PCM format, A-law without any conversions directly from the ISDN. After the validation process these recordings included defined periods of leading and trailing silence (approximately 200 ms). Certain subjective criteria concerning the quality were registered, e.g. kind of noise, recording level and intelligibility.

All features regarding the speakers, telephones, environments and subjective evaluations on the quality of the speech samples enable in-depth analyses of the recogniser's accuracy.

4.2. Design and Transcription of ASR Vocabulary

Design of ASR vocabularies in telecommunication applications depends on the possibilities of implementing synonyms for mission critical words or expressions.

The constant part normally consist of digits and names while the variety of word design only exists for dialogue control and interaction.

For the creation of a unique testing vocabulary a set of critical words was selected (s. *table 1*).

The phonetic transcriptions were processed automatically by the lexicon software. Differences between the transcriptions are shown in bold characters.

Characters representing silence or tone group boundaries were left out at the beginning and the end of the transcriptions strings.

If the specific lexicon software used product specific characters or IPA characters for the transcription, a translation to SAMPA was done according to [4].

Differences within the transcription alphabets existed as follows:

- ◆ Syllable boundary \$ (RC),
- ◆ Separator - (RC),
- ◆ Tone group boundary | (RB),
- ◆ Primary stress “ (RB, RC)
- ◆ glottal stop ? (RB, RC)
- ◆ near-open central unrounded /r/ **6** (RB)
- ◆ Syllabic /n/ =**n** (RB)
- ◆ trill/fricative **R** (RC, RB).

4.3. Accuracy Results

The tests of accuracy were carried out by 4050 German speech samples without using out-of-vocabulary data.

Automatically generated recogniser results are shown in *table 1*.

According to the provided tests of three different ASR technologies we concentrated on the differences in transcription of their associated vocabularies.

It turned out that the occurrence of phonetic specifics did not correlate to the deviations in recognition accuracy results.

Word	Recogniser A		Recogniser B		Recogniser C	
	Transcription	Wacc	Transcription	Wacc	Transcription	Wacc
Amerika	ame:ri:ka:	99,33%	?am“e: ri: ka:	97,33%	a\$“me:\$rI\$:ka	98,67%
Australien	aUstra:li:@n	100,00%	?aU stRa: li@n	100,00%	aUs\$“tra:\$lj@n	100,00%
Belgien	bElgi:@n	99,33%	b“El gi: @n	96,00%	“bEl\$gi:\$@n	97,33%
Brasilien	brazi:li:@n	100,00%	bRa zi: li@n	98,67%	bra\$“zi:l\$ j@n	98,00%
Dänemark	dE:n@mark	100,00%	d“E: n@ mark	96,00%	“de:\$n@\$maRk	100,00%
Deutschland	dOYtSlant	100,00%	d“OYtS lant	100,00%	“dOYtS\$ lant	100,00%
England	ENlant	85,33%	?“EN lant	83,33%	“EN\$ lant	93,33%
Finnland	fInlant	100,00%	f“In lant	95,33%	“fIn\$ lant	91,33%
Frankreich	fraNkraIC	100,00%	f“aNk raIC	98,67%	“fraNk\$raIC	93,33%
Griechenland	gri:C@n lant	100,00%	gr“i: C@n lant	96,00%	“gri:\$C@n\$ lant	100,00%
Großbritannien	grOsbritani:@n	100,00%	gRo:s bRi ta: ni@n	98,67%	gro:s\$brI\$“tan\$ j@n	99,33%
Holland	hOlant	100,00%	h“O lant	98,00%	“hOI\$ lant	99,33%
Irland	Irlant	95,33%	?I6 lant	64,67%	“IR\$ lant	72,67%
Italien	i:ta:li:@n	99,33%	?i: “a: li: @n	98,67%	i:\$“ta:\$ j@n	98,67%
Japan	ja:pan	100,00%	j“a: pan	100,00%	“ja:\$pa:n	98,67%
Kanada	kanada:	100,00%	k“a na da:	98,67%	“ka\$na\$da	99,33%
Luxemburg	lUksEmbUrk	100,00%	lU ks@m bU6k	99,33%	“lUk\$S@m\$bURk	100,00%
Niederlande	ni:d@rland@	100,00%	ni: d@6 lan d@	99,33%	“ni:\$d@R\$lan\$d@	100,00%
Norwegen	nOrve:g@n	100,00%	n“Or ve: g=n	96,67%	“nOR\$ve:\$g@n	99,33%
Österreich	2:st@raIC	96,67%	?“2:s t@ raIC	96,67%	“2:s\$t@RraIC	99,33%
Portugal	pOrtu:gal	98,67%	p“Or tu: gal	95,33%	“pOR\$tu:\$gal	99,33%
Rußland	rUslant	99,33%	r“Us lant	92,00%	“rUs\$ lant	89,33%
Schweden	Sve:d@n	100,00%	Sv“e: d=n	99,33%	“Sve:\$d@n	100,00%
Schweiz	Svalts	98,67%	Sv“alts	98,67%	“SvaItS	98,67%
Spanien	Spa:ni:@n	100,00%	Sp“a: ni: @n	100,00%	“Spa:n\$ j@n	100,00%
USA	u:Esa:	98,67%	?“u: ?“Es ?“a:	94,00%	u:\$Es\$?a:	100,00%
Vereinigte Staaten	fEraInIkt@sta:t@n	99,33%	fE6 ?a nIC t@ St“a: t=n	97,33%	fEr\$“a nIC\$t@-“Sta:t@n	99,33%

Table 1: Phonetic Transcription of critical Words for a special Test Set vs. Word Accuracy

5. CONCLUSIONS

The research work at Deutsche Telekom Speech Systems, Computer Telephony Systems section, showed that the unified ASR interface with standardised functionality remains the only way of getting reasonable results in ASR assessment and evaluation of different technologies.

Due to the fact that no direct relation between the phonetic specifics and the deviations in accuracy could be found further experiments based on vocabularies transcribed in a unique manner will be provided. If for each recogniser these new test results will not differ significantly from the accuracy shown in *table 1* then there is no need to insist on technology provider dependent transcription schemes.

Deploying different recognition technologies for one application would only require one single vocabulary creation environment.

Today's ASR real-time applications only meet the user acceptance, if the described well defined vocabularies for voice user interfaces are integrated.

6. REFERENCES

- [1] Enterprise Computer Telephony Forum (1996), S.100 Revision 1.0, Media Services "C" Language, Application Programming Interfaces.
- [2] David S. Pallett, Adrian Fourcin (1995), Speech Input: Assessment and Evaluation. In A. Cole (editor), Survey of the State of the Art in Human Language Technology Ronald. <http://cslu.cse.ogi.edu/HLTsurvey/>.
- [3] A. Fourcin, G. Harland, W. Barry & V. Hazan, (1989), Speech input and output assessment. Multilingual methods and standards. Ellis Horwood Ltd., Chichester.
- [4] Dafydd Gibbon, Roger Moore and Richard Winski (eds.) (1997), Handbook of Standards and Resources for Spoken Language Systems. Mouton de Gruyter, Berlin, New York.