

CSE 726 Spring 2012  
Seminar  
Multilingual Text Mining  
[www.cedar.buffalo.edu/~rohini/cse726](http://www.cedar.buffalo.edu/~rohini/cse726)

Time: Thursday 1:30 – 3:30 pm, Davis 113A  
Registration #: 14035  
Instructor: Rohini K. Srihari

**Description:**

This course will explore various aspects of multilingual text mining ranging from cross-lingual information retrieval to machine translation, and multilingual information extraction. Lectures will focus on state of the art techniques for multilingual text mining and will include material from conference proceedings and journals. Techniques for mining multilingual social media will be a focus, as will the use of multilingual Wikipedia as a critical resource. More information on selected readings will be forthcoming. Specific topics to be covered include:

- Infrastructure for Multilingual Information Systems
  - Unicode, Tokenization, lexicons
  - Language Identification
  - Hadoop infrastructure
- Cross-Lingual Information Retrieval
  - CLEF
  - Indexing, retrieval
- Machine Translation
  - Transliteration
  - Statistical, rule-based translation, phrase translation using comparable corpora
- Multilingual Information Extraction
  - Named entity tagging, relationship extraction
  - NIST/TAC Multilingual Knowledge Base Population (KBP) task
- Multilingual Social media mining
  - Sentiment analysis
  - Trend analysis, meme discovery
  - Code mixing, switching
- Object Ranking: new search paradigm
- Web as Corpus: Multilingual Resource Creation

Students are expected to make a class presentation based on selected papers, work on a group programming project, as well as participate actively in class discussions. Students will be required to submit a written report that describes their programming project and experiments. The final grade will be based on all of the above, i.e. (i) class presentation, (ii) project and demonstration, (iii) final paper.

The project this semester will focus on leveraging Multilingual Wikipedia. The first step involves indexing Wikipedia in English, Chinese, Arabic, Spanish, Tamil etc. using the Hadoop cluster and Lucene/SOLR infrastructure that we have created. Based on this, several tasks will be considered, ranging from (i) multilingual KBP, (ii) multilingual sentiment analysis (by topic), (iii) discovering translations of terminology, phrases by exploiting comparable documents in multilingual Wikipedia, and (iv) trend identification and meme tracking in multilingual social

media.

The required background is some combination of : information retrieval, machine learning, computational linguistics, and general interest in the area.

Prerequisites: machine learning, information retrieval, computational linguistics background helpful

**Instructor:** Rohini K. Srihari, CEDAR (UB Commons), 645-1602  
email: [rohini@cedar.buffalo.edu](mailto:rohini@cedar.buffalo.edu) office hours: TBD

**Course URL:** [www.cedar.buffalo.edu/~rohini/CSE726](http://www.cedar.buffalo.edu/~rohini/CSE726)

### Course Schedule (Tentative)

Week	Readings	Topics	Lecture Dates
1		Introduction to multilingual (ML) information systems, project descriptions	Week of Jan 16
2		Unicode, tokenization, Language identification, language characteristics	Week of Jan 23
3		Cross Lingual IR, CLEF	Week of Jan 30
4		Machine Translation	Week of Feb 6
5		Machine Translation	Week of Feb 13
6		Multilingual Information Extraction: NE	Week of Feb 20
7		ML KBP	Week of Feb 27
8		ML sentiment analysis	Week of March 5
9		SPRING BREAK	Week of March 12
10		Mining Wikipedia for ML Resource Creation	Week of March 19
11		Code Mixing, Code Switching	Week of March 26
12		Object/Entity Ranking	Week of April 2
13		TBD	Week of April 9
14		Project Discussion	Week of April 23
15		Project Discussion	Week of April 30