

Textual Entailment

Mark Sammons, University of Illinois

Idan Szpektor, Yahoo!

V.G. Vinod Vydiswaran, University of Illinois

The NLP and ML communities are rising to grander, larger-scale challenges such as Machine Reading, Learning by Reading, and Learning to Read, challenges requiring deeper and more integrated natural language understanding capabilities.

The task of Recognizing Textual Entailment (RTE) requires automated systems to identify when two spans of text share a common meaning -- for example, that "Alphaville Inc.'s attempted acquisition of Bauhaus led to a jump in both companies' stock prices" entails "Bauhaus' stock rose", but not "Alphaville acquired Bauhaus". This general capability would be a solid proxy for Natural Language Understanding, and has direct relevance to the grand challenges named above. Moreover, it could be used to improve performance in a large range of Natural Language Processing tasks such as Information Extraction, Question Answering, Exhaustive Search, Machine Translation and many others. The operational definition of Textual Entailment used by researchers in the field avoids commitment to any specific knowledge representation, inference method, or learning approach, thus encouraging application of a wide range of techniques to the problem.

Techniques developed for RTE have now been successfully applied in the domains of Question Answering, Relation Extraction, and Machine translation, and RTE systems continue to improve their performance even as the corpora on which they are evaluated (provided first by PASCAL, and now by NIST TAC) have become progressively more challenging. Over the sequence of RTE challenges from PASCAL and NIST TAC, the more successful systems seem to have converged in their overall approach.

The goal of this tutorial is to introduce the task of Recognizing Textual Entailment to researchers from other areas of NLP. We will identify and analyze common inference and learning approaches from a range of the more successful RTE systems, and investigate the role of knowledge resources. We will examine successful applications of RTE techniques to Question Answering and Machine Translation, and identify key research challenges that must be overcome to continue improving RTE systems.

Tutorial Outline

1. Introduction (35 minutes)

Define and motivate the Recognizing Textual Entailment (RTE) task. Introduce the RTE evaluation framework. Define the relationship between RTE and other major NLP tasks. Identify (some of) the semantic challenges inherent in the RTE task, including the introduction of 'contradiction' as an entailment category. Describe the use of RTE components/techniques in Question Answering, Machine Translation, and Relation Extraction.

2. The State of the Art (35 minutes)

Outline the basic structure underlying RTE systems. With reference to recent publications on RTE: cover the range of preprocessing/analysis that may be used; define representations/data structures typically used; outline inference procedures and machine learning techniques. Identify challenging aspects of the RTE problem in the context of system successes and failures.

3. Machine Learning for Recognizing Textual Entailment (35 minutes)

Describe the challenges involved in applying machine learning techniques to the Textual Entailment problem. Describe in more detail the main approaches to inference, which explicitly or implicitly use the concept of alignment. Show how alignment fits into assumptions of semantic compositionality, how it facilitates machine learning approaches, and how it can accommodate phenomena-specific resources. Show how it can be used for contradiction detection.

4. Knowledge Acquisition and Application in Textual Entailment (35 minutes)

Establish the role of knowledge resources in Textual Entailment, and the consequent importance of Knowledge Acquisition. Identify knowledge resources currently used in RTE systems, and their limitations. Describe existing knowledge acquisition approaches, emphasizing the need for learning directional semantic relations. Define suitable representations and algorithms for using knowledge, including context-sensitive knowledge application. Discuss the problem of noisy data, and the prospects for new knowledge resources/new acquisition approaches.

5. Key Challenges for Recognizing Textual Entailment (15 minutes)

Identify the key challenges in improving textual entailment systems: more reliable inputs (when is a solved problem not solved), domain adaptation, missing knowledge, scaling up. The need for a common entailment infrastructure to promote resource sharing and development.

Biographical Information of the Presenters

Mark Sammons
University of Illinois
201 N. Goodwin Ave.
Urbana, IL 61801 USA
Phone: 1-217-265-6759
Email: mssammon@illinois.edu

Mark Sammons is a Principal Research Scientist working with the Cognitive Computation Group at the University of Illinois. His primary interests are in Natural Language Processing and Machine Learning, with a focus on integrating diverse information sources in the context of Textual Entailment. His work has focused on developing a Textual Entailment framework that can easily incorporate new resources; designing appropriate inference procedures for recognizing entailment; and identifying and developing automated approaches to recognize and represent implicit content in natural language text. Mark received his MSC in Computer Science from the University of Illinois in 2004, and his PhD in Mechanical Engineering from the University of Leeds, England, in 2000.

Idan Szpektor
Yahoo! Research, Building 30 Matam Park, Haifa 31905, ISRAEL.
Phone: + 972-74-7924666; Email: idan@yahoo-inc.com

Idan Szpektor is a Research Scientist at Yahoo! Research. His primary research interests are in natural language processing, machine learning and information retrieval. Idan recently submitted his PhD thesis at Bar-Ilan University where he worked on unsupervised acquisition and application of broad-coverage knowledge-bases for Textual Entailment. He has been a main organizer of the second PASCAL Recognizing Textual Entailment Challenge and an advisor for the third RTE Challenge. He served on the program committees of EMNLP and TextInfer and reviewed papers for ACL, COLING and EMNLP. Idan Szpektor received his M.Sc. from Tel-Aviv University in 2005, where he worked on unsupervised knowledge acquisition for Textual Entailment.

V.G.Vinod Vydiswaran
University of Illinois
201 N. Goodwin Ave.
Urbana, IL 61801 USA

Phone: 1-217-333-2584

Email: vgvinodv@illinois.edu

V.G.Vinod Vydiswaran is a 3rd year Ph.D. student in the Department of Computer Science at the University of Illinois at Urbana-Champaign. His research interests include text informatics, natural language processing, machine learning, and information extraction. His work has included developing a Textual Entailment system, and applying Textual Entailment to relation extraction and information retrieval. He received his Masters degree from Indian Institute of Technology Bombay, India in 2004, where he worked on Conditional models for Information Extraction. Later, he worked at Yahoo! Research & Development Center at Bangalore, India, on scaling Information Extraction technologies over the Web.