

# Overview of the CLEF-2007 Cross-Language Speech Retrieval Track

Pavel Pecina and Petra Hoffmannová  
Institute of Formal and Applied Linguistics, Charles University  
Malostranske namesti 25, 118 00 Praha 1, Czech Republic  
pecina@ufal.mff.cuni.cz  
hoffmannova@knih.mff.cuni.cz

Gareth J.F. Jones and Ying Zhang  
School of Computing  
Dublin City University, Dublin 9, Ireland  
gjones@computing.dcu.ie  
yzhang@computing.dcu.ie

Douglas W. Oard  
College of Information Studies and Institute for Advanced Computer Studies  
University of Maryland, College Park, MD 20742, USA  
oard@umd.edu

## Abstract

The CLEF-2007 Cross-Language Speech Retrieval (CL-SR) track included two tasks: to identify topically coherent segments of English interviews in a known-boundary condition, and to identify time stamps marking the beginning of topically relevant passages in Czech interviews in an unknown-boundary condition. Six teams participated in the English evaluation, performing both monolingual and cross-language searches of ASR transcripts, automatically generated metadata, and manually generated metadata. Four teams participated in the Czech evaluation, performing monolingual searches of automatic speech recognition transcripts.

## Categories and Subject Descriptors

H.3 [Information Storage and Retrieval]: H.3.3 Information Search and Retrieval

## General Terms

Measurement, Performance, Experimentation

## Keywords

Speech Retrieval, Evaluation

# 1 Introduction

The 2007 Cross-Language Evaluation Forum (CLEF) Cross-Language Speech Retrieval (CL-SR) track is the third and final year for evaluation of ranked retrieval from spontaneous conversational speech from an oral history collection at CLEF. As in the CLEF 2006 CL-SR task [2], automatically transcribed interviews conducted in English could be searched using queries in one of six languages, and automatically transcribed interviews conducted in Czech could be searched using queries in one of two languages. New relevance judgments for additional topics were created to expand the Czech collection in 2007. The English collection used in 2007 was the same as that used in 2006. As in CLEF 2005 and CLEF 2006, the English task was based on a known-boundary condition for topically coherent segments. The Czech task was based on a unknown-boundary condition in which participants were required to identify a time stamp for the beginning of each distinct topically relevant passage.

The remainder of this paper is organized as follows. Section 2 describes the English task and summarizes the results for the submitted runs. Section 3 does the same for the Czech task. The paper concludes in Section 4 with a brief recap of what has been learned across all three years of the CLEF CL-SR track.

## 2 English Task

The structure of the CLEF 2007 CL-SR English task was identical to that used in 2006, which we review here briefly (see [2] for more details).

### 2.1 Segments

The “documents” searched in the English task are 8,104 segments that were designated by professional indexers as topically coherent. A detailed description of the structure and fields of the English segment collection is given in the 2005 track overview paper [3]. Automatically generated transcripts from two Automatic Speech Recognition (ASR) systems are available. The ASR-TEXT2006B field contains a transcript generated using the best presently available ASR system, which has a mean word error rate of 25% on held-out data. Only 7,378 segments have text in this field. For the remaining 726 segments, no ASR output was available from that system, so in those cases the ASRTEXT2006B field includes content identical to the ASRTEXT2004A field (which has a 35% mean word error rate) which was generated using an earlier less accurate transcription system. An extensive set of manually and automatically generated metadata is also available for each segment.

### 2.2 Topics

The same 63 training topics and 33 evaluation topics were used for the English task this year as had been used in 2006. Participating teams were asked not to use the evaluation topics for system tuning. Translations into Czech, Dutch, French, German, and Spanish had been created by native speakers of those languages. Participating teams were asked to submit runs for 105 topics (the 63 training topics, the 33 evaluation topics, and 9 further topics for which relevance data is not currently available, to support possible future construction of new relevance assessment pools), but results are reported only for the 33 evaluation topics.

### 2.3 Evaluation Measure

As in the CLEF-2006 CL-SR track, we report uninterpolated Mean Average Precision (MAP) as the principal measure of retrieval effectiveness. Version 8.0 of the `trec_eval` program was used to

compute this measure.<sup>1</sup> The Wilcoxon signed-rank signed test was employed for evaluation of significance.

## 2.4 Relevance Judgments

Subject matter experts created multi-scale and multi-level relevance assessments in the same manner as was done for the CLEF-2005 CL-SR track [3]. These were then conflated into binary judgments using the same procedure as was used for CLEF-2005: the union of direct and indirect relevance judgments with scores of 2, 3, or 4 (on a 0–4 scale) were treated as topically relevant, and any other case as non-relevant. This resulted in a total of 20,560 binary judgments across the 33 topics, among which 2,449 (12%) are relevant. Results from 2007 may not be strictly comparable with results from 2006; both were generated from the same initial set of relevance judgments, but those judgments were filtered at different sites in 2006 and 2007, in both cases to remove judgments for segments that are not contained in the distributed collection, and we have not yet done a detailed comparison of the results of those filtering process.

## 2.5 Techniques

This section gives a brief description of the methods used by each team participating in the English task. Additional details are available in each team’s paper.

### 2.5.1 Brown University (BLLIP)

The Brown Laboratory for Linguistic Information Processing (BLLIP) team extended the basic Dirichlet-smoothed unigram IR model to incorporate bigram mixing and collection smoothing. In their enhanced language model, the bigram and unigram models were mixed using a tunable mixture weight over all documents. They attempted linearly mixing the test collection with two larger text corpora, 40,000 sentences from the Wall Street Journal and 450,000 sentences from the North American News Corpus, in order to alleviate the sparse data problems in the case of small collections. They observed that bigram statistics appeared to have greater impact with pseudo-relevance feedback than without. The collection smoothing approach clearly provided a substantial improvement.

### 2.5.2 Dublin City University (DCU)

Dublin City University concentrated on the issues of topic translation, combining this with search field combination and pseudo-relevance feedback methods used for their CLEF 2006 submissions. Non-English topics were translated into English using the Yahoo! BabelFish free online translation service combined with domain-specific translation lexicons gathered automatically from Wikipedia. The combination of multiple fields using the BM25F variant of Okapi weights was explored. Additionally, they integrated their information retrieval methods based on the Okapi model with summary-based pseudo-relevance feedback.

### 2.5.3 University of Amsterdam (UVA)

The University of Amsterdam explored the use of character  $n$ -gram tokenization to improve the retrieval of documents using automatically generated text, as well as the combination of manually generated with automatically generated text. They reported that  $n = 4$  provided the best retrieval effectiveness when the cross-word overlapping  $n$ -gram tokenization strategy is used. The field combination was done using the Indri query language, in which varying weights were assigned to different fields. Cross-language experiments were conducted using manually created Dutch topics donated by the University of Twente. Dutch topics were automatically translated into English using two different online tools, SYSTRAN and FreeTranslation. The translations generated from each MT system were then combined as a ‘bag-of-words’ English query.

---

<sup>1</sup>The trec\_eval program is available from [http://trec.nist.gov/trec\\_eval/](http://trec.nist.gov/trec_eval/).

#### 2.5.4 University of Chicago (UC)

The University of Chicago team focused on the contribution of automatically assigned thesaurus terms to retrieval effectiveness and the utility of different query translation strategies. For French–English cross-language retrieval, they adopted two query translation strategies: MT-based translation using the publicly available translation tool provided by Google, and dictionary-based translation. Their dictionary-based translation procedure applied a backoff stemming strategy in order to support matching with highest precision between the query terms and the bilingual word list. They noted that 27% of the French query terms remained untranslated and were thus retained.

#### 2.5.5 University of Jaén (SINAI)

The SINAI group at the University of Jaén investigated the effect of selection of different fields (referred to as “labels” in their paper) on retrieval effectiveness. The Information Gain measure was employed to select the best XML tags in the document collection. The tags with higher values of Information Gain were selected to compose the final collection. Their experiments were conducted with the Lemur retrieval information system by applying the KL-divergence weighing scheme. The French, German and Spanish topics were translated to English using a translation module, SINTRAM, which works with different online machine translators and combines the different translations based on heuristics.

#### 2.5.6 University of Ottawa (UO)

The University of Ottawa used weighted summation of normalized similarity measures to combine 15 different weighting schemes from two IR systems (Terrier and SMART). Two query expansion techniques, one based on the thesaurus and the other one on blind relevance feedback, were examined. In their cross-language experiments, the queries were automatically translated from French and Spanish into English by combining the results of multiple online machine translation tools. Results for an extensive set of locally scored runs were also reported.

## 2.6 Results

Table 1 summarizes the evaluation results for all 29 official runs averaged over the 33 evaluation topics, listed in descending order of MAP. These 29 runs were further categorized into four groups based on the query language used (English or non-English) and the document fields (automatic-only or at least one manual assigned) indexed: 9 automatic-only monolingual runs, 6 automatic-only cross-language runs, 9 monolingual runs with manually assigned metadata, and 5 cross-language runs with manually assigned metadata.

### 2.6.1 Automatic-Only Monolingual Runs

Teams were required to run at least one monolingual condition using the title (T) and description (D) fields of the topics and indexing only automatically generated fields; the best “required runs” are shown in bold in Table 2 as a basis for cross-system comparisons. The University of Ottawa (0.0855), Dublin City University (0.0787), and the BLLIP team (0.0785) reported comparable results (no significant difference at the 95% confidence level). These results are statistically significant better than those reported by the next two teams, the University of Chicago (0.0571) and the University of Amsterdam (0.0444), which were statistically indistinguishable from each other.

### 2.6.2 Automatic-Only Cross-Language Runs

As shown in Table 3, the best result (0.0636) for cross-language runs on automatically generated indexing data (a French–English run from Dublin City University) achieved 81% of the monolingual retrieval effectiveness with comparable conditions (0.0787 as shown in Table 2).

Run ID	MAP	Lang	Query	Document Fields	Site
dcuEnTDNmanual	0.2847	EN	TDN	MK,SUM	DCU
uoEnTDtManF1	0.2761	EN	TD	MK,SUM	UO
brown.TDN.man	0.2577	EN	TDN	MK,SUM	BLLIP
dcuEnTDmanualauto	0.2459	EN	TD	MK,SUM,ASR06B	DCU
brown.TD.man	0.2366	EN	TD	MK,SUM	BLLIP
brown.T.man	0.2348	EN	T	MK,SUM	BLLIP
UvA_4.enopt	0.2088	EN	TD	MK,SUM,ASR06B	UVA
dcuFrTDmanualauto	0.1980	FR	TD	MK,SUM,ASR06B	DCU
UvA_5.nlopt	0.1408	NL	TD	MK,SUM,AK2,ASR06B	UVA
<b>uoEnTDtQExF1</b>	0.0855	EN	TD	AK1,AK2,ASR04	UO
uoEnTDtQExF2	0.0841	EN	TD	AK1,AK2,ASR04	UO
brown.TDN.auto	0.0831	EN	TDN	AK1,AK2,ASR06B	BLLIP
<b>dcuEnTDauto</b>	0.0787	EN	TD	AK1,AK2,ASR06B	DCU
<b>brown.TD.auto</b>	0.0785	EN	TD	AK1,AK2,ASR06B	BLLIP
SinaiSp100	0.0737	ES	TD	ALL	SINAI
dcuFrTDauto	0.0636	FR	TD	AK1,AK2,ASR06B	DCU
uoEsTDtF2	0.0619	ES	TD	AK1,AK2,ASR04	UO
uoFrTDtF2	0.0603	FR	TD	AK1,AK2,ASR04	UO
SinaiFr100	0.0597	FR	TD	ALL	SINAI
SinaiEn100	0.0597	EN	TD	ALL	SINAI
SinaiSp050	0.0579	ES	TD	SUM,AK1,AK2,ASR04,ASR06A,ASR06B	SINAI
<b>UCkwENTD</b>	0.0571	EN	TD	AK1,AK2,ASR06B	UC
SinaiEn050	0.0515	EN	TD	SUM,AK1,AK2,ASR04,ASR06A,ASR06B	SINAI
<b>UCbaseENTD1</b>	0.0512	EN	TD	ASR06B	UC
<b>UvA_2.en4g</b>	0.0444	EN	TD	AK2,ASR06B	UVA
UvA_1.base	0.0430	EN	TD	ASR06B	UVA
UCkwFRTD1	0.0406	FR	TD	AK1,AK2,ASR06B	UC
UvA_3.nl4g	0.0400	NL	TD	AK2,ASR06B	UVA
UCbaseFRTD1	0.0322	FR	TD	ASR06B	UC

Table 1: Evaluation results for all English official runs. MK = MANUALKEYWORD (Manual metadata), SUM = SUMMARY (Manual metadata), AK1 = AUTOKEYWORD2004A1 (Automatic), AK2 = AUTOKEYWORD2004A2, ASR03 = ASRTEXT2003A (Automatic), ASR04 = ASRTEXT2004A (Automatic), ASR06A = ASRTEXT2006A (Automatic), ASR06B = ASRTEXT2006B (Automatic), and ALL = all fields.

### 2.6.3 Monolingual Runs With Manual Metadata

For monolingual TD runs on manually generated indexing data, the University of Ottawa achieved the best result (0.2761), which is statistically significantly better than all other runs under comparable conditions, as shown in Table 4. For TDN runs, the DCU result (0.2847) is not statistically significantly better than that obtained by BLLIP (0.2577).

### 2.6.4 Cross-Language Runs With Manual Metadata

The evaluation results for cross-language runs on manually generated indexing data are shown in Table 5. The best cross-language result (0.1980), representing 81% of monolingual retrieval effectiveness under comparable conditions (0.2459 shown in Table 4), was achieved by DCU's French-English run.

Run ID	MAP	Lang	Query	Document Fields	Site
<b>uoEnTDtQExF1</b>	<b>0.0855</b>	EN	TD	AK1,AK2,ASR04	UO
uoEnTDtQExF2	0.0841	EN	TD	AK1,AK2,ASR04	UO
brown.TDN.auto	0.0831	EN	TDN	AK1,AK2,ASR06B	BLLIP
<b>dcuEnTDauto</b>	<b>0.0787</b>	EN	TD	AK1,AK2,ASR06B	DCU
<b>brown.TD.auto</b>	<b>0.0785</b>	EN	TD	AK1,AK2,ASR06B	BLLIP
<b>UCkwENTD</b>	<b>0.0571</b>	EN	TD	AK1,AK2,ASR06B	UC
<b>UCbaseENTD1</b>	<b>0.0512</b>	EN	TD	ASR06B	UC
<b>UvA_2_en4g</b>	<b>0.0444</b>	EN	TD	AK2,ASR06B	UVA
UvA_1_base	0.0430	EN	TD	ASR06B	UVA

Table 2: Evaluation results for automatic English monolingual runs. Bold runs are the required condition. AK1 = AUTOKEYWORD2004A1, AK2 = AUTOKEYWORD2004A2, ASR03 = ASRTEXT2003A, ASR04 = ASRTEXT2004A, ASR06A = ASRTEXT2006A, and ASR06B = ASRTEXT2006B.

Run ID	MAP	Lang	Query	Document Fields	Site
dcuFrTDauto	0.0636	FR	TD	AK1,AK2,ASR06B	DCU
uoEsTDtF2	0.0619	ES	TD	AK1,AK2,ASR04	UO
uoFrTDtF2	0.0603	FR	TD	AK1,AK2,ASR04	UO
UCkwFRTD1	0.0406	FR	TD	AK1,AK2,ASR06B	UC
UvA_3_nl4g	0.0400	NL	TD	AK2,ASR06B	UVA
UCbaseFRTD1	0.0322	FR	TD	ASR06B	UC

Table 3: Evaluation results for automatic cross-language runs. AK1 = AUTOKEYWORD2004A1, AK2 = AUTOKEYWORD2004A2, ASR04 = ASRTEXT2004A, and ASR06B = ASRTEXT2006B.

Run ID	MAP	Lang	Query	Document Fields	Site
dcuEnTDNmanual	0.2847	EN	TDN	MK,SUM	DCU
uoEnTDtManF1	0.2761	EN	TD	MK,SUM	UO
brown.TDN.man	0.2577	EN	TDN	MK,SUM	BLLIP
dcuEnTDmanualauto	0.2459	EN	TD	MK,SUM,ASR06B	DCU
brown.TD.man	0.2366	EN	TD	MK,SUM	BLLIP
brown.T.man	0.2348	EN	T	MK,SUM	BLLIP
UvA_4.enopt	0.2088	EN	TD	MK,SUM,ASR06B	UVA
SinaiEn100	0.0597	EN	TD	ALL	SINAI
SinaiEn050	0.0515	EN	TD	SUM,AK1,AK2,ASR04,ASR06A,ASR06B	SINAI

Table 4: Evaluation results for monolingual English runs with manual metadata. MK = MANUALKEYWORD, SUM = SUMMARY, AK1 = AUTOKEYWORD2004A1, AK2 = AUTOKEYWORD2004A2, ASR04 = ASRTEXT2004A, ASR06A = ASRTEXT2006A, ASR06B = ASRTEXT2006B, and ALL = all fields.

### 3 Czech Task

The structure of the Czech task was quite similar to the one used in the 2006 with differences which we describe in the following subsections. Further details can be found in the 2006 track

Run ID	MAP	Lang	Query	Document Fields	Site
dcuFrTDmanualauto	0.1980	FR	TD	MK,SUM,ASR06B	DCU
UvA_5_nlopt	0.1408	NL	TD	MK,SUM,AK2,ASR06B	UVA
SinaiSp100	0.0737	ES	TD	ALL	SINAI
SinaiFr100	0.0597	FR	TD	ALL	SINAI
SinaiSp050	0.0579	ES	TD	SUM,AK1,AK2,ASR04,ASR06A,ASR06B	SINAI

Table 5: Evaluation results for cross-language runs with manual metadata. MK = MANU-ALKEYWORD, SUM = SUMMARY, AK1 = AUTOKEYWORD2004A1, AK2 = AUTOKEYWORD2004A2, ASR04 = ASRTEXT2004A, ASR06A = ASRTEXT2006A, ASR06B = ASRTEXT2006B, and ALL = all fields.

overview paper [2].

### 3.1 Interviews

The default “quickstart” collection was generated from the same set of 357 Czech interviews as in 2006. It contained 11,377 overlapping passages with the following fields:

**DOCNO** containing a unique document number in the same format as the start times that systems were required to produce in a ranked list.

**INTERVIEWDATA** containing the first name and last initial for the person being interviewed. This field is identical for every passage that was generated from the same interview.

**ASRSYSTEM** specifying the type of the ASR transcript, where “2004” and “2006” denote colloquial and formal Czech transcripts respectively.

**CHANNEL** specifying which recorded channel (left or right) was used to produce the transcript.

**ASRTEXT** containing words in order from the transcript selected by ASRSYSTEM and CHANNEL for a passage beginning at the start time indicated in DOCNO.

The average passage duration in the default 2007 quickstart collection is 3.75 minutes, and each passage has a 33% overlap with the subsequent passage (i.e., passages begin about every 2.5 minutes).

No thesaurus terms (neither manual nor automatic, neither English nor Czech) were distributed with the collection this year. This step was taken as an expedient because it was not practical to correct the time misalignment that was present in the 2006 quickstart collection for the manually assigned thesaurus terms (and because automatically assigned thesaurus terms had not proven to be useful in 2006, perhaps because of poorly matched training data having been used to train the classifier).

### 3.2 Topics

This year we released a total of 118 topics: 105 original English topics from 2006, 10 broadened Czech topics from 2006, and 3 new broadened topics that were constructed this year. All topics were originally created in English and then translated into Czech by native speakers. Some minor errors in the Czech translations from last year were corrected.<sup>2</sup> Translations into other languages were not distributed with the collection. No teams used the English topics this year; all official runs with the Czech collection were monolingual.

<sup>2</sup>The corrected topics were 1259, 1282, 1551, 14313, and 24313. Of these, only topic 14313 was used in the 2006 Czech task evaluation, and none have been used for reported official results in the English task to date.

Topic	# rel						
1192	18	2265	113	3019	14	4005	68
1345	12	2358	126	3021	16	4006	135
1554	46	2384	37	3022	29	4007	51
1829	6	2404	8	3023	78	4009	10
1897	31	3000	41	3024	105	4011	132
1979	17	3001	102	3026	33	4012	61
2000	114	3002	95	3027	86	14313	17
2006	63	3007	107	3028	199	15601	108
2012	90	3008	53	3032	9	15602	25
2185	25	3010	18	4001	35		
2224	63	3016	40	4004	13		

Table 6: Number of relevant passages identified for each of the evaluation topics.

Participating teams were asked to run all 118 available topics. Two of the 118 topics were used as assessment training topics and excluded from the evaluation, 29 topics were available for training systems (with relevance judgments from 2006), 50 of the remaining 87 topics were selected as possible evaluation topics (with at least 6 relevant passages identified during the search-guided assessment phase), highly-ranked (i.e., “pooled”) assessment was completed for 42 of those 50 topics, and those 42 were used as the evaluation topics.

### 3.3 Evaluation Measure

The evaluation measure used in the Czech task is the same as in 2006. It’s based on the mean Generalized Average Precision (mGAP) measure, which was originally introduced to deal with human assessments of partial relevance [1]. In our case, the human assessments are binary but the degree of match to those assessments can be partial. The Wilcoxon signed-rank signed test was employed for evaluation of significance.

### 3.4 Relevance Judgments

Relevance judgments were completed at Charles University in Prague for 42 topics this year under the same conditions as in 2006 by six relevance assessors. Evaluation topics had been selected to have at least six relevant start times in the Czech collection in order to minimize the effect of quantization noise on the computation of mGAP. A total of 2,389 start times for relevant passages were identified, thus yielding an average of 56 relevant passages per topic (minimum 6, maximum 199). Table 3.4 shows the number of relevant start times for each of the 42 topics. To support future experiments on searching a bilingual speech collection, 34 of the 2007 CLEF CL-SR Czech task evaluation topics also present in the 2007 CLEF CL-SR English task collection (as training, evaluation, or unused topics).<sup>3</sup>

### 3.5 Techniques

All participating teams employed existing information retrieval systems to perform monolingual retrieval and submitted total of 15 runs for official scoring. Each team submitted at least one run in the required title+description condition. The narrative field was used only in two runs by University of West Bohemia. Most of the teams used only automatically generated queries. Manual query construction was performed only by Charles University. All teams used the provided quickstart collection for at least some runs. The University of West Bohemia also used the quickstart scripts with different parameters to generate another collection for some experiments.

<sup>3</sup>The exceptions being the broadened topics, which are the 4000-series.

Run name	mGAP score	Query language	Query construction	Topic fields	Document fields	Term normalization	Site name
UWB_2-1_tdn_l	0.0264	CZ	Auto	TDN	ASR2006	lemma	UWB
UWB_3-1_tdn_l	0.0237	CZ	Auto	TDN	ASR2006	lemma	UWB
UWB_2-1_tdn_s	0.0228	CZ	Auto	TD	ASR2006	stem	UWB
UCcsaTD2	0.0203	CZ	Auto	TD	ASR2006	aggressive stem	UC
prague04	0.0190	CZ	Auto	TD	ASR2006	lemma	CUNI
UCcslTD1	0.0189	CZ	Auto	TD	ASR2006	light stem	UC
prague01	0.0187	CZ	Auto	TD	ASR2006	lemma	CUNI
prague02	0.0181	CZ	Manual	TD	ASR2006	lemma	CUNI
UWB_3-1_tdn_l	0.0131	CZ	Auto	TD	ASR2006	lemma	UWB
UWB_2-1_tdn_w	0.0129	CZ	Auto	TD	ASR2006	none	UWB
UCunstTD3	0.0126	CZ	Auto	TD	ASR2006	none	UC
brown.s.f	0.0114	CZ	Auto	TD	ASR2006	light stem	BLLIP
brown.sA.f	0.0106	CZ	Auto	TD	ASR2006	aggressive stem	BLLIP
prague03	0.0102	CZ	Manual	TD	ASR2006	none	CUNI
brown.f	0.0052	CZ	Auto	TD	ASR2006	none	BLLIP

Table 7: Czech official runs.

### 3.5.1 Brown University (BLLIP)

The system of Brown University was based on the language model paradigm for retrieval and implemented in the Indri system. A unigram language model, Czech-specific stemming, and pseudo-relevance feedback were applied in three officially submitted runs.

### 3.5.2 Charles University (CUNI)

The Charles University team performed experiments with Indri retrieval model from the Lemur project with pseudo-relevance feedback, stopwords removal, and morphological lemmatization obtained by in-house morphological analysis and a part-of-speech tagger. The team submitted four official runs; two of them employed manual query construction.

### 3.5.3 University of Chicago (UC)

The University of Chicago employed the InQuery information retrieval system with stop-word removal and three different stemming approaches: no stemming, light stemming, and aggressive stemming. Three runs were submitted for official scoring.

### 3.5.4 University of West Bohemia (UWB)

The University of West Bohemia employed a TF\*IDF model with blind relevance feedback implemented in Lemur. Five runs submitted for official scoring differed in methods used for word normalization (none, lemmatization, stemming), in formulas used for term weighting (Raw TF, BM25), and in topic fields used (TDN, TD).

### 3.5.5 Results

The results of all official runs evaluated on 42 topics are reported in Table 7. The effect of term normalization handling the rich Czech morphology is quite significant. The runs employing any type of term normalization (stemming or lemmatization) outperform systems indexing only original word forms with no normalization by 61–119%. The scores of directly comparable runs are given in Table 8, all the differences are statistically significant at a 95% confidence level.

The second collection generated by the University of West Bohemia generated some interesting insights. They used the quickstart scripts distributed with the test collection to decrease the

Run name	mGAP score	mGAP improvement	Query construction	Topic fields	Term normalization	Site name
UWB_2-1.td_s	0.0228	+76.7%	Auto	TD	stem	UWB
UWB_2-1.td_w	0.0129		Auto	TD	none	UWB
UCcsaTD2	0.0203	+61.1%	Auto	TD	aggressive stem	UC
UCunstTD3	0.0126		Auto	TD	none	UC
prague02	0.0181	+77.5%	Manual	TD	lemma	CUNI
prague03	0.0102		Manual	TD	none	CUNI
brown.s.f	0.0114	+119.2%	Auto	TD	light stem	BLLIP
brown.f	0.0052		Auto	TD	none	BLLIP

Table 8: Comparison of systems with term normalization and without normalization.

average passage duration (from 3.75 minutes to 2.5 minutes) and to increase the overlap between subsequent passages (from 33% to 50%). This had the effect of substantially decreasing the average start-time spacing between passages (from 2.5 for 1.25 minutes). This resulted in an apparent improvement in mGAP (compare UWB\_2-1.tdn.l: mGAP=0.0264 and UWB\_3-1.tdn.l: mGAP=0.0237) that turned out not to be statistically significant. The two-sided width of the scoring window is set at 5 minutes in our evaluation script, so this range of start time spacings is well within the scorable range, but more closely spaced passages offer some potential for reducing quantization noise in the evaluation script. Although we compute evaluation results only from start times, our assessors marked both start and end times. Their average duration of a marked relevant passage is 2.83 minutes, which seems to be somewhat better matched to the 2.5 minutes passages used in the University of West Bohemia’s alternate condition (2.5 minutes for UWB\_2-1.tdn.l, 3.75 minutes for UWB\_3-1.tdn.l and all runs from other sites).

The Charles University team reported on the first experiments with interactive use of the Czech collection. One of their runs based on manual query construction turned out to be statistically indistinguishable from a run under comparable conditions from the same team with queries that were generated automatically, and a second run with manually formed queries did not do well at all (probably because lemmatization was not used in that second run).

## 4 Conclusion and Future Plans

Like all CLEF tracks, the CL-SR track had three key goals: (1) to develop evaluation methods and reusable evaluation resources for an important information access problem in which cross-language information access is a natural part of the task, (2) to generate results that can provide a strong baseline against which future research results with the same evaluation resources can be compared, and (3) to foster the development of a research community with the experience and expertise to make those future advances. In the case of the CL-SR track, those goals have now been achieved. Over 3 years, research teams from 14 universities in 6 countries submitted 123 runs for official scoring, and additional locally scored runs have been reported in papers published by those research teams. The resulting English and Czech collections are the first information retrieval test collections of substantial size for spontaneous conversational speech, unique characteristics of the English collection have fostered new research comparing searches based on automatic speech recognition and manually assigned metadata, and unique characteristics of the Czech collection have inspired new research on evaluation of information retrieval from unsegmented speech.

Now that the track has been completed, these new CLEF test collections will be made available to nonparticipants through the Evaluations and Language Resources Distribution Agency (ELDA). The training data for the automatic speech retrieval systems that were used to generate the transcripts in those collections is also expected to become available soon, most likely through the Linguistic Data Consortium (LDC). It is our hope that these resources will be used together to investigate more closely coupled techniques than have been possible to date with just the

present CLEF CL-SR test collections. Looking further forward, we believe that it is now time for the information retrieval research community to look beyond oral history to other instances of spontaneous conversational speech such as that found in recordings of meetings, historically significant telephone conversations, and broadcast conversations (e.g., call-in radio “talk shows”). We also believe that it would be productive to begin to explore the application of some of the technology developed for this track to improve access to a broad range oral history collections and similar cultural heritage materials (e.g., interviews contained in broadcast archives). Together, these directions for future work will likely continue to extend the legacy and impact of this initial investment in exploring the retrieval of information from spontaneous conversational speech.

## Acknowledgments

This year’s track would not have been possible without the efforts of a great many people. Our heartfelt thanks go to the dedicated group of relevance assessors in Prague without whom the Czech collection simply would not exist, to Scott Olsson for helping to prepare the English collection this year, to Ayelet Goldin and Jianqiang Wang for their timely help with critical details of the Czech relevance assessment and scoring process, to Jan Hajic for his support and advice throughout, and to Carol Peters for her seemingly endless patience. This work has been supported in part by NSF IIS award 0122466 (MALACH) and by the Ministry of Education of the Czech Republic, projects MSM 0021620838 and #1P05ME786.

## References

- [1] Jaana Kekalainen and Kalervo Jarvelin. Using graded relevance assessments in IR evaluation. In *Journal of the American Society for Information Science and Technology*, 2002.
- [2] Douglas W. Oard, Jianqiang Wang, Gareth J. F. Jones, Ryen W. White, Pavel Pecina, Dagobert Soergel, Xiaoli Huang, and Izhak Shafran. Overview of the CLEF-2006 cross-language speech retrieval track. In *Proceedings of the CLEF 2006 Workshop on Cross-Language Information Retrieval and Evaluation*, 2006.
- [3] Ryen W. White, Douglas W. Oard, Gareth J. F. Jones, Dagobert Soergel, and Xiaoli Huang. Overview of the CLEF-2005 cross-language speech retrieval track. In *Proceedings of the CLEF 2005 Workshop on Cross-Language Information Retrieval and Evaluation*, 2005.