

# Overview of the CLEF-2006 Cross-Language Speech Retrieval Track

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

Jianqiang Wang  
Department of Library and Information Studies  
State University of New York at Buffalo, Buffalo, NY 14260, U.S.A.  
jw254@buffalo.edu

Gareth J.F. Jones  
School of Computing  
Dublin City University, Dublin 9, Ireland  
Gareth.Jones@computing.dcu.ie

Ryen W. White  
Microsoft Research  
One Microsoft Way, Redmond, WA 98052, U.S.A.  
ryenw@microsoft.com

Pavel Pecina  
MFF UK, Malostranske namesti 25, Room 422  
Charles University, 118 00 Praha 1, Czech Republic  
pecina@ufal.mff.cuni.cz

Dagobert Soergel and Xiaoli Huang  
College of Information Studies  
University of Maryland, College Park, MD 20742, U.S.A.  
{dsoergel,xiaoli}@umd.edu

Izhak Shafran  
OGI School of Science & Engineering, Oregon Health and Sciences University  
20000 NW Walker Rd, Portland, OR 97006, U.S.A.  
zak@cslu.ogi.edu

## Abstract

The CLEF-2006 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. Five teams participated in the English evaluation, performing both monolingual and cross-language searches of ASR transcripts, automatically generated metadata, and manually generated metadata. Results indicate that the 2006 evaluation topics are more challenging than those used in 2005, but that cross-language searching continued to pose no unusual challenges when compared with collections of character-coded text. Three teams participated in the Czech evaluation, but no team achieved results comparable to those obtained with English interviews. The reasons for this outcome are not yet clear.

## 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, Generalized Average Precision

# 1 Introduction

The 2006 Cross-Language Evaluation Forum (CLEF) Cross-Language Speech Retrieval (CL-SR) track continues last year's effort to support research on ranked retrieval from spontaneous conversational speech. Automatically transcribing spontaneous speech has proven to be considerably more challenging than transcribing the speech of news anchors for the Automatic Speech Recognition (ASR) techniques on which fully-automatic content-based search systems are based.

The CLEF 2005 CL-SR task focused on searching English interviews. For CLEF 2006, 30 new search topics were developed for the same collection, and an improved ASR transcript with better accuracy for the same set of testimonies was added. This made it possible to validate the retrieval techniques that were shown to be effective with last year's topics, and to further explore the influence of ASR accuracy on the retrieval effectiveness. The CLEF 2006 CL-SR track also added a new task of searching Czech interviews.

Similar to CLEF 2005, the English task is again based on a known-boundary condition for topically coherent segments. The Czech search task is based on a unknown-boundary condition where participants are required to identify a time stamp for the beginning of each distinct topically relevant passage.

The first part of this paper describes the English language CL-SR task and summarizes the participants' submitted results. This is followed by a description of the Czech language task with corresponding details of submitted runs.

# 2 English Task

The structure of the CLEF 2006 CL-SR English task was identical to that used in 2005. Two English collections were released this year. The first release (March 14, 2006) contained all material that was now available for training (i.e., both the training and the test topics from last year's CLEF 2005 CL-SR evaluation). There was one small difference from the original 2005 data release: each person's last name that appears in the NAME field (or in the associated XML data files) was

reduced into its initial followed by three dots (e.g., “Smith” became “S...”). This collection contains a total of 63 search topics, 8,104 topically coherent segments (the equivalent of “documents” in a classic IR evaluation), and 30,497 relevance judgments.

The second release (June 5, 2006) included a re-release of all the training materials (unchanged) and an additional 42 candidate evaluation topics (30 new topics, plus 12 other topics for which relevance judgments had not previously been released) and two new fields based on an improved ASR transcript from the IBM T. J. Watson Research Center.

## 2.1 Segments

Other than the changes described above, the segments used for the CLEF 2006 CL-SR task were identical to those used for CLEF 2005. Two new fields contain ASR transcripts of higher accuracy than were available in 2005 (ASRTEXT2006A and ASRTEXT2006B). The ASRTEXT2006A 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. Because of time constraints, however, only 7,378 segments have text in this field. For the remaining 726 segments, no ASR output was available from the 2006A system at the time the collection was distributed. The ASRTEXT2006B field seeks to avoid this no-content condition by including content identical to the ASRTEXT2006A field when available, and content identical to the ASRTEXT2004A field otherwise. Since ASRTEXT2003A, ASRTEXT2004A, and ASRTEXT2006B contain ASR text that was automatically generated for all 8,104 segments, any (or all) of them can be used for the required run based on automatic data. A detailed description of the structure and fields of the English segment collection is given in last year’s track overview paper [11].

## 2.2 Topics

The limited size of the collection would likely make it impractical to continue to do new topic development for the same set of segments in future years, so we elected to use every previously unreleased topic for the CLEF-2006 CLSR English task. A total of 30 new topics were created for this year’s evaluation from actual requests received by the USC Shoah Foundation Institute for Visual History and Education.<sup>1</sup> These were combined with 12 topics that had been developed in previous years, but for which relevance judgments had not been released. This resulted in a set of 42 topics that were candidates for use in the evaluation.

All topics were initially prepared in English. Translations into Czech, Dutch, French, German, Spanish were created by native speakers of those languages, and the same process was used to prepare French translations of the narrative field for all topics in the training collection (which had not been produced in 2005 due to resource constraints). With the exception of Dutch, all translations were checked for reasonableness by a second native speaker of the language.<sup>2</sup>

A total of 33 of the 42 candidate topics were used as a basis for the official 2006 CL-SR evaluation; the remaining 9 topics were rejected because they had either too few known relevant segments (fewer than 5) or too high a density of known relevant segments among the available judgments (over 48%, suggesting that many relevant segments may not have been found). Participating teams were asked to submit results for all 105 available topics (the 63 topics in the 2006 training set and the 42 topics in the 2006 evaluation candidate set) so that new pools could be formed to perform additional judgments on the development set if additional assessment resources become available.

---

<sup>1</sup>On January 1, 2006 the University of Southern California (USC) Shoah Foundation Institute for Visual History and Education was established as the successor to the Survivors of the Shoah Visual History Foundation, which had originally assembled and manually indexed the collection used in the CLEF CL-SR track.

<sup>2</sup>A subsequent quality assurance check for Dutch revealed only a few minor problems. Both the as-run and the final corrected topics will therefore be released for Dutch.

## 2.3 Evaluation Measure

As in the CLEF-2005 CL-SR track, we report Mean uninterpolated 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>3</sup>

## 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 [11]. These were then conflated into binary judgments using the same procedure as was used for CLEF-2005: the union of direct 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 28,223 binary judgments across the 33 topics, among which 2,450 (8.6%) are relevant.

## 2.5 Techniques

The following gives a brief description of the methods used by the participants in the English task. Additional details are available in each team’s paper.

### 2.5.1 University of Alicante (UA)

The University of Alicante used the MINIPAR parser to produce an analysis of syntactic dependencies in the topic descriptions and in the automatically generated portion of the collection. They then used these results in combination with their locally developed IR-n system to produce overlapping passages. Their experiments focused on combining these sources of evidence and on optimizing search effectiveness using pruning techniques.

### 2.5.2 Dublin City University (DCU)

Dublin City University used two systems based on the Okapi retrieval model. One version used Okapi with their summary-based pseudo relevance feedback method. The other system explored combination of multiple segment fields using the method introduced in [8]. This system also explored the use of a field-based method for term selection in query expansion with pseudo-relevance feedback.

### 2.5.3 University of Maryland (UMD)

The University of Maryland team tried two techniques, using the InQuery system in both cases [1]. Four fields of automatic data were combined to create a segment index. Retrieval results from this index were compared with results from index based on individual automatic data field, showing that combining the four automatic data fields could slightly help, although the observed improvement is not statistically significant. Manual metadata fields were also combined in the same, but no comparative results were reported. In addition, the team also applied the so-called “meaning matching” technique to French-English cross-language retrieval. Although there is some sign showing the technique helps marginally, the CLIR effectiveness is significantly worse than monolingual performance.

### 2.5.4 Universidad Nacional de Educacin a Distancia (UNED)

The UNED team compared the utility of the 2006 ASR with manually generated summaries and manually assigned keywords. A CLIR experiment was performed using Spanish queries with the 2006 ASR.

---

<sup>3</sup>The `trec_eval` program is available from [http://trec.nist.gov/trec\\_eval/](http://trec.nist.gov/trec_eval/). The DCU results reported in this paper are based on a subsequent re-submission that corrected a formatting error.

### 2.5.5 University of Ottawa (UO)

The University of Ottawa used two information retrieval systems in their experiments: SMART [2] and Terrier [7]. The two systems were used with many different weighting schemes for indexing the segments and the queries, and with several query expansion techniques (including a new proposed method based on log-likelihood scores for collocations). For the English collection, different Automatic Speech Recognition transcripts (with different estimated word error rates) were used for indexing the segments, and also several combinations of automatic transcripts. Cross-language experiments were run after the topics were automatically translated into English by combining the results of several online machine translation tools. The manual summaries and manual keywords were used for indexing in the manual run.

### 2.5.6 University of Twente (UT)

The University of Twente employed a locally developed XML retrieval system that supports Narrowed Extended XPath (NEXI) queries to search the collection. They also prepared Dutch translations of the topics that they used as a basis for CLIR experiments.

## 2.6 English evaluation results

Table 1 summarizes the results for all 30 official runs averaged over the 33 evaluation topics, listed in descending order of MAP. Required runs are shown in bold. The best results for the required condition (title plus description queries, automatically generated data, from Dublin City University) of 0.0747 are considerably below (i.e., just 58% of) last year’s best results. A similar effect was not observed when manually generated metadata were indexed, however with this year’s best result (0.2902) being 93% of last year’s best manually generated metadata result. From this we conclude that this year’s topic set seems somewhat less well matched with the ASR results, but that the topics are not otherwise generally much harder for information retrieval techniques based on term matching. CLIR also seemed to pose no unusual challenges with this year’s topic set, with the best CLIR on automatically generated indexing data (a French run from the University of Ottawa) achieving 83% of the MAP achieved by a comparable monolingual run. Similar effects were observed with manually generated metadata (at 80% of the corresponding monolingual MAP for Dutch queries, from the University of Twente).

## 3 Czech Task

The goal of the Czech task was to automatically identify the start points of topically-relevant passages in interviews. Ranked lists for each topic were submitted by each system in the same form as the CLEF ad hoc task, with the single exception that a system-generated starting point was specified rather than a document identifier. The format for this was “VHF[IntCode].[starting-time],” where “IntCode” is the five-digit interview code (with leading zeroes added) and “starting-time” is the system-suggested replay starting point (in seconds) with reference to the beginning of the interview. Lists were to be ranked by systems in the order that they would suggest for listening to passages beginning at the indicated points.

### 3.1 Interviews

The Czech task was broadly similar to the English task in that the goal was to design systems that could help searchers identify sections of an interview that they might wish to listen to. The processing of the Czech interviews was, however, different from that used for English in three important ways:

- No manual segmentation was performed. This alters the format of the interviews (which for Czech is time-oriented rather than segment-oriented), it alters the nature of the task (which

Run name	MAP	Lang	Query	Doc field	Site
uoEnTDNtMan	0.2902	EN	TDN	MK,SUM	UO
3d20t40f6sta5flds	0.2765	EN	TDN	ASR06B,AK1,AK2,N,SUM,MK	DCU
umd.manu	0.2350	EN	TD	N,MK,SUM	UMD
UTsummkENor	0.2058	EN	T	MK,SUM	UT
dcuEgTDall	0.2015	EN	TD	ASR06B,AK1,AK2,N,SUM,MK	DCU
uneden-manualkw	0.1766	EN	TD	MK	UNED
UTsummkNl2or	0.1654	NL	T	MK,SUM	UT
dcuFchTDall	0.1598	FR	TD	ASR06B,AK1,AK2,N,SUM,MK	DCU
umd.manu.fr.0.9	0.1026	FR	TD	N,MK,SUM	UMD
umd.manu.fr.0	0.0956	FR	TD	N,MK,SUM	UMD
unedes-manualkw	0.0904	ES	TD	MK	UNED
unedes-summary	0.0871	ES	TD	SUM	UNED
uoEnTDNsQEx04A	0.0768	EN	TDN	ASR04,AK1,AK2	UO
<b>dcuEgTDauto</b>	<b>0.0733</b>	EN	TD	ASR06B,AK1,AK2	DCU
uoFrTDNs	0.0637	FR	TDN	ASR04,AK1,AK2	UO
uoSpTDNs	0.0619	ES	TDN	ASR04,AK1,AK2	UO
<b>uoEnTDt04A06A</b>	<b>0.0565</b>	EN	TD	ASR04,ASR06B,AK1,AK2	UO
<b>umd.auto</b>	<b>0.0543</b>	EN	TD	ASR04,ASR06B,AK1,AK2	UMD
UTasr04aEN	0.0495	EN	T	ASR04	UT
dcuFchTDauto	0.0462	FR	TD	ASR06B,AK1,AK2	DCU
UA_TDN_FL_ASR06BA1A2	0.0411	EN	TDN	ASR06B,AK1,AK2	UA
UA_TDN_ASR06BA1A2	0.0406	EN	TDN	ASR06B,AK1,AK2	UA
UA_TDN_ASR06BA2	0.0381	EN	TDN	ASR06B,AK2	UA
UTasr04aNl2	0.0381	NL	T	ASR04	UT
<b>UTasr04aEN-TD</b>	<b>0.0381</b>	EN	TD	ASR04	UT
<b>uneden</b>	<b>0.0376</b>	EN	TD	ASR06B	UNED
<b>UA_TD_ASR06B</b>	<b>0.0375</b>	EN	TD	ASR06B	UA
UA_TD_ASR06BA2	0.0365	EN	TD	ASR06B,AK2	UA
unedes	0.0257	ES	TD	ASR06B	UNED
umd.auto.fr.0.9	0.0209	FR	TD	ASR04,ASR06B,AK1,AK2	UMD

Table 1: English official runs. Bold runs are required. N = Name (Manual metadata), MK = Manual Keywords (Manual metadata), SUM = Summary (Manual metadata), ASR04 = ASR-TEXT2004A (Automatic) AK1 = AUTOKEYWORD2004A1 (Automatic), AK2 = AUTOKEYWORD2004A2. See [11] for descriptions of these fields. (Automatic)

for Czech is to identify replay start points rather than to select among predefined segments), and it alters the nature of the manually assigned metadata (there are no manually written summaries for Czech and the meaning of a manual thesaurus term assignment for Czech is that discussion of a topic started at that time).

- The two available Czech ASR transcripts were generated using different ASR systems. In both cases, the acoustic models were trained using 15-minute snippets from 336 speakers, all of whom are present in the test set as well. However, the language model was created by interpolating two models—an in-domain model from transcripts, and an out-of-domain model from selected portions of Czech National Corpus. For details, see the baseline systems described in [9, 10]. Apart from the improvement in transcription accuracy, the 2006 system differs from the 2004 system in that the transcripts are produced in formal Czech, rather than the colloquial Czech that was produced in 2004. Since the topics were written in formal Czech, the 2006 ASR transcripts may yield better matching. Interview-specific vocabulary priming (adding proper names to the recognizer vocabulary based on names present in a pre-interview questionnaire) was not done for either Czech system. Thus, a somewhat higher error rate on named entities might be expected for the Czech systems than for the two English systems (2004 and 2006) in which vocabulary priming was included.

- ASR is available for both the left and right stereo channels (which usually were recorded from microphones with different positions and orientations).

Because the task design for Czech is not directly compatible with the design of document-oriented IR systems, we provided a “quickstart” package containing the following:

- A quickstart script for generating overlapping passages directly from the ASR transcripts. The passage duration (in seconds), the spacing between passage start times (also in seconds), and the desired ASR system (2004 or 2006) could be specified. The default settings (180, 60, and 2006) result in 3-minute passages in which one minute on each end overlaps with the preceding or subsequent passage.
- A quickstart collection created by running the quickstart script with the default settings. This collection contains 11,377 overlapping passages.

The quickstart collection contains the following automatically generated fields:

**DOCNO** The DOCNO field contains a unique document number in the same format as the start times that systems were required to produce in a ranked list. This design allowed the output of a typical IR system to be used directly as a list of correctly formatted (although perhaps not very accurate) start times for scoring purposes.

**ASRSYSTEM** specifying the source of the ASR text collection (either “2004” for the colloquial Czech system developed by the University of West Bohemia and Johns Hopkins University in 2004 or “2006” for an updated and possibly more accurate formal Czech system provided by the same research groups in 2006).

**CHANNEL** The CHANNEL field specifies which recorded channel (left or right) was used to produce the transcript. The channel that produced the greatest number of total words over the entire transcript (which is usually the channel that produced the best ASR accuracy for words spoken by the interviewee) was automatically selected by default. This automatic selection process was hardcoded in the script, although the script could be modified to generate either or both channels.

**ASRTEXT** The ASRTEXT field contains words in order from the transcript selected by ASRSYSTEM and CHANNEL for a passage beginning at the start time indicated in DOCNO. When the selected transcript contains no words at all from that time period, words are drawn from one alternate source that is chosen in the following priority order: (1) the same ASRSYSTEM from the other CHANNEL, (2) the same CHANNEL from the other ASRSYSTEM, or (3) the other CHANNEL from the other ASRSYSTEM.

**ENGLISHAUTOKEYWORD** The ENGLISHAUTOKEYWORD field contains a set of thesaurus terms that were assigned automatically using a k-Nearest Neighbor (kNN) classifier using only words from the ASRTEXT field of the passage; the top 20 thesaurus terms are included in best-first order. Thesaurus terms (which may be phrases) are separated with a vertical bar character. The classifier was trained using English data (manually assigned thesaurus terms and manually written segment summaries) and run using automatically produced English translations of the 2006 Czech ASRTEXT [6]. Two types of thesaurus terms are present, but not distinguished: (1) terms that express a subject or concept; (2) terms that express a location, often combined with time in one precombined term [5]. Because the classifier was trained on the English collection, in which thesaurus terms were assigned with segments, the natural interpretation of an automatically assigned thesaurus term is that the classifier believes the indicated topic is associated with the word spoken in this passage. Note that this differs from the way in which the presence of a manually assigned thesaurus term (described below) should be interpreted.

**CZECHAUTOKEYWORD** The CZECHAUTOKEYWORD field contains Czech translations of the ENGLISHAUTOKEYWORD field. These translations were obtained from three sources: (1) professional translation of about 3,000 thesaurus terms, (2) volunteer translation of about 700 thesaurus terms, and (3) a custom-built machine translation system that reused words and phrases from manually translated thesaurus terms to produce additional translations. Some words (e.g., foreign place names) remained untranslated when none of the three sources yielded a usable translation.

Three additional fields containing data produced by human indexers at the Survivors of the Shoah Visual History Foundation were also available for use in contrastive conditions:

**INTERVIEWDATA** The INTERVIEWDATA field contains 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.

**ENGLISHMANUKEYWORD** The ENGLISHMANUALKEYWORD field contains thesaurus terms that were manually assigned with one-minute granularity from a custom-built thesaurus by subject matter experts at the Survivors of the Shoah Visual History Foundation while viewing the interview. The format is the same as that described for the ENGLISHAUTOKEYWORD field, but the meaning of a keyword assignment is different. In the Czech collection, manually assigned thesaurus terms are used as onset marks—they appear only once at the point where the indexer recognized that a discussion of a topic or location-time pair had started; continuation and completion of discussion are not marked.

**CZECHMANUKEYWORD** The CZECHMANUALKEYWORD field contains Czech translations of the English thesaurus terms that were produced from the ENGLISHMANUALKEYWORD field using the process described above.

All three teams used the quickstart collection; no other approaches to segmentation and no other settings for passage length or passage start time spacing were tried.

## 3.2 Topics

At the time the Czech evaluation topics were released, it was not yet clear which of the available topics were likely to yield a sufficient number of relevant passages in the Czech collection. Participating teams were therefore asked to run 115 topics—every available topic at that time. This included the full 105 topic set that was available this year for English (including all training and all evaluation candidate topics) and 10 adaptations of topics from that set in which geographic restrictions had been removed (as insurance against the possibility that the smaller Czech collection might not have adequate coverage for exactly the same topics).

All 115 topics had originally been constructed in English and then translated into Czech by native speakers. Since translations into languages other than Czech were not available for the 10 adapted topics, only English and Czech topics were distributed with the Czech collection. No teams used the English topics this year; all official runs this year with the Czech collection were monolingual.

Two additional topics were created as part of the process of training relevance assessors, and those topics were distributed to participants along with a (possibly incomplete) set of relevance judgments. This distribution occurred too late to influence the design of any participating system.

## 3.3 Evaluation Measure

The evaluation measure that we chose for Czech is designed to be sensitive to errors in the start time, but not in the end time, of system-recommended passages. It is computed in the same manner as mean average precision, but with one important difference: partial credit is awarded in a way that rewards system-recommended start times that are close to those chosen by assessors.



After a simulation study, we chose a symmetric linear penalty function that reduces the credit for a match by 0.1 (absolute) for every 15 seconds of mismatch (either early or late) [4]. This results in the same computation as the well-known mean Generalized Average Precision (mGAP) measure that was introduced to deal with human assessments of partial relevance [3]. In our case, the human assessments are binary; it is the degree of match to those assessments that can be partial. Relevance judgments are drawn without replacement so that only the highest ranked match (including partial matches) can be scored for any relevance assessment; other potential matches receive a score of zero. Differences at or beyond a 150 second error are treated as a no-match condition, thus not “using up” a relevance assessment.

### 3.4 Relevance Judgments

Relevance judgments were completed at Charles University in Prague for a total of 29 Czech topics by subject matter experts who were native speakers of Czech. All relevance assessors had good English reading skills. Topic selection was performed by individual assessors, subject to the following factors:

- At least five relevant start times in the Czech collection were required in order to minimize the effect of quantization noise on the computation of mGAP.
- The greatest practical degree of overlap with topics for which relevance judgments were available in the English collection was desirable.

Once a topic was selected, the assessor iterated between topic research (using external resources) and searching the collection. A new search system was designed to support this interactive search process. The best channel of the Czech ASR and the manually assigned English thesaurus terms were indexed as overlapping passages, and queries could be formed using either or both. Once a promising interview was found, an interactive search within the interview could be performed using either type of term and promising regions were identified using a graphical depiction of the retrieval status value. Assessors could then scroll through the interview using these indications, the displayed English thesaurus terms, and the displayed ASR transcript as cues. They could then replay the audio from any point in order to confirm topical relevance. As they did this, they could indicate the onset and conclusion of the relevant period by designating points on the transcript that were then automatically converted to times with 15-second granularity.<sup>4</sup> Only the start times are used for computation of the mGAP measure, but both start and end times are available for future research.

Once that search-guided relevance assessment process was completed, the assessors were provided with a set of additional points to check for topical relevance that were computed using a pooling technique similar to that used for English. The top 50 start times from every official run were pooled, duplicates (at one minute granularity) were removed, and the results were inserted into the assessment system as system recommendations. Every system recommendation was checked, although assessors exercised judgment regarding when it would be worthwhile to actually listen to the audio in order to limit the cost of this “highly ranked” assessment process. Relevant passages identified in this way were added to those found using search-guided assessment to produce the final set of relevance judgments (topic 4000 was generalized from a pre-existing topic).

A total of 1,322 start times for relevant passages were identified, thus yielding an average of 46 relevant passages per topic (minimum 8, maximum 124). Table 2 shows the number of relevant start times for each of the 29 topics, 28 of which are the same as topics used in the English test collection.

---

<sup>4</sup>Several different types of time spans arise when describing evaluation of speech indexing systems. For clarity, we have tried to stick to the following terms when appropriate: manually defined segments (for English indexing), 15-minute snippets (for ASR training), 15-second increments (for the start and end time of Czech relevance judgments), relevant passages (identified by Czech relevance assessors), and automatically generated passages (for the quickstart collection).

topid	#rel	topid	#rel	topid	#rel	topid	#rel	topid	#rel
1166	8	1181	21	1185	50	1187	26	1225	9
1286	70	1288	9	1310	20	1311	27	1321	27
14312	14	1508	83	1620	35	1630	17	1663	34
1843	52	2198	18	2253	124	3004	43	3005	84
3009	77	3014	87	3015	50	3017	83	3018	26
3020	67	3025	51	3033	45	4000	65		

Table 2: Number of the relevant passages identified for each of the 29 topics in the Czech collection.

## 3.5 Techniques

The participating teams all employed existing information retrieval systems to perform monolingual searches of the quickstart collection.

### 3.5.1 University of Maryland (UMD)

The University of Maryland submitted three official runs in which they tried combining all the fields (Czech ASR text, Czech (manual and automatic) keyword, and the English translations of the keywords) to form a unified passage index using Inquiry. They compared the retrieval results based on this index with those based on ASR alone or the combination of automatic keywords and ASR text.

### 3.5.2 University of Ottawa (UO)

Three runs were submitted from the University of Ottawa for the Czech task using SMART and one run was submitted using Terrier.

### 3.5.3 University of West Bohemia (UWB)

The University of West Bohemia was the only team to apply morphological normalization and stopword removal for Czech. A classic TF\*IDF model was implemented in Lemur, along with the Lemur implementation of blind relevance feedback. Five runs were submitted for official scoring, and one additional run was scored locally.

### 3.5.4 Results

With two exceptions, the mean Generalized Average Precision (mGAP) values were between 0.0003 and 0.0005. In a side experiment reported in the UWB paper, random permutation of the possible start times was found to yield a mGAP of 0.0005 in a simulation study. We therefore conclude that none of those runs demonstrated any useful degree of system support for the task.

Two runs yielded more interesting results. The best official run, from UO, achieved a mGAP of 0.0039, and a run that was locally scored at UWB achieved a mGAP of 0.0015. Interestingly, these are two of the three runs in which the ENGLISHMANUALKEYWORD field was used. A positive influence from that factor would require that untranslated English terms (e.g., place names) match terms that were present in the topic descriptions (either with or without morphological normalization). The UWB paper provides an analysis that suggests that the beneficial effect of using that field may be limited to a single topic.

The use of overlapping passages in the quickstart collection probably reduced mGAP values substantially because the design of the measure tends to penalize duplication. Specifically, the start time of the highest-ranking passage that matches a passage start time in the relevance judgments will “use up” that judgment. Subsequent passages in which the same matching terms were present would then receive no credit at all (even if they were closer matches). We had

Run name	mGAP	Lang	Query	Doc field	Site
uoCzEnTDNsMan	0.0039	CZ,EN	TDN	CAK,CMK,EAK,EMK	UO
uoCzTDNsMan	0.0005	CZ	TDN	ASR,CAK,CMK	UO
uoCzEnTDt	0.0005	CZ,EN	TD	ASR,CAK	UO
umd.asr	0.0005	CZ	TD	ASR	UMD
uoCzTDNs	0.0004	CZ	TDN	ASR,CAK	UO
uoCzTDs	0.0004	CZ	TD	ASR,CAK	UO
UWB_mk_aTD	0.0004	CZ	TD	ASR,CMK	UWB
UWB_mk_a_akTD	0.0004	CZ	TD	ASR,CAK,CMK	UWB
UWB_mk_a_akTDN	0.0004	CZ	TDN	ASR,CAK,CMK	UWB
umd.akey.asr	0.0004	CZ	TD	ASR,CAK,EAK	UMD
UWB_aTD	0.0003	CZ	TD	ASR	UWB
UWB_a_akTD	0.0003	CZ	TD	ASR,CAK	UWB
umd.all	0.0003	CZ	TD	ASR,CAK,CMK,EAK,EMK	UMD

Table 3: Czech official runs. Bold runs are required. CAK = CZECHAUTOKEYWORD (Automatic), EAK = ENGLISHAUTOKEYWORD (Automatic), CMK = CZECHMANUKEYWORD (Manual metadata), EMK = ENGLISHMANUKEYWORD (Manual metadata)

originally intended the quickstart collection to be used only for out-of-the-box sanity checks, with the idea that teams would either modify the quickstart scripts or create new systems outright to explore a broader range of possible system designs. Time pressure and a lack of a suitable training collection precluded that sort of experimentation, however, and the result was that this undesirable effect of passage overlap affected every system.

Other possible explanations for the relatively poor results also merit further investigation. This is the first time that mGAP has been used in this way to evaluate actual system results, so it is possible that the measure is poorly designed or that there is a bug in the scoring script. Simulation studies suggest that is not likely to be the case, however. This is also the first time that Czech ASR has been used, and it is the first time that relevance assessment has been done in Czech (using a newly designed system). So there are many possible factors that need to be explored. This year’s Czech collection is exactly what we need for such an investigation, so it should be possible to make significant progress over the next year.

## 4 Conclusion and Future Plans

The CLEF 2006 CL-SR track extended the previous year’s work on the English task by adding new topics, and introduced a new Czech task with a new unknown-boundary evaluation condition. The results of the English task suggest that the evaluation topics this year posed somewhat greater difficulty for systems doing fully automatic indexing. Studying what made these topics more difficult would be an interesting scope for future work. However, the most significant achievement of this year’s track was the development of a CL-SR test collection based on a more realistic unknown-boundary condition. Now that we have both that collection and an initial set of system designs, we are in a good position to explore issues of system and evaluation design that clearly have not yet been adequately resolved.

We expect that it would be possible to continue the CLEF CL-SR track in 2007 if there is sufficient interest. For Czech, it may be possible to obtain relevance judgments for additional topics, perhaps increasing to a total of 50 the number of topics that the track can leave as a legacy for use by future researchers. Developing additional topics for English seems to be less urgent (and perhaps less practical), but we do expect to be able to provide additional automatically generated indexing data (either ASR for additional interviews, word lattices in some form, or both) if there is interest in further work with the English collection. Some unique characteristics of the CL-SR collection may also be of interest to other tracks, including domain-specific retrieval and geoCLEF.

We look forward to discussing these and other issues when we meet in Alicante!

## 5 Acknowledgments

This 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 Maryland and Prague, to the Dutch, French and Spanish teams that helped with topic translation, and to Bill Byrne, Martin Cetkovsky, Bonnie Dorr, Ayelet Goldin, Sam Gustman, Jan Hajic, Jimmy Lin, Baolong Liu, Craig Murray, Scott Olsson, Bhuvana Ramabhadran and Deborah Wallace for their help with creating the techniques, software, and data sets on which we have relied.

## References

- [1] Broglio, J., Callan, J. P. and Croft, W. B.: INQUERY System Overview. In Proceedings of the Annual Meeting of the Association for Computational Linguistics, pages 47–67, 1993.
- [2] Buckley, C., Salton, G., and Allan, J.: Automatic retrieval with locality information using SMART. In Proceedings of the First Text REtrieval Conference (TREC-1), pages 59–72. NIST Special Publication 500-207, 1993.
- [3] Kekalainen, J. and Jarvelin, K.: Using graded relevance assessments in IR evaluation. *Journal of the American Society for Information Science and Technology*, 53(13)1120–1129, 2002.
- [4] Liu, B. and Oard, D. W.: One-sided measures for evaluating ranked retrieval effectiveness with spontaneous conversational speech. In Proceedings of the 29th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pages 673–674, 2006.
- [5] Murray, C., Dorr, B. J., Lin, J., Hajic, J. and Pecina, P.: Leveraging Reusability: Cost-effective Lexical Acquisition for Large-scale Ontology Translation. In Proceedings of the Annual Meeting of the Association for Computational Linguistics, 2006.
- [6] Olsson, J. S., Oard, D. W. and Hajic, J.: Cross-language text classification. In Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval: 645–646, 2005.
- [7] Ounis, I., Amati, G., Plachouras, A., He, B., Macdonald, C. and Johnson, D.: Terrier Information Retrieval Platform. In Proceedings of the 27th European Conference on Information Retrieval (ECIR 05), 2005.
- [8] Robertson, S. E., Zaragoza, H., and Taylor, M.: Simple BM25 Extension to Multiple Weighted Fields, Proceedings of the 13th ACM International Conference on Information and Knowledge Management, pages 42–49, 2004.
- [9] Shafran, I. and Byrne, W.: Task-Specific Minimum Bayes-Risk Decoding using Learned Edit Distance, In Proceedings of INTERSPEECH2004-ICSLP, vol. 3, pages 1945–1948, 2004.
- [10] Shafran, I. and Hall, K.: Corrective Models for Speech Recognition of Inflected Languages, Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP), 2006.
- [11] White, R. W., Oard, D. W., Jones, G. J. F., Soergel, D., and Huang, X.: Overview of the CLEF-2005 Cross-Language Speech Retrieval Track, Proceedings of the CLEF 2005 Workshop on Cross-Language Information Retrieval and Evaluation, 2005.