

# Unsupervised Morpheme Analysis Evaluation by IR experiments – Morpho Challenge 2007

Mikko Kurimo, Mathias Creutz, Ville Turunen  
Adaptive Informatics Research Centre, Helsinki University of Technology  
P.O.Box 5400, FIN-02015 TKK, Finland  
Mikko.Kurimo@tkk.fi

## Abstract

This paper presents the evaluation of Morpho Challenge Competition 2 (information retrieval). The Competition 1 (linguistic gold standard) is described in a companion paper. In Morpho Challenge 2007, the objective was to design statistical machine learning algorithms that discover which morphemes (smallest individually meaningful units of language) words consist of. Ideally, these are basic vocabulary units suitable for different tasks, such as text understanding, machine translation, information retrieval, and statistical language modeling. In this paper the morpheme analysis submitted by the Challenge participants were evaluated by performing information retrieval (IR) experiments, where the words in the documents and queries were replaced by their proposed morpheme representations and the search was based on morphemes instead of words. The IR evaluations were provided for three languages: Finnish, German, and English and the participants were encouraged to apply their algorithm to all of them. The challenge organizers performed the IR experiments using the queries, texts, and relevance judgments available in CLEF forum and morpheme analysis methods submitted by the challenge participants. The results show that the morpheme analysis has a significant effect in IR performance in all languages, and that the performance of the best unsupervised methods can be superior to the supervised reference methods. The challenge was part of the EU Network of Excellence PASCAL Challenge Program and organized in collaboration with CLEF.

## Categories and Subject Descriptors

H.3 [Information Storage and Retrieval]: H.3.1 Content Analysis and Indexing; H.3.3 Information Search and Retrieval

## General Terms

Algorithms, Performance, Experimentation

## Keywords

Morphological analysis, Machine learning

## 1 Introduction

The scientific objectives of the Morpho Challenge 2007 were: to learn of the phenomena underlying word construction in natural languages, to advance machine learning methodology, and to discover approaches suitable for a wide range of languages. The suitability for a wide range of languages is

becoming increasingly important, because language technology methods need to be quickly and as automatically as possible extended to new languages that have limited previous resources. That is why learning the morpheme analysis directly from large text corpora using unsupervised machine learning algorithms is such an attractive approach and a very relevant research topic today.

Morpho Challenge 2007 is a follow-up to our previous Morpho Challenge 2005 (Unsupervised Segmentation of Words into Morphemes) [8]. In Morpho Challenge 2005 the focus was in the segmentation of data into units that are useful for statistical modeling. The specific task for the competition was to design an unsupervised statistical machine learning algorithm that segments words into the smallest meaning-bearing units of language, morphemes. In addition to comparing the obtained morphemes to a linguistic "gold standard", their usefulness was evaluated by using them for training statistical language models for speech recognition.

In Morpho Challenge 2007 a more general focus was chosen to not only to segment words into smaller units, but also to perform *morpheme analysis* of the word forms in the data. For instance, the English words "boot, boots, foot, feet" might obtain the analyses "boot, boot + plural, foot, foot + plural", respectively. In linguistics, the concept of morpheme does not necessarily directly correspond to a particular word segment but to an abstract class. For some languages there exist carefully constructed linguistic tools for this kind of analysis, although not for many, but using statistical machine learning methods we may still discover interesting alternatives that may rival even the most careful linguistically designed morphologies.

The problem of learning the morphemes directly from large text corpora using an unsupervised machine learning algorithm is clearly a difficult one. First the words should be somehow segmented into meaningful parts, and then these parts should be clustered in the abstract classes of morphemes that would be useful for modeling. It is also challenging to learn to generalize the analysis to rare words, because even the largest text corpora are very sparse, a significant portion of the words may occur only once. Many important words, for example proper names and their inflections or some forms of long compound words, may also not exist in the training material at all, and their analysis is often even more challenging. However, benefits for successful morpheme analysis, in addition to obtaining a set of basic vocabulary units for modeling, can be seen for many important tasks in language technology. The additional information included in the units can provide support for building more sophisticated language models, for example, in speech recognition [1], machine translation [9], and information retrieval [12].

The evaluation of the unsupervised morpheme analysis was in this challenge solved by developing two complementary evaluations, one including a comparison to linguistic morpheme analysis gold standard, and another including a practical real-world application where morpheme analysis might be used. This paper presents how the application-oriented evaluation called *Competition 2* was performed and the companion paper [7] describes the linguistic evaluation called *Competition 1*. As a practical real-world application domain we chose the problem of finding useful index terms for information retrieval tasks in multiple languages. Traditionally, and especially in processing English texts, stemming algorithms have been used to reduce the different inflected word forms into the common roots or stems for indexing. However, to achieve best results when ported to new languages the development of stemming algorithms requires a considerable amount of special development work. In many highly-inflecting, compounding, and agglutinative European languages the amount of different word forms is huge and the task of extracting the useful index terms becomes both more complex and more important.

The same IR tasks that were attempted using the Morpho Challenge participants' morpheme analysis, were also tested by a number of reference methods to see how the unsupervised morpheme analysis performed in comparison to them. These references included the organizers' public Morfessor Categories-Map [3] and Morfessor Baseline [2, 4], the Morfessor analysis improved by a hybrid method [11], grammatical morpheme analysis based on the linguistic gold standards [5], the traditional Porter stemming [10] of words and also by the words as such without any processing.

## 2 Task

The Morpho Challenge 2007 task was to return the unsupervised morpheme analysis of every word form contained in a long word list supplied by the organizers for each test language [7]. The participants were pointed to corpora [7] in which the words occur, so that the algorithms may utilize information about word context. The information retrieval (IR) experiments were performed by the organizers based on the morpheme analyses submitted by the participants.. The words in the documents and the test queries were first replaced by their proposed morpheme representations and the search was then based on morphemes instead of words. To achieve the goal of designing language independent methods, the participants were encouraged to submit results in all test languages: Finnish, German and English.

## 3 Data sets

The data sets for testing the IR performance in each test language consisted of news paper articles as the source documents, test queries and the binary relevance judgments regarding to the queries. The organizers performed the IR experiments based on the morpheme analyses submitted by the participants, so it was not necessary for the participants to get these data sets. However, all the data was available for registered participants in the Cross-Language Evaluation Forum (CLEF)<sup>1</sup>.

The source documents were news articles collected from different news papers selected as follows:

- In Finnish: 55K documents from short articles in Aamulehti 1994-95, 50 test queries on specific news topics and 23K binary relevance assessments (CLEF 2004)
- In English: 170K documents from short articles in Los Angeles Times 1994 and Glasgow Herald 1995, 50 test queries on specific news topics and 20K binary relevance assessments (CLEF 2005).
- In German: 300K documents from short articles in Frankfurter Rundschau 1994, Der Spiegel 1994-95 and SDA German 1994-95, 60 test queries with 23K binary relevance assessments (CLEF 2003).

When performing the indexing and retrieval experiments for Competition 2, it turned out that the test data contained quite many new words in addition to those that were provided as training data for the Competition 1 [7]. Thus, the participants were offered a chance to improve the retrieval results of their morpheme analyses by providing them a list of the new words found in all test languages. The participants then had the choice to either run their algorithms to analyze as many of the new words as they could or liked, or to provide no extra analyses. No text data resources to find context for the new words were provided, but it was made possible to register to CLEF to use the text data available in there or any other data the participants could get.

## 4 Participants and their submissions

By the deadline in May, 2007, 6 research groups had submitted the segmentation results obtained by their algorithms. A total of 12 different algorithms were submitted, 8 of them ran experiments on all four test languages. All the submitted algorithms are listed in Table 1. In general, the submissions were all interesting and relevant and all of them met the exact specifications given and were able to get properly evaluated. In addition to the competitors' 12 morpheme analysis algorithms, we evaluated a number of reference methods described in Section 5.

The outputs of the submitted algorithms are analyzed closer in [7]. From the IR point of view it is interesting to note that only Monson and Zeman decided to provide several alternative analysis for most words instead of just the most likely one. McNamee's algorithms did not attempt to

---

<sup>1</sup><http://www.clef-campaign.org/>

Table 1: The submitted algorithms.

Algorithm	Authors	Affiliation
“Bernhard 1”	Delphine Bernhard	TIMC-IMAG, F
“Bernhard 2”	Delphine Bernhard	TIMC-IMAG, F
“Bordag 5”	Stefan Bordag	Univ. Leipzig, D
“Bordag 5a”	Stefan Bordag	Univ. Leipzig, D
“McNamee 3”	Paul McNamee and James Mayfield	JHU, USA
“McNamee 4”	Paul McNamee and James Mayfield	JHU, USA
“McNamee 5”	Paul McNamee and James Mayfield	JHU, USA
“Zeman ”	Daniel Zeman	Karlova Univ., CZ
“Monson Morfessor”	Christian Monson et al.	CMU, USA
“Monson ParaMor”	Christian Monson et al.	CMU, USA
“Monson ParaMor-Morfessor”	Christian Monson et al.	CMU, USA
“Pitler”	Emily Pitler and Samarth Keshava	Univ. Yale, USA
“Morfessor Categories-MAP”	The organizers	Helsinki Univ. Tech, FI
“Morfessor Baseline”	The organizers	Helsinki Univ. Tech, FI
“dummy”	The organizers	Helsinki Univ. Tech, FI
“grammatical”	The organizers	Helsinki Univ. Tech, FI
“Porter”	The organizers	Helsinki Univ. Tech, FI
“Tepper”	Michael Tepper	Univ. Washington, USA

provide a real morpheme analysis, but focused directly on finding a representative substring for each word type that would be likely to perform well in the IR evaluation.

## 5 Reference methods

To study and understand how the different morpheme analysis performed in the IR tasks, we attempted the same tasks with different reference methods. This also revealed us whether the unsupervised morpheme analysis (or even a supervised one) could really be useful in the IR tasks compared to simple word based indexing.

1. *Morfessor Categories-Map*: The same Morfessor Categories-Map (or here just “Morfessor MAP”, for short) as described in Competition 1 [7] was used for the unsupervised morpheme analysis. The stem vs. suffix tags were kept, but did not receive any special treatment in the indexing, because we did not want to favor this particular tagging.
2. *Morfessor Baseline*: All the words were simply split into smaller pieces without any morpheme analysis. This means that the obtained subword units were directly used as index terms as such. This was performed using the Morfessor Baseline algorithm as in Morpho Challenge 2005 [8]. We expected that this would not be optimal for IR, but because the unsupervised morpheme analysis is such a difficult task, this simple method would probably do quite well.
3. *dummy*: No words were split nor any morpheme analysis provided. This means that all were directly used as index terms as such without any stemming or tags. We expected that although the morpheme analysis should provide helpful information for IR, all the submissions would not probably be able to beat this brute force baseline. However, if some morpheme analysis method would consistently beat this baseline in all languages and task, it would mean that the method were probably useful in a language and task independent way.

4. *grammatical*: The words were analyzed using the gold standard in each language that were utilized as the “ground truth” in the Competition 1 [7]. Besides the stems and suffixes, the gold standard analyses typically consist of all kinds of grammatical tags which we decided to simply include as index terms, as well. For many words the gold standard analyses included several alternative interpretations that were all included in the indexing. However, we decided to also try the method adopted in the morpheme segmentation for Morpho Challenge 2005 [8] that only the first interpretation of each word is applied. This was here called “grammatical first” whereas the default was called “grammatical all”. Because our gold standards are quite small, 60k (English) - 600k (Finnish), compared to the amount of words that the unsupervised methods can analyze, we did not expect “grammatical” to perform particularly well, even though it would probably capture some useful indexing features to beat the “dummy”, at least.
5. *Porter*: No real morpheme analysis was performed, but the words were stemmed by the Porter stemming, an option provided by the Lemur toolkit. Because this is quite standard procedure in IR, especially for English text material, we expected this to provide the best results, at least for English. For the other languages the default Porter stemming was not likely to perform very well.
6. *Tepper*: A hybrid method developed by Michael Tepper [11] was utilized to improve the morpheme analysis reference obtained by our Morfessor Categories-MAP. Based on the obtained performance in Competition 1 [7], we expected that this could provide some interesting results here, as well.

## 6 Evaluation

In this evaluation, the organizers applied the analyses provided by the participants in information retrieval experiments. The words in the queries and source documents were replaced by the corresponding morpheme analyses provided by the participants, and the search was then based on morphemes instead of words. Any word that did not have a morpheme analysis was left un-replaced.

The evaluation was performed using a state-of-the-art retrieval method (the latest version of the freely available LEMUR toolkit<sup>2</sup>). We utilized two standard retrieval methods: Tfidf and Okapi term weighting. The Tfidf implementation in LEMUR applies term frequency weights for both query and document based on the BM25 weighting and the Euclidean dot-product as similarity measure. Okapi in LEMUR is an implementation of the BM25 retrieval function as described in [6].

The evaluation criterion was Uninterpolated Average Precision. There were several different categories and the winner with the highest Average Precision was selected separately for each language and each category:

1. All morpheme analyses from the training data are used as index terms “*withoutnew*” vs. additionally using also the morpheme analyses for new words that existed in the IR data but not in the training data “*withnew*”.
2. Tfidf term weighting was utilized for all index terms without any stoplists vs. Okapi term weighting for all index terms excluding an automatic stoplist constructed for each run separately and consisting of the most common terms (frequency threshold was 75,000 for Finnish and 150,000 for German and English). The stoplist was developed for the Okapi weighting, because otherwise Okapi weights were not suitable for the indexes that had many very common terms.

---

<sup>2</sup><http://www.lemurproject.org/>

## 7 Results

Table 2: The obtained average precision (AP%) in the information retrieval task for the submitted segmentations in **Finnish** (Competition 2 participants in bold and reference methods in normal font). Indexing is performed using Tfidf weighting for all morphemes (left) and Okapi weighting for all morphemes except the most common ones (stoplist) with frequency higher than 75,000 (right).

Tfidf weighting for all morphemes			Okapi weighting and a stoplist		
METHOD	WORDLIST	AP%	METHOD	WORDLIST	AP%
Morfessor baseline	withnew	0.4105	<b>Bernhard 2</b>	withnew	0.4915
<b>Bernhard 1</b>	withoutnew	0.4016	<b>Bernhard 1</b>	withnew	0.4681
grammatical first	withoutnew	0.3995	<b>Bernhard 2</b>	withoutnew	0.4425
<b>Bernhard 2</b>	withoutnew	0.3984	Morfessor baseline	withnew	0.4412
Morfessor baseline	withoutnew	0.3978	Morfessor MAP	withnew	0.4353
grammatical all	withoutnew	0.3952	<b>Bordag 5a</b>	withnew	0.4309
Morfessor MAP	withnew	0.3913	<b>Bordag 5</b>	withnew	0.4308
<b>Bernhard 1</b>	withnew	0.3896	grammatical all	withoutnew	0.4307
<b>Bordag 5</b>	withnew	0.3831	grammatical first	withoutnew	0.4216
Morfessor MAP	withoutnew	0.3814	<b>Bernhard 1</b>	withoutnew	0.4183
<b>Bernhard 2</b>	withnew	0.3811	grammatical first	withnew	0.4176
<b>Bordag 5</b>	withoutnew	0.3802	<b>Bordag 5a</b>	withoutnew	0.4147
grammatical first	withnew	0.3760	<b>Bordag 5</b>	withoutnew	0.4095
grammatical all	withnew	0.3734	grammatical all	withnew	0.4066
<b>Bordag 5a</b>	withoutnew	0.3721	Morfessor baseline	withoutnew	0.3820
<b>Bordag 5a</b>	withnew	0.3673	<b>McNamee 5</b>	withnew	0.3684
<b>McNamee 5</b>	withoutnew	0.3646	Morfessor MAP	withoutnew	0.3632
<b>McNamee 5</b>	withnew	0.3618	<b>McNamee 5</b>	withoutnew	0.3620
porter	withnew	0.3566	<b>McNamee 4</b>	withnew	0.3603
dummy	withnew	0.3559	<b>McNamee 4</b>	withoutnew	0.3567
<b>McNamee 4</b>	withoutnew	0.3518	porter	withnew	0.3517
<b>McNamee 4</b>	withnew	0.3257	<b>McNamee 3</b>	withoutnew	0.3386
<b>McNamee 3</b>	withoutnew	0.2941	dummy	withnew	0.3274
<b>Zeman</b>	withoutnew	0.2494	<b>McNamee 3</b>	withnew	0.3243
<b>McNamee 3</b>	withnew	0.2182	<b>Zeman</b>	withoutnew	0.2813

The results of the information retrieval evaluations are shown in Tables 2 – 4. In the Finnish task, the highest average precision was obtained by the “Bernhard 2” algorithm, which was also won the Competition 1 [7]. The highest average precision 0.49 was obtained using the Okapi weighting and stoplist for both the originally submitted morpheme analysis (for Competition 1) and the morpheme analysis for the new words added for Competition 2. For the competition category without the new words the winner was the same algorithm and without stoplist (using Tfidf) “Bernhard 1”.

The “Bernhard 1” algorithm obtained the highest average precision 0.47 for the German task using the new words, Okapi and stoplist. The same algorithm won also the categories for using Tfidf without stoplist with and without new words, but in using Okapi and stoplist without the new words the winner was the “Bernhard 2”. However, the difference between “Bernhard 1” and “Bernhard 2” was very small in all categories.

For English, the highest average precision was obtained by the “Bernhard 2” algorithm, which was also won the Competition 1 [7]. As in Finnish and German, the highest average precision 0.39 was obtained with the new words and using the Okapi weighting and stoplist. The same algorithm won also the category without the new words and using the Okapi weighting and stoplist, but in

Table 3: The obtained average precision (AP%) in the information retrieval task for the submitted segmentations in **German** (Competition 2 participants in bold and reference methods in normal font). Indexing is performed using TfIdf weighting for all morphemes (left) and Okapi weighting for all morphemes except the most common ones (stoplist) with frequency higher than 150,000 (right).

Tfidf weighting for all morphemes			Okapi weighting and a stoplist		
METHOD	WORDLIST	AP%	METHOD	WORDLIST	AP%
Morfessor baseline	withnew	0.3874	<b>Bernhard 1</b>	withnew	0.4729
Morfessor baseline	withoutnew	0.3826	<b>Bernhard 2</b>	withoutnew	0.4676
<b>Bernhard 1</b>	withoutnew	0.3777	<b>Bernhard 2</b>	withnew	0.4625
<b>Bernhard 2</b>	withoutnew	0.3731	<b>Bernhard 1</b>	withoutnew	0.4611
porter	withnew	0.3725	<b>Monson</b> Morfessor	withnew	0.4602
<b>Bernhard 1</b>	withnew	0.3720	Morfessor MAP	withnew	0.4571
<b>Bernhard 2</b>	withnew	0.3703	Morfessor baseline	withnew	0.4486
<b>Monson</b> Morfessor	withnew	0.3520	<b>Monson</b> Morfessor	withoutnew	0.4481
<b>Monson</b> Morfessor	withoutnew	0.3502	Morfessor MAP	withoutnew	0.4447
dummy	withnew	0.3496	Morfessor baseline	withoutnew	0.4417
<b>Bordag 5a</b>	withnew	0.3496	<b>Bordag 5</b>	withnew	0.4308
Morfessor MAP	withoutnew	0.3480	<b>Bordag 5</b>	withoutnew	0.4303
<b>McNamee 5</b>	withoutnew	0.3442	<b>Bordag 5a</b>	withnew	0.4259
Morfessor MAP	withnew	0.3397	<b>Bordag 5a</b>	withoutnew	0.4257
<b>McNamee 5</b>	withnew	0.3327	<b>Monson</b> ParaMor-M	withnew	0.4012
<b>Bordag 5</b>	withoutnew	0.3273	<b>Monson</b> ParaMor-M	withoutnew	0.3989
grammatical first	withoutnew	0.3223	porter	withnew	0.3866
<b>Monson</b> ParaMor-M	withnew	0.3200	<b>McNamee 5</b>	withoutnew	0.3617
grammatical first	withnew	0.3196	<b>McNamee 5</b>	withnew	0.3527
<b>Monson</b> ParaMor-M	withoutnew	0.3184	grammatical first	withoutnew	0.3467
<b>Bordag 5</b>	withnew	0.3156	<b>McNamee 4</b>	withoutnew	0.3453
grammatical all	withnew	0.3128	grammatical first	withnew	0.3445
grammatical all	withoutnew	0.3126	<b>McNamee 4</b>	withnew	0.3351
<b>McNamee 4</b>	withoutnew	0.3091	<b>Monson</b> ParaMor	withnew	0.3241
<b>McNamee 4</b>	withnew	0.3049	dummy	withnew	0.3228
<b>Monson</b> ParaMor	withnew	0.2887	<b>Monson</b> ParaMor	withoutnew	0.3224
<b>Monson</b> ParaMor	withoutnew	0.2861	grammatical all	withnew	0.3004
<b>Zeman</b>	withoutnew	0.2828	<b>McNamee 3</b>	withoutnew	0.2953
<b>McNamee 3</b>	withoutnew	0.2023	grammatical all	withoutnew	0.2926
<b>McNamee 3</b>	withnew	0.1945	<b>McNamee 3</b>	withnew	0.2868
			<b>Zeman</b>	withoutnew	0.2568

Table 4: The obtained average precision (AP%) in the information retrieval task for the submitted segmentations in **English** (Competition 2 participants in bold and reference methods in normal font). Indexing is performed using Tfidf weighting for all morphemes (left) and Okapi weighting for all morphemes except the most common ones (stoplist) with frequency higher than 150,000 (right).

Tfidf weighting for all morphemes			Okapi weighting and a stoplist		
METHOD	WORDLIST	AP%	METHOD	WORDLIST	AP%
porter	withnew	0.3052	porter	withnew	0.4083
<b>McNamee 5</b>	withoutnew	0.2888	<b>Bernhard 2</b>	withnew	0.3943
<b>McNamee 5</b>	withnew	0.2885	<b>Bernhard 2</b>	withoutnew	0.3922
Morfessor baseline	withnew	0.2863	<b>Bernhard 1</b>	withnew	0.3900
Morfessor baseline	withoutnew	0.2851	Morfessor baseline	withnew	0.3882
<b>McNamee 4</b>	withoutnew	0.2842	<b>Bernhard 1</b>	withoutnew	0.3881
<b>McNamee 4</b>	withnew	0.2838	Morfessor baseline	withoutnew	0.3869
Tepper	withoutnew	0.2784	grammatical first	withoutnew	0.3774
dummy	withnew	0.2783	grammatical first	withnew	0.3756
Morfessor MAP	withnew	0.2782	Tepper	withoutnew	0.3728
<b>Bernhard 1</b>	withoutnew	0.2781	<b>Monson</b> Morfessor	withoutnew	0.3721
<b>Bernhard 1</b>	withnew	0.2777	Morfessor MAP	withnew	0.3716
Morfessor MAP	withoutnew	0.2774	Morfessor MAP	withoutnew	0.3714
<b>Bernhard 2</b>	withnew	0.2682	<b>Monson</b> Morfessor	withnew	0.3703
<b>Monson</b> Morfessor	withoutnew	0.2676	<b>Pitler</b>	withoutnew	0.3652
<b>Bernhard 2</b>	withoutnew	0.2673	<b>Pitler</b>	withnew	0.3648
<b>Monson</b> Morfessor	withnew	0.2667	grammatical all	withoutnew	0.3621
<b>Pitler</b>	withoutnew	0.2666	grammatical all	withnew	0.3592
<b>Pitler</b>	withnew	0.2639	<b>McNamee 4</b>	withoutnew	0.3577
<b>Monson</b> ParaMor-M	withnew	0.2628	<b>McNamee 4</b>	withnew	0.3576
<b>Monson</b> ParaMor-M	withoutnew	0.2624	<b>McNamee 5</b>	withoutnew	0.3438
grammatical all	withoutnew	0.2619	<b>Monson</b> ParaMor-M	withnew	0.3435
grammatical first	withoutnew	0.2612	<b>McNamee 5</b>	withnew	0.3433
grammatical all	withnew	0.2602	<b>Bordag 5</b>	withoutnew	0.3427
grammatical first	withnew	0.2599	<b>Monson</b> ParaMor-M	withoutnew	0.3426
<b>Monson</b> ParaMor	withnew	0.2400	<b>Bordag 5</b>	withnew	0.3421
<b>Monson</b> ParaMor	withoutnew	0.2390	<b>Bordag 5a</b>	withoutnew	0.3409
<b>Zeman</b>	withoutnew	0.2297	<b>Bordag 5a</b>	withnew	0.3395
<b>Bordag 5</b>	withoutnew	0.2210	dummy	withnew	0.3123
<b>Bordag 5</b>	withnew	0.2202	<b>McNamee 3</b>	withoutnew	0.3047
<b>Bordag 5a</b>	withoutnew	0.2169	<b>McNamee 3</b>	withnew	0.3030
<b>Bordag 5a</b>	withnew	0.2165	<b>Monson</b> ParaMor	withnew	0.2835
<b>McNamee 3</b>	withoutnew	0.1695	<b>Monson</b> ParaMor	withoutnew	0.2821
<b>McNamee 3</b>	withnew	0.1677	<b>Zeman</b>	withoutnew	0.2674

categories using Tfidf without stoplist the winner was “McNamee 5”.

As expected, the “grammatical” reference method based on linguistic Gold Standard morpheme analysis [7] did not perform very well. However, with stoplist and Okapi term weighting it did achieve better results than the “dummy” method in all languages. In Finnish and English the performance was better than average, but quite poor in German. The “grammatical first” that utilized only the first of the alternative analysis in indexing was at least as good or better than the “grammatical all”, which seems to indicate that the alternative analysis are not very useful here.

For the “Morfessor” references it is interesting to note that they always performed better than the “grammatical”, which seems to suggest that the coverage of the analysis (“Morfessor” does not have any out-of-vocabulary words) is more important for IR than the grammatical correctness. The reason for “Morfessor Baseline” being almost always better than the more sophisticated “Morfessor Categories-MAP” is not clear, and also the hybrid “Tepper” improvement for Morfessor does not seem to affect the IR results. In general, the old “Morfessor Baseline” seems to provide a very good baseline in all tested languages also for the IR tasks as it did for the language modeling and speech recognition in [8]. Here, only the “Bernhard 1” and “Bernhard 2” methods managed to beat it.

## 8 Discussions

The comparison of the results in the Tfidf and Okapi categories show that the Okapi with stoplist performed significantly better for all languages. We also run Tfidf with stoplist (the results not included here) which achieved results that were better than the plain Tfidf and only slightly inferior to Okapi with stoplist. However, we decided to rather report the original Tfidf, since we wanted to show what is the performance and the relative ranking of the methods without the stoplist.

When comparing the results in the “withnew” and “withoutnew” categories, we see that with stoplist (and Okapi) the addition of the analysis of the new words helps in Finnish, but in German and in English it does not seem to affect the results. Probably this just indicates that in Finnish the vocabulary explosion is more severe and the new corpus introduced a significant amount of important new words. In general, the new words can be analyzed in two different ways: either use the trained analyzer method as such, or train it first with the new words. In this evaluation both ways were actually possible for the participants, and many of them probably already applied the second one.

The Porter stemming that is a standard word preprocessing tool in IR remained unbeaten (by a narrow margin) in our evaluations in English, but in German and especially in Finnish, the unsupervised morpheme analysis methods clearly dominated the evaluation. There might exist better stemming algorithms for those languages, but because of the more complex morphology, their development might not be an easy task.

As future work in this field it should be relatively straight-forward to evaluate the unsupervised morpheme analysis in several other interesting languages, because it is not bounded to only those languages where rule-based grammatical analysis can be performed. It would also be interesting to try to combine the rival analysis to produce something better.

## 9 Conclusions

The objective of Morpho Challenge 2007 was to design a statistical machine learning algorithm that discovers which morphemes (smallest individually meaningful units of language) words consist of. Ideally, these are basic vocabulary units suitable for different tasks, such as text understanding, machine translation, information retrieval, and statistical language modeling. The current challenge was a successful follow-up to our previous Morpho Challenge 2005 (Unsupervised Segmentation of Words into Morphemes). This time the task was more general in that instead of

looking for an explicit segmentation of words, the focus was in the morpheme analysis of the word forms in the data.

The scientific goals of this challenge were to learn of the phenomena underlying word construction in natural languages, to discover approaches suitable for a wide range of languages and to advance machine learning methodology. The analysis and evaluation of the submitted machine learning algorithm for unsupervised morpheme analysis showed that these goals were quite nicely met. There were several novel unsupervised methods that achieved good results in several test languages, both with respect to finding meaningful morphemes and useful units for information retrieval. The IR results also revealed that the morpheme analysis has a significant effect in IR performance in all languages, and that the performance of the best unsupervised methods can be superior to the supervised reference methods.

12 different segmentation algorithms from 6 research groups were submitted and evaluated. The IR evaluations included 3 different languages: Finnish, German and English. The algorithms and results were presented in Morpho Challenge Workshop, arranged in connection with other CLEF 2007 Workshop, September 19-21, 2007. Morpho Challenge 2007 was part of the EU Network of Excellence PASCAL Challenge Program and organized in collaboration with CLEF.

## Acknowledgments

We thank all the participants for their submissions and enthusiasm. We owe great thanks as well to the organizers of the PASCAL Challenge Program and CLEF who helped us organize this challenge and the challenge workshop. Especially, we would like to thank Carol Peters from CLEF for helping us to get Morpho Challenge in CLEF 2007 and organize a great workshop there. We thank also Krista Lagus for comments of the manuscript. Our work was supported by the Academy of Finland in the projects *Adaptive Informatics* and *New adaptive and learning methods in speech recognition*. This work was supported in part by the IST Programme of the European Community, under the PASCAL Network of Excellence, IST-2002-506778. This publication only reflects the authors' views. We acknowledge that access rights to data and other materials are restricted due to other commitments.

## References

- [1] Jeff A. Bilmes and Katrin Kirchhoff. Factored language models and generalized parallel backoff. In *Proceedings of the Human Language Technology, Conference of the North American Chapter of the Association for Computational Linguistics (HLT-NAACL)*, pages 4–6, Edmonton, Canada, 2003.
- [2] Mathias Creutz and Krista Lagus. Unsupervised discovery of morphemes. In *Proceedings of the Workshop on Morphological and Phonological Learning of ACL-02*, pages 21–30, 2002.
- [3] Mathias Creutz and Krista Lagus. Inducing the morphological lexicon of a natural language from unannotated text. In *Proceedings of the International and Interdisciplinary Conference on Adaptive Knowledge Representation and Reasoning (AKRR'05)*, pages 106–113, Espoo, Finland, 2005.
- [4] Mathias Creutz and Krista Lagus. Unsupervised morpheme segmentation and morphology induction from text corpora using Morfessor. Technical Report A81, Publications in Computer and Information Science, Helsinki University of Technology, 2005. URL: <http://www.cis.hut.fi/projects/morpho/>.
- [5] Mathias Creutz and Krister Linden. Morpheme segmentation gold standards for finnish and english. Technical Report A77, Publications in Computer and Information Science, Helsinki University of Technology, 2004. URL: <http://www.cis.hut.fi/projects/morpho/>.

- [6] S. Robertson et al. Okapi at TREC-3. In *Proceedings of the Third Text Retrieval Conference (TREC-3)*, pages 109–126, 1994.
- [7] Mikko Kurimo, Mathias Creutz, and Matti Varjokallio. Unsupervised morpheme analysis evaluation by a comparison to a linguistic Gold Standard – Morpho Challenge 2007. In *Working Notes for the CLEF 2007 Workshop*, Budapest, Hungary, 2007.
- [8] Mikko Kurimo, Mathias Creutz, Matti Varjokallio, Ebru Arisoy, and Murat Saraclar. Unsupervised segmentation of words into morphemes - Challenge 2005, an introduction and evaluation report. In *PASCAL Challenge Workshop on Unsupervised segmentation of words into morphemes*, Venice, Italy, 2006.
- [9] Y.-S. Lee. Morphological analysis for statistical machine translation. In *Proceedings of the Human Language Technology, Conference of the North American Chapter of the Association for Computational Linguistics (HLT-NAACL)*, Boston, MA, USA, 2004.
- [10] M. Porter. An algorithm for suffix stripping. *Program*, 14(3):130–137, July 1980.
- [11] Michael Tepper. *A Hybrid Approach to the Induction of Underlying Morphology*. PhD thesis, University of Washington, 2007.
- [12] Y.L. Ziemann and H.L. Bleich. Conceptual mapping of user’s queries to medical subject headings. In *Proceedings of the 1997 American Medical Informatics Association (AMIA) Annual Fall Symposium*, October 1997.