

Report of NEWS 2010 Transliteration Mining Shared Task

A Kumaran

Microsoft Research India
Bangalore, India

Mitesh M. Khapra

Indian Institute of Technology Bombay
Mumbai, India

Haizhou Li

Institute for Infocomm
Research, Singapore

Abstract

This report documents the details of the Transliteration Mining Shared Task that was run as a part of the Named Entities Workshop (NEWS 2010), an ACL 2010 workshop. The shared task featured mining of name transliterations from the paired Wikipedia titles in 5 different language pairs, specifically, between English and one of Arabic, Chinese, Hindi Russian and Tamil. Totally 5 groups took part in this shared task, participating in multiple mining tasks in different languages pairs. The methodology and the data sets used in this shared task are published in the Shared Task White Paper [Kumaran et al, 2010]. We measure and report 3 metrics on the submitted results to calibrate the performance of individual systems on a commonly available Wikipedia dataset. We believe that the significant contribution of this shared task is in (i) assembling a diverse set of participants working in the area of transliteration mining, (ii) creating a baseline performance of transliteration mining systems in a set of diverse languages using commonly available Wikipedia data, and (iii) providing a basis for meaningful comparison and analysis of trade-offs between various algorithmic approaches used in mining. We believe that this shared task would complement the NEWS 2010 transliteration generation shared task, in enabling development of practical systems with a small amount of seed data in a given pair of languages.

1 Introduction

Proper names play a significant role in Machine Translation (MT) and Information Retrieval (IR) systems. When the systems involve multiple languages, The MT and IR system rely on Machine Transliteration systems, as the proper names are not usually available in standard translation lexicons. The quality of the Machine Transliteration systems plays a significant part in determining the overall quality of the system, and hence, they are critical for most multilingual application systems. The importance of Machine Transliteration systems has been well understood

by the community, as evidenced by significant publication in this important area.

While research over the last two decades has shown that reasonably good quality Machine Transliteration systems may be developed easily, they critically rely on parallel names corpora for their development. The Machine Transliteration Shared Task of the NEWS 2009 workshop (NEWS 2009) has shown that many interesting approaches exist for Machine Transliteration, and about 10-25K parallel names is sufficient for most state of the art systems to provide a practical solution for the critical need. The traditional source for crosslingual parallel data – the bilingual dictionaries – offer only limited support as they do not include proper names (other than ones of historical importance). The statistical dictionaries, though they contain parallel names, do not have sufficient coverage, as they depend on some threshold statistical evidence¹. New names and many variations of them are introduced to the vocabulary of a language every day that need to be captured for any good quality end-to-end system such as MT or CLIR. So there is a perennial need for harvesting parallel names data, to support end-user applications and systems well and accurately.

This is the specific focus of the Transliteration Mining Shared Task in NEWS 2010 workshop (an ACL 2010 Workshop): To mine accurately parallel names from a popular, ubiquitous source, the Wikipedia. Wikipedia exists in more than 250 languages, and every Wikipedia article has a link to an equivalent article in other languages². We focused on this specific resource – the Wikipedia titles in multiple languages and the inter-linking between them – as the source of parallel names. Any successful mining of parallel names from title would signal copious availability of parallel names data, enabling transliteration generation systems in many languages of the world.

¹ In our experiments with Indian Express news corpora over 2 years shows that 80% of the names occur less than 5 times in the *entire* corpora.

² Note that the titles contain concepts, events, dates, etc., in addition to names. Even when the titles are names, parts of them may not be transliterations.

2 Transliteration Mining Shared Task

In this section, we provide details of the shared task, and the datasets used for the task and results evaluation.

2.1 Shared Task: Task Details

The task featured in this shared task was to develop a mining system for identifying single word transliteration pairs from the standard inter-linked Wikipedia topics (aka, Wikipedia Inter-Language Links, or WIL³) in one or more of the specified language pairs. The WIL’s link articles on the same topic in multiple languages, and are traditionally used as a parallel language resource for many natural language processing applications, such as Machine Translation, Crosslingual Search, *etc.* Specific WIL’s of interest for our task were those that contained proper names – either wholly or partly – which can yield rich transliteration data.

The task involved transliteration mining in the language pairs summarized in Table 1.

Source Language	Target Language	Track ID
English	Chinese	WM-EnCn
English	Hindi	WM-EnHi
English	Tamil	WM-EnTa
English	Russian	WM-EnRu
English	Arabic	WM-EnAr

Table 1: Language Pairs in the shared task

Each WIL consisted of a topic in the source and target language pair, and the task was to identify parts of the topic (in the respective language titles) that are transliterations of each other. A seed data set (of about 1K transliteration pairs) was provided for each language pair, and was the only resource to be used for developing a mining system. The participants were expected to produce a paired list of source-target single word named entities, for every WIL provided. At the evaluation time, a random subset of WIL’s (about 1K WIL’s) in each language pair were hand labeled, and used to test the results produced by the participants.

Participants were allowed to use only the 1K seed data provided by the organizers to produce “standard” results; this restriction is imposed to provide a meaningful way of comparing the ef-

³ Wikipedia’s Interlanguage Links: http://en.wikipedia.org/wiki/Help:Interlanguage_links

fective methods and approaches. However, “non-standard” runs were permitted where participants were allowed to use more seed data or any language-specific resource available to them.

2.2 Data Sets for the Task

The following datasets were used for each language pair, for this task.

Training Data	Size	Remarks
Seed Data (Parallel names)	~1K	Paired names between source and target languages.
To-be-mined Wikipedia Inter-Wiki-Link Data (Noisy)	Variable	Paired named entities between source and target languages obtained directly from Wikipedia
Test Data	~1K	This was a subset of Wikipedia Inter-Wiki-Link data, which was hand labeled for evaluation.

Table 2: Datasets created for the shared task

The first two sets were provided by the organizers to the participants, and the third was used for evaluation.

Seed transliteration data: In addition we provided approximately 1K parallel names in each language pair as seed data to develop any methodology to identify transliterations. For standard run results, only this seed data was to be used, though for non-standard runs, more data or other linguistics resources were allowed.

English Names	Hindi Names
village	विलेज
linden	लिनडन
market	मार्केट
mysore	मैसूर

Table 3: Sample English-Hindi seed data

English Names	Russian Names
gregory	Григорий
hudson	Гудзон
victor	Виктор
baranowski	барановский

Table 4: Sample English-Russian seed data

To-Mine-Data WIL data: All WIL’s were extracted from the Wikipedia around January 2010,

and provided to the participants. The extracted names were provided *as-is*, with no hand verification about their correctness, completeness or consistency. As sample of the WIL data for English-Hindi and English-Russian is shown in Tables 5 and 6 respectively. Note that there are 0, 1 or more single-word transliterations from each WIL.

#	English Wikipedia Title	Hindi Wikipedia Title
1	Indian National Congress	भारतीय राष्ट्रीय कांग्रेस
2	University of Oxford	ऑक्सफर्ड विश्वविद्यालय
3	Indian Institute of Science	भारतीय विज्ञान संस्थान
4	Jawaharlal Nehru University	जवाहरलाल नेहरू विश्वविद्यालय

Table 5: English-Hindi Wikipedia title pairs

#	English Wikipedia Title	Russian Wikipedia Title
1	Mikhail Gorbachev	Горбачёв, Михаил Сергеевич
2	George Washington	Вашингтон, Джордж
3	Treaty of Versailles	Версальский договор
4	French Republic	Франция

Table 6: English-Russian Wikipedia title pairs

Test set: We randomly selected ~1000 wikipedia links (from the large noisy Inter-wiki-links) as test-set, and manually extracted the single word transliteration pairs associated with each of these WILs. Please note that a given WIL can provide 0, 1 or more single-word transliteration pairs. To keep the task simple, it was specified that only those transliterations would be considered correct that were clear transliterations word-per-word (morphological variations one or both sides are not considered transliterations) These 1K test set was be a subset of Wikipedia data provided to the user. The gold dataset is as shown in Tables 7 and 8.

WIL#	English Names	Hindi Names
1	Congress	कांग्रेस
2	Oxford	ऑक्सफर्ड
3	<Null>	<Null>
4	Jawaharlal	जवाहरलाल
4	Nehru	नेहरू

Table 7: Sample English-Hindi transliteration pairs mined from Wikipedia title pairs

WIL#	English Names	Russian Names
1	Mikhail	Михаил
1	Gorbachev	Горбачёв
2	George	Джордж
2	Washington	Вашингтон
3	Versailles	Версальский
4	<Null>	<Null>

Table 8: Sample English-Russian transliteration pairs mined from Wikipedia title pairs

2.3 Evaluation:

The participants were expected to mine such single-word transliteration data for every specific WIL, though the evaluation was done only against the randomly selected, hand-labeled test set. A participant may submit a maximum of 10 runs for a given language pair (including a minimum of one mandatory “standard” run). There could be more standard runs, without exceeding 10 (including the non-standard runs).

At evaluation time, the task organizers checked every WIL in test set from among the user-provided results, to evaluate the quality of the submission on the 3 metrics described later.

3 Evaluation Metrics

We measured the quality of the mining task using the following measures:

1. Precision_{CorrectTransliterations} (P_{Trans})
2. Recall_{CorrectTransliteration} (R_{Trans})
3. F-Score_{CorrectTransliteration} (F_{Trans}).

Please refer to the following figures for the explanations:

A = True Positives (TP) = Pairs that were identified as "Correct Transliterations" by the participant and were indeed "Correct Transliterations" as per the gold standard

B = False Positives (FP) = Pairs that were identified as "Correct Transliterations" by the participant but they were "Incorrect Transliterations" as per the gold standard.

C = False Negatives (FN) = Pairs that were identified as "Incorrect Transliterations" by the participant but were actually "Correct Transliterations" as per the gold standard.

D = True Negatives (TN) = Pairs that were identified as "Incorrect Transliterations" by the participant and were indeed "Incorrect Transliterations" as per the gold standard.

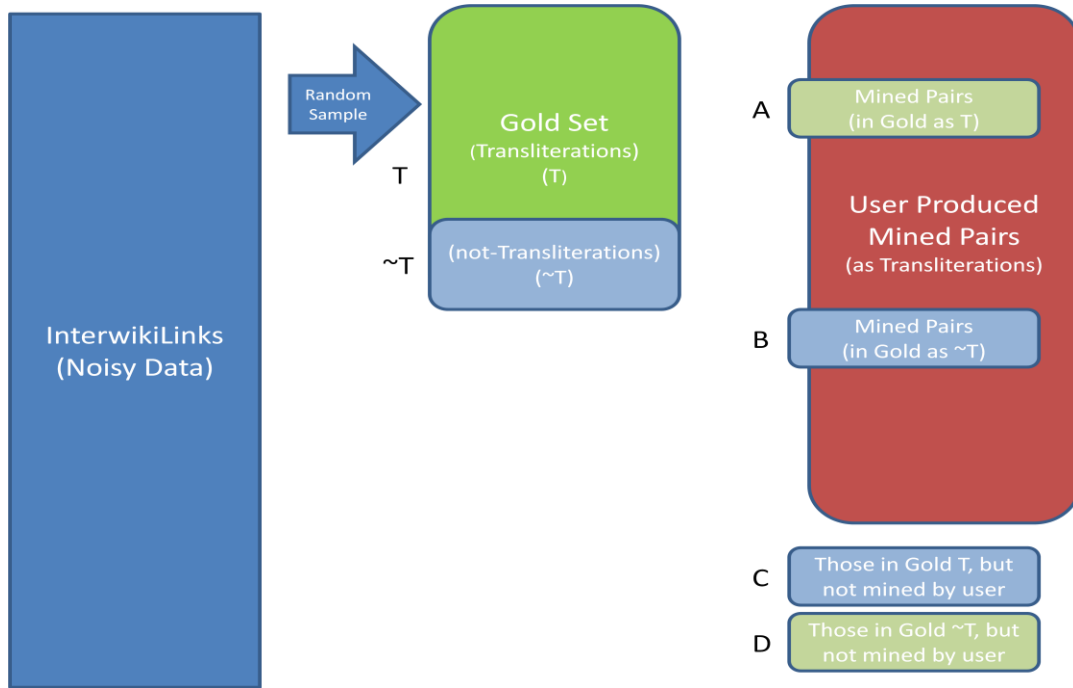


Figure 1: Overview of the mining task and evaluation

1. **Recall**_{CorrectTransliteration} (R_{Trans})

The recall was computed using the sample as follows:

$$R_{Trans} = \frac{TP}{TP + FN} = \frac{A}{A + C} = \frac{A}{T}$$

2. **Precision**_{CorrectTransliteration} (P_{Trans})

The precision was computed using the sample as follows:

$$P_{Trans} = \frac{TP}{TP + FP} = \frac{A}{A + B}$$

3. **F-Score** (F)

$$F = \frac{2 * P_{Trans} * R_{Trans}}{P_{Trans} + R_{Trans}}$$

4 Participants & Approaches

The following 5 teams participated in the Transliteration Mining Task*:

#	Team	Organization
1	Alberta	University of Alberta, Canada
2	CMIC	Cairo Microsoft Innovation Centre, Egypt
3	Groningen	University of Groningen, Netherlands
4	IBM Egypt	IBM Egypt, Cairo, Egypt
5	MINT*	Microsoft Research India, India

* Non-participating system, included for reference.

Table 9: Participants in the Shared Task

The approaches used by the 4 participating groups can be broadly classified as discriminative and generation based approaches. Discriminative approaches treat the mining task as a binary classification problem where the goal is to build a classifier that identifies whether a given pair is a valid transliteration pair or not. Generation based approaches on the other hand generate transliterations for each word in the source title and measure their similarity with the candidate words in the target title. Below, we give a summary of the various participating systems.

The CMIC team (Darwish et. al., 2010) used a generative transliteration model (*HMM*) to transliterate each word in the source title and compared the transliterations with the words appearing in the target title. For example, for a given word E_i in the source title if the model generates a transliteration F_j which appears in the target title then (E_i, F_j) are considered as transliteration pairs. The results are further improved by using phonetic conflation (*PC*) and iteratively training (*IterT*) the generative model using the mined transliteration pairs. For phonetic conflation a modified SOUNDEX scheme is used wherein vowels are discarded and phonetically similar characters are conflated. Both, phonetic conflation and iterative training, led to an increase in

recall which was better than the corresponding decline in precision.

The Alberta team (Jiampoamarn et. al., 2010) fielded 5 different systems in the shared task. The first system uses a simple edit distance based method where a pair of strings is classified as a transliteration pair if the Normalized Edit Distance (*NED*) between them is above a certain threshold. To calculate the NED, the target language string is first Romanized by replacing each target grapheme by the source grapheme having the highest conditional probability. These conditional probabilities are obtained by aligning the seed set of transliteration pairs using an M2M-aligner approach (Jiampoamarn et. al., 2007). The second system uses a SVM based discriminative classifier trained using an improved feature representation (*BK 2007*) (Bergsma and Kondrak, 2007). These features include all substring pairs up to a maximum length of three as extracted from the aligned word pairs. The transliteration pairs in the seed data provided for the shared task were used as positive examples. The negative examples were obtained by generating all possible source-target pairs in the seed data and taking those pairs which are not transliterations but have a longest common subsequence ratio above a certain threshold. One drawback of this system is that longer substrings cannot be used due to the combinatorial explosion in the number of unique features as the substring length increases. To overcome this problem they propose a third system which uses a standard n-gram string kernel (*StringKernel*) that implicitly embeds a string in a feature space that has one coordinate for each unique n-gram (Shawe-Taylor and Cristianini, 2004). The above 3 systems are essentially discriminative systems. In addition, they propose a generation based approach (*DI-RECTL+*) which determines whether the generated transliteration pairs of a source word and target word are similar to a given candidate pair. They use a state-of-the-art online discriminative sequence prediction model based on many-to-many alignments, further augmented by the incorporation of joint n-gram features (Jiampoamarn et. al., 2010). Apart from the four systems described above, they propose an additional system for English Chinese, wherein they formulate the mining task as a matching problem (*Matching*) and greedily extract the pairs with highest similarity. The similarity is calculated using the alignments obtained by training a generation model (Jiampoamarn et. al., 2007) using the seed data.

The IBM Cairo team (Noemans et. al., 2010) proposed a generation based approach which takes inspiration from Phrase Based Statistical Machine Translation (PBSMT) and learns a character-to-character alignment model between the source and target language using GIZA++. This alignment table is then represented using a finite state automaton (*FSA*) where the input is the source character and the output is the target character. For a given word in the source title, candidate transliterations are generated using this FST and are compared with the words in the target title. In addition they also submitted a baseline run which used phonetic edit distance.

The Groningen (Nabende et. al., 2010) team used a generation based approach that uses pair HMMs (*P-HMM*) to find the similarity between a given pair of source and target strings. The proposed variant of pair HMM uses transition parameters that are distinct between each of the edit states and emission parameters that are also distinct. The three edit states are substitution state, deletion state and insertion state. The parameters of the pair HMM are estimated using the Baum-Welch Expectation Maximization algorithm (Baum et. al. 1970).

Finally, as a reference, results of a previously published system – MINT (Udupa et. al., 2009) – were also included in this report as a reference. MINT is a large scalable mining system for mining transliterations from comparable corpora, essentially multilingual news articles in the same timeline. While MINT takes a two step approach – first aligning documents based on content similarity, and subsequently mining transliterations based on a name similarity model – for this task, only the transliteration mining step is employed. For mining transliterations a logistic function based similarity model (*LFS*) trained discriminatively with the seed parallel names data was employed. It should be noted here that the MINT algorithm was used *as-is* for mining transliterations from Wikipedia paired titles, with no fine-tuning. While the standard runs used only the data provided by the organizers, the non-standard runs used about 15K (*Seed+*) parallel names between the languages.

5 Results & Analysis

The results for EnAr, EnCh, EnHi, EnRu and EnTa are summarized in Tables 10, 11, 12, 13 and 14 respectively. The results clearly indicate that there is no single approach which performs well across all languages. In fact, there is even

no single genre (discriminative v/s generation based) which performs well across all languages. We, therefore, do a case by case analysis of the results and highlight some important observations.

- The discriminative classifier using string kernels proposed by Jiampojarn *et. al.* (2010) consistently performed well in all the 4 languages that it was tested on. Specifically, it gave the best performance for EnHi and EnTa.
- The simple discriminative approach based on Normalized Edit Distance (NED) gave the best result for EnRu. Further, the authors report that the results of StringKernel and BK-2007 were not significantly better than NED.
- The use of phonetic conflation consistently performed better than the case when phonetic conflation was not used.
- The results for EnCh are significantly lower when compared to the results for other language pairs. This shows that mining transliteration pairs between alphabetic languages (EnRu, EnAr, EnHi, EnTa) is relatively easier as compared to the case when one of the languages is non-alphabetic (EnCh)

6 Plans for the Future Editions

This shared task was designed as a complementary shared task to the popular NEWS Shared Tasks on Transliteration Generation; successful mining of transliteration pairs demonstrated in this shared task would be a viable source for generating data for developing a state of the art transliteration generation system.

We intend to extend the scope of the mining in 3 different ways: (i) extend mining to more language pairs, (ii) allow identification of near transliterations where there may be changes do to the morphology of the target (or the source) languages, and, (iii) demonstrate an end-to-end transliteration system that may be developed starting with a small seed corpora of, say, 1000 paired names.

References

- Baum, L., Petrie, T., Soules, G. and Weiss, N. 1970. *A Maximization Technique Occurring in the Statistical Analysis of Probabilistic Functions of Markov Chains*. In *The Annals of Mathematical Statistics*, 41 (1): 164-171.
- Bergsma, S. and Kondrak, G. 2007. *Alignment Based Discriminative String Similarity*. In *Proceedings of the 45th Annual Meeting of the ACL, 2007*.
- Darwish, K. 2010. *Transliteration Mining with Phonetic Conflation and Iterative Training*. *Proceedings of the 2010 Named Entities Workshop: Shared Task on Transliteration Mining, 2010*.
- Jiampojarn, S., Dwyer, K., Bergsma, S., Bhargava, A., Dou, Q., Kim, M. Y. and Kondrak, G. 2010. *Transliteration generation and mining with limited training resources*. *Proceedings of the 2010 Named Entities Workshop: Shared Task on Transliteration Mining, 2010*.
- Shawe-Taylor, J and Cristianini, N. 2004. *Kernel Methods for Pattern Analysis*. Cambridge University Press.
- Klementiev, A. and Roth, D. 2006. *Weakly supervised named entity transliteration and discovery from multilingual comparable corpora*. *Proceedings of the 44th Annual Meeting of the ACL, 2006*.
- Knight, K. and Graehl, J. 1998. *Machine Transliteration*. Computational Linguistics.
- Kumaran, A., Khapra, M. and Li, Haizhou. 2010. *Whitepaper on NEWS 2010 Shared Task on Transliteration Mining*. *Proceedings of the 2010 Named Entities Workshop: Shared Task on Transliteration Mining, 2010*.
- Nabende, P. 2010. *Mining Transliterations from Wikipedia using Pair HMMs*. *Proceedings of the 2010 Named Entities Workshop: Shared Task on Transliteration Mining, 2010*.
- Noeman, S. and Madkour, A. 2010. *Language independent Transliteration mining system using Finite State Automata framework*. *Proceedings of the 2010 Named Entities Workshop: Shared Task on Transliteration Mining, 2010*.
- Udupa, R., Saravanan, K., Kumaran, A. and Jagarlamudi, J. 2009. *MINT: A Method for Effective and Scalable Mining of Named Entity Transliterations from Large Comparable Corpora*. *Proceedings of the 12th Conference of the European Chapter of Association for Computational Linguistics, 2009*.

Participant	Run Type	Description	Precision	Recall	F-Score
IBM Egypt	Standard	FST, edit distance 2 with normalized characters	0.887	0.945	0.915
IBM Egypt	Standard	FST, edit distance 1 with normalized characters	0.859	0.952	0.903
IBM Egypt	Standard	Phonetic distance, with normalized characters	0.923	0.830	0.874
CMIC	Standard	HMM + IterT	0.886	0.817	0.850
CMIC	Standard	HMM + PC	0.900	0.796	0.845
CMIC	Standard	(HMM + IterT) + PC	0.818	0.827	0.822
Alberta	Non- Standard		0.850	0.780	0.820
Alberta	Standard	BK-2007	0.834	0.798	0.816
Alberta	Standard	NED+	0.818	0.783	0.800
CMIC	Standard	(HMM + PC + IterT) + PC	0.895	0.678	0.771
Alberta	Standard	DirectTL+	0.861	0.652	0.742
CMIC	Standard	HMM	0.966	0.587	0.730
CMIC	Standard	HMM + PC + IterT	0.952	0.588	0.727
IBM Egypt	Standard	FST, edit distance 2 without normalized characters	0.701	0.747	0.723
IBM Egypt	Standard	FST, edit distance 1 without normalized characters	0.681	0.755	0.716
IBM Egypt	Standard	Phonetic distance, without normalized characters	0.741	0.666	0.702

Table 10: Results of the English Arabic task

Participant	Run Type	Description	Precision	Recall	F-Score
Alberta	Standard	Matching	0.698	0.427	0.530
Alberta	Non-Standard		0.700	0.430	0.530
CMIC	Standard	(HMM + IterT) + PC	1	0.030	0.059
CMIC	Standard	HMM + IterT	1	0.026	0.05
CMIC	Standard	HMM + PC	1	0.024	0.047
CMIC	Standard	(HMM + PC + IterT) + PC	1	0.022	0.044
CMIC	Standard	HMM	1	0.016	0.032
CMIC	Standard	HMM + PC + IterT	1	0.016	0.032
Alberta	Standard	DirectTL+	0.045	0.005	0.009

Table 11: Results of the English Chinese task

Participant	Run Type	Description	Precision	Recall	F-Score
MINT*	Non-Standard	LFS + Seed ⁺	0.967	0.923	0.944
Alberta	Standard	StringKernel	0.954	0.895	0.924
Alberta	Standard	NED+	0.875	0.941	0.907
Alberta	Standard	DirectTL+	0.945	0.866	0.904
CMIC	Standard	(HMM + PC + IterT) + PC	0.953	0.855	0.902
Alberta	Standard	BK-2007	0.883	0.880	0.882
CMIC	Standard	(HMM + IterT) + PC	0.951	0.812	0.876
CMIC	Standard	HMM + PC	0.959	0.786	0.864
Alberta	Non-Standard		0.890	0.820	0.860
MINT*	Standard	LFS	0.943	0.780	0.854
MINT*	Standard	LFS	0.946	0.773	0.851

* Non-participating system

CMIC	Standard	HMM + PC + IterT	0.981	0.687	0.808
CMIC	Standard	HMM + IterT	0.984	0.569	0.721
CMIC	Standard	HMM	0.987	0.559	0.714

Table 10: Results of the English Hindi task

Participant	Run Type	Description	Precision	Recall	F-Score
Alberta	Standard	NED+	0.880	0.869	0.875
CMIC	Standard	HMM + PC	0.813	0.839	0.826
MINT*	Non-Standard	LFS + Seed ⁺	0.797	0.853	0.824
Groningen [^]	Standard	P-HMM	0.780	0.834	0.806
Alberta	Standard	StringKernel	0.746	0.889	0.811
CMIC	Standard	HMM	0.868	0.748	0.804
CMIC	Standard	HMM + PC + IterT	0.843	0.747	0.792
Alberta	Non-Standard		0.730	0.870	0.790
Alberta	Standard	DirectL+	0.778	0.795	0.786
CMIC	Standard	HMM + IterT	0.716	0.868	0.785
MINT*	Standard	LFS	0.822	0.752	0.785
CMIC	Standard	(HMM + PC + IterT) + PC	0.771	0.794	0.782
Alberta	Standard	BK-2007	0.684	0.902	0.778
CMIC	Standard	(HMM + IterT) + PC	0.673	0.881	0.763
Groningen	Standard	P-HMM	0.658	0.334	0.444

Table 11: Results of the English Russian task

Participant	Run Type	Description	Precision	Recall	F-Score
Alberta	Standard	StringKernel	0.923	0.906	0.914
MINT*	Non-Standard	LFS + Seed ⁺	0.910	0.897	0.904
MINT*	Standard	LFS	0.899	0.814	0.855
MINT*	Standard	LFS	0.913	0.790	0.847
Alberta	Standard	BK-2007	0.808	0.852	0.829
CMIC	Standard	(HMM + IterT) + PC	0.939	0.741	0.828
Alberta	Non-Standard		0.820	0.820	0.820
Alberta	Standard	DirectL+	0.919	0.710	0.801
Alberta	Standard	NED+	0.916	0.696	0.791
CMIC	Standard	HMM + IterT	0.952	0.668	0.785
CMIC	Standard	HMM + PC	0.963	0.604	0.743
CMIC	Standard	(HMM + PC + IterT) + PC	0.968	0.567	0.715
CMIC	Standard	HMM + PC + IterT	0.975	0.446	0.612
CMIC	Standard	HMM	0.976	0.407	0.575

Table 12: Results of the English Tamil task

* Non-participating system

[^] Post-deadline submission of the participating system