

Open Information Extraction: the Second Generation

Authors:

Oren Etzioni, Anthony Fader, Janara Christensen, Stephen Soderland, and Mausam

Published In:

International Joint Conference on Artificial Intelligence, 2011



How to Scale IE?

1970s-1980s: heuristic, hand-crafted clues

- Facts from earnings announcements
- **Narrow** domains; **brittle** clues

1990s: IE as supervised learning

“**Mary** was named to the post of **CFO**, succeeding **Joe** who retired abruptly.”



**Does “IE as supervised learning”
scale to reading the Web?**

No.



Critique of IE=supervised learning

- Relation specific
- Genre specific
- Hand-craft training examples

Does not scale to the Web!



Semi-Supervised Learning

- Few hand-labeled examples **per relation!**
- → Limit on the number of relations
- → relations are **pre-specified**
- → **Still does not scale to the Web**



Machine Reading at Web Scale

- A “universal schema” is impossible
- Global consistency is like world peace
- **Ontological “glass ceiling”**
 - Limited vocabulary
 - Pre-determined predicates
 - Swamped by reading at scale!





Motivation

- General purpose
 - hundreds of thousands of relations
 - thousands of domains
- Scalable: computationally efficient
 - huge body of text on Web and elsewhere
- Scalable: minimal manual effort
 - large-scale human input impractical
- Knowledge needs not anticipated in advance
 - rapidly retargetable





Open IE Guiding Principles

- Domain independence
 - Training for each domain/fact type not feasible
- Scalability
 - Ability to process large number of documents fast
- Coherence
 - Readability important for human interactions



Open vs. Traditional IE

	<i>Traditional IE</i>	<i>Open IE</i>
Input:	Corpus + <i>Hand-labeled Data</i>	Corpus + Existing resources
Relations:	<i>Specified in Advance</i>	Discovered Automatically
Complexity:	$O(D * R)$ <i>R relations</i>	$O(D)$ <i>D documents</i>
Output:	<i>relation-specific</i>	Relation-independe nt



TextRunner



First Web-scale, Open IE system (Banko, IJCAI '07)

1,000,000,000 distinct extractions

Peak of 0.9 precision (but low recall)



Demo

- <http://openie.cs.washington.edu>



Outline

Inference

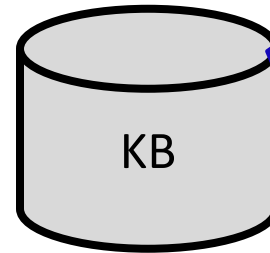


End-user applications



Extraction

Fact




*Downstream
NLP/AI Tasks*





Open Information Extraction

- 2007: Texrunner (~Open IE 1.0)
 - CRF and self-training
- 2010: ReVerb (~Open IE 2.0)
 - POS-based relation pattern
- 2012: OLLIE (~Open IE 3.0)
 - Dep-parse based extraction; nouns; attribution
- 2014: Open IE 4.0
 - SRL-based extraction; temporal, spatial...
- 2016 [@IITD]: Open IE 5.0
 - compound noun phrases, numbers, lists



increasing
precision,
recall,
expressiveness



Fundamental Hypothesis

- \exists *semantically tractable* subset of English
- Characterized relations & arguments via POS
- Characterization is compact, domain independent
- Covers 85% of binary relations in sample



ReVerb

Identify **Relations** from **Verbs**.

1. Find longest phrase matching a simple syntactic constraint:

$$V \mid VP \mid VW^*P$$

V = verb particle? adv?

W = (noun | adj | adv | pron | det)

P = (prep | particle | inf. marker)

Sample of ReVerb Relations

<i>invented</i>	<i>acquired by</i>	<i>has a PhD in</i>
<i>inhibits tumor growth in</i>	<i>voted in favor of</i>	<i>won an Oscar for</i>
<i>has a maximum speed of</i>	<i>died from complications of</i>	<i>mastered the art of</i>
<i>gained fame as</i>	<i>granted political asylum to</i>	<i>is the patron saint of</i>
<i>was the first person to</i>	<i>identified the cause of</i>	<i>wrote the book on</i>



Lexical Constraint

Problem: “overspecified” relation phrases

Obama is offering only modest greenhouse gas reduction targets at the conference.

Solution: must have many distinct args in a large corpus

is offering only modest ...

Obama

the conference

} ≈ 1

is the patron saint of

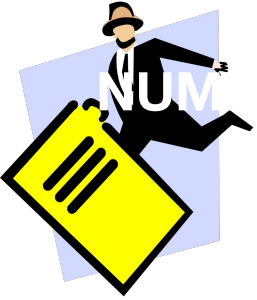
Anne mothers

George England

Hubbins quality footwear

....

100s \approx



Number of Relations

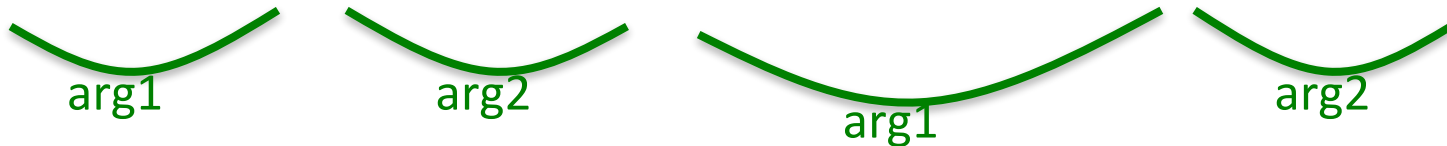
<i>DARPA MR Domains</i>	<i><50</i>
<i>NYU, Yago</i>	<i><100</i>
<i>NELL</i>	<i>~500</i>
<i>DBpedia 3.2</i>	<i>940</i>
<i>PropBank</i>	<i>3,600</i>
<i>VerbNet</i>	<i>5,000</i>
<i>WikiPedia InfoBoxes, $f > 10$</i>	<i>~5,000</i>
<i>TextRunner (phrases)</i>	<i>100,000+</i>
<i>TextRunner (sentences)</i>	<i>1,500,000</i>



ReVerb Extraction Algorithm

1. Identify longest **relation phrases** satisfying constraints

Hudson was born in Hampstead, which is a suburb of London.



2. Heuristically identify **arguments** for relation phrase



(Hudson, was born in, Hampstead)

(Hampstead, is a suburb of, London)



ReVerb Strength

- Outputs more meaningful & informative relations

Homer made a deal with the devil.

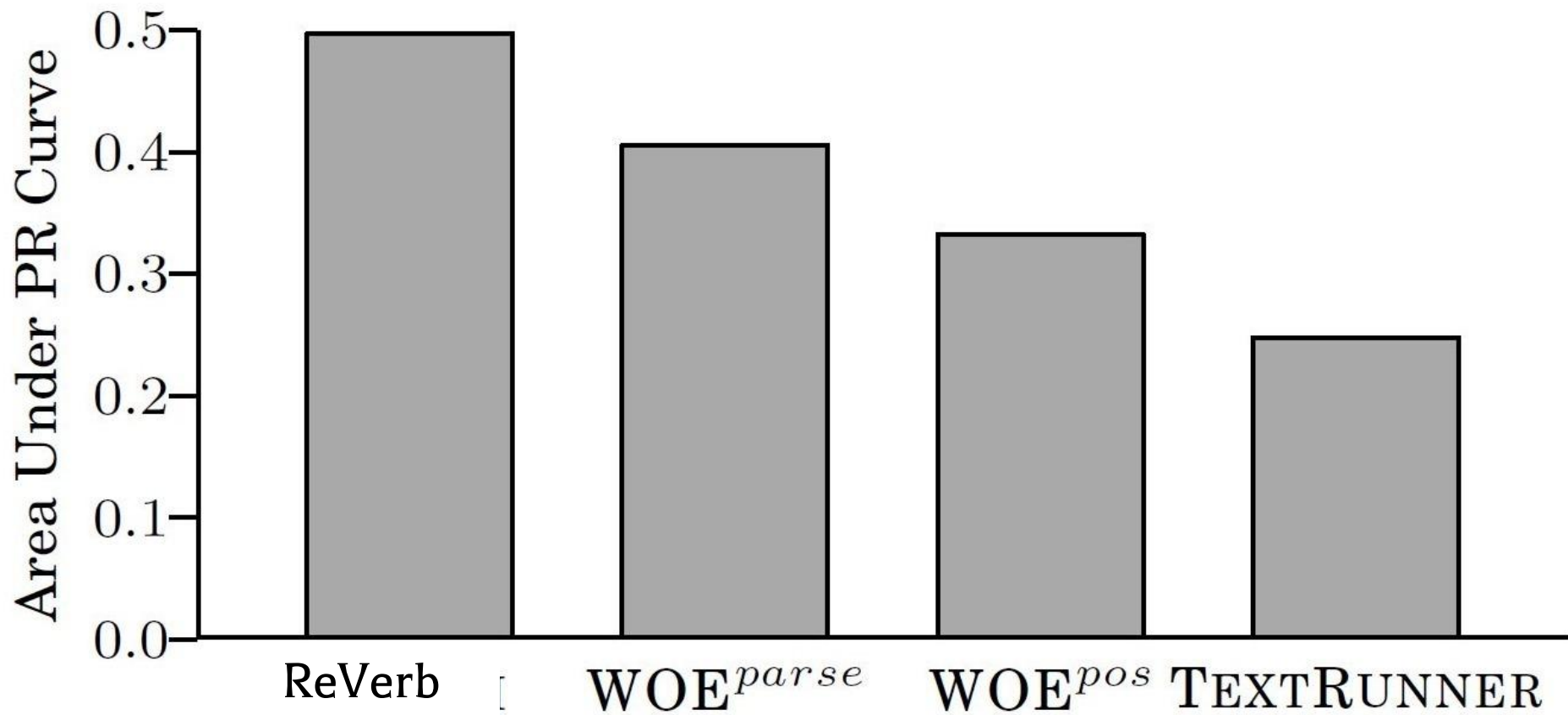


TR → (Homer, made, deal)

RVerb → (Homer, made a deal with, devil)



Experiments: Relation Phrases





ReVerb Error Analysis:

- 65% cases where a relation phrase was correctly identified, but the argument-finding heuristics failed.
- Remaining cases were n-ary relation mistaken as a binary relation.

For eg. extracting (I, gave, him) from the sentence “I gave him 15 photographs”.

- False negatives (52%) were due to the argument-finding heuristics choosing the wrong arguments, or failing to extract all possible arguments



ArgLearner: Motivating Examples

“The assassination of Franz Ferdinand, improbable as it may seem, began WWI.”

(it, began, WWI)

“Republicans in the Senate filibustered an effort to begin debate on the jobs bill.”

(the Senate, filibustered, an effort)

“The plan would reduce the number of teenagers who begin smoking.”

(The plan, would reduce the number of, teenagers)



Analysis – arg1 substructure

<i>Category</i>	<i>Pattern</i>	<i>Freq</i>
Basic Noun Phrases <i>Chicago was founded in 1833</i>	<i>NN, JJ NN, etc</i>	65%
Prepositional Attachments <i>The forest in Brazil is threatened by ranching.</i>	<i>NP PP NP</i>	19%
List <i>Google and Apple are headquartered in Silicon Valley.</i>	<i>NP, (NP,)* CC NP</i>	15%
Relative Clause <i>Chicago, which is located in Illinois, has three million residents.</i>	<i>NP (that WP WDT)? NP? VP NP</i>	<1%



Analysis – arg2 substructure

Category	Pattern	Freq
Basic Noun Phrases <i>Calcium prevents osteoporosis</i>	<i>NN, JJ NN, etc</i>	60%
Prepositional Attachments <i>Barack Obama is one of the presidents of the United States</i>	<i>NP PP NP</i>	18%
List <i>A galaxy consists of stars and stellar remnants</i>	<i>NP, (NP,)* CC NP</i>	15%
Independent Clause <i>Scientists estimate that 80% of oil remains a threat.</i>	<i>(that WP WDT)? NP? VP NP</i>	8%
Relative Clause <i>The shooter killed a woman who was running from the scene.</i>	<i>NP (that WP WDT)? NP? VP NP</i>	6%



Argument Extraction Methodology

- Break problem into four parts:

- Identify arg1 right bound

... TOK TOK TOK TOK TOK rel TOK TOK TOK ...



Classifier (Weka's REPTree)

- Identify arg1 left bound

... TOK TOK TOK TOK TOK rel TOK TOK TOK ...



Classifier (CRF Mallet)

- Identify arg2 left bound

... TOK TOK TOK TOK TOK rel TOK TOK TOK ...



- Identify arg2 right bound

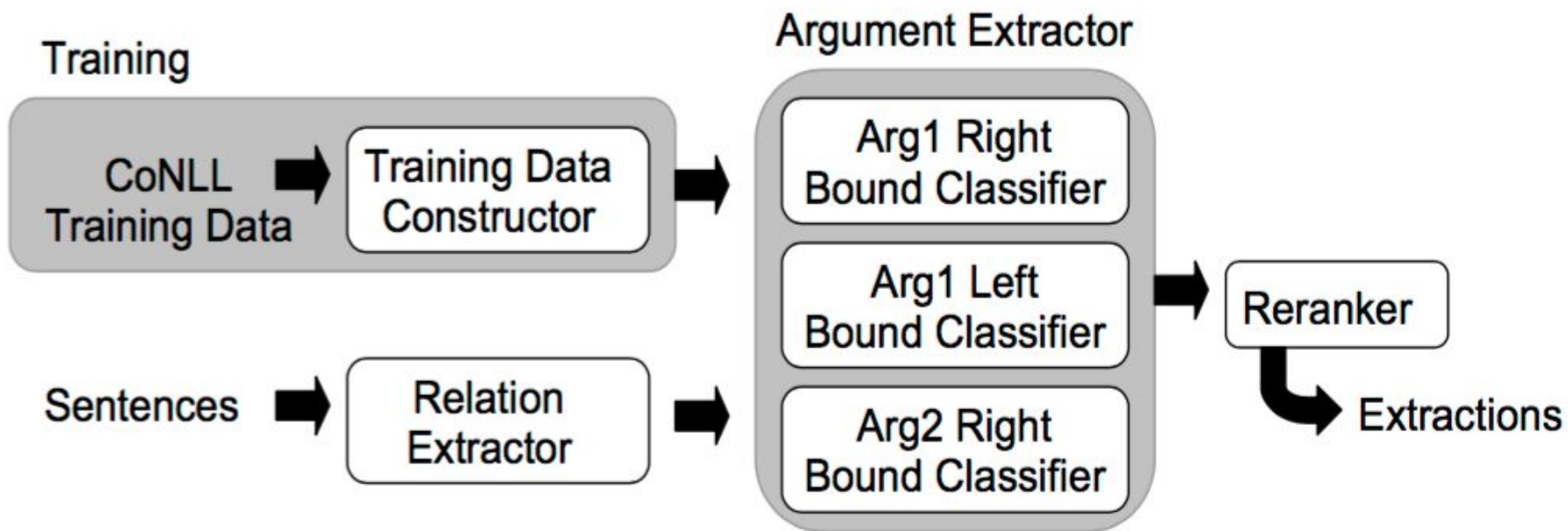
... TOK TOK TOK TOK TOK rel TOK TOK TOK ...



Classifier (CRF Mallet)

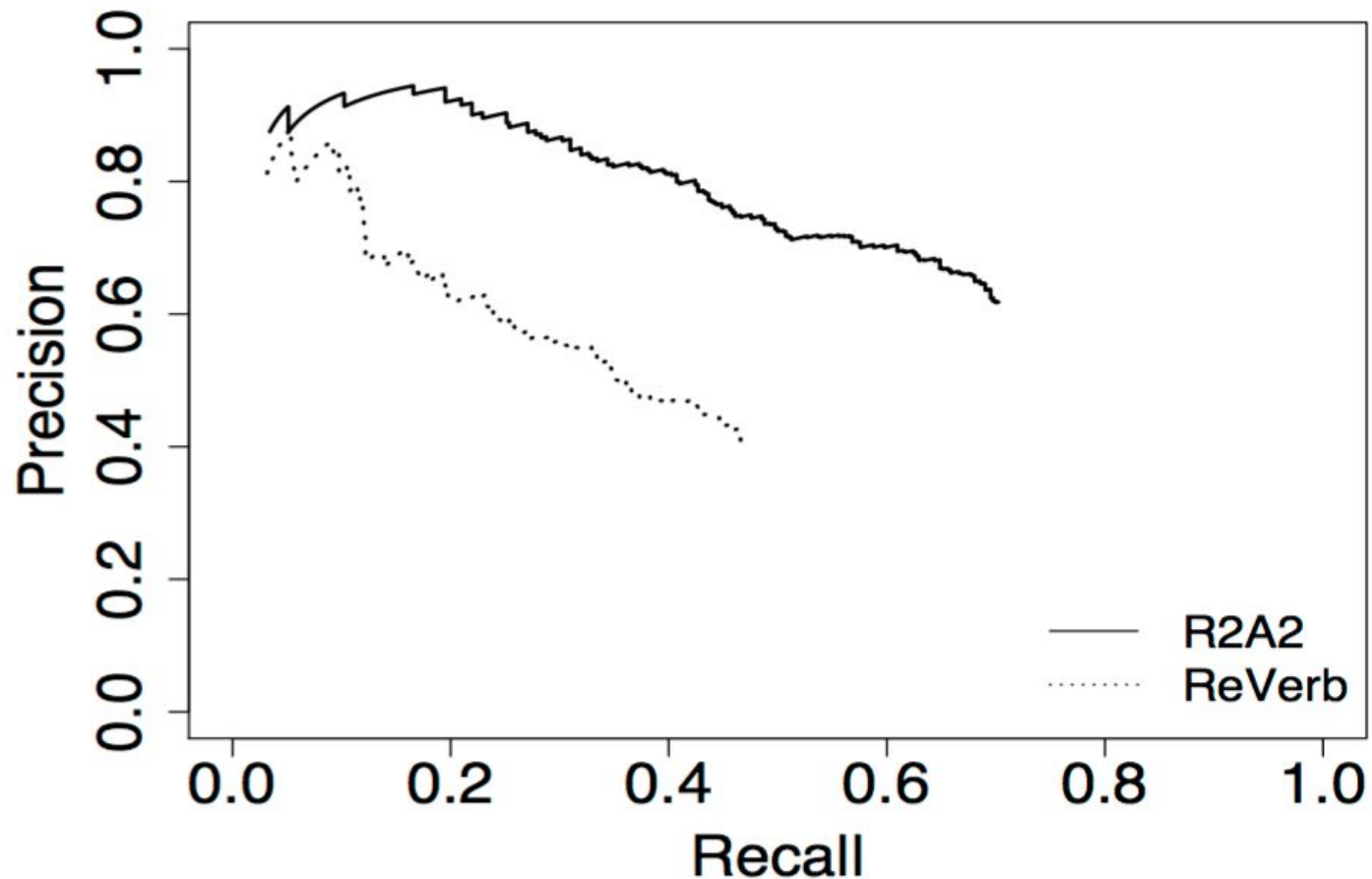


ArgLearner's System Architecture





Evaluation



R2A2 has substantially higher recall and precision than REVERB.



Possible Extension:

“relation discovery from REVERB can be used as a component in NELL to get a NELL-REVERB hybrid that is better at extending its ontology. In contrast to REVERB, NELL has an aspect of temporality and can extract new/update existing entries from an evolving corpus.” - Surag

“Temporality and context not addressed. Ollie incorporates context, but if something was factual at one point but is no longer factual, Ollie will still see it as factual, so temporality needs to be explored.” - Akshay

“Ignores dependency parse information which can be used to provide long range context.” - Akshay

“Many of the observations are for grammatically correct sentences, something which may not be taken for granted in Social Network platforms like Twitter. Extending this method to work on them might be an interesting task” - Barun

“Confidence for extractions could possibly be based on similarity of their Word2Vec vectors” - Gagan

“n-ary relations and relations not limited to verb. (addressed in OPENIE4)
Using more than POS and other syntactic features (SRL used in openIE4)” - Nupur



Thank You!



Error Analysis

