

Supervised Open Information Extraction

Gabriel Stanovsky^{2,3}, Julian Michael², Luke Zettlemoyer², Ido Dagan¹

github.com/gabrielStanovsky/supervised-oie



Open Information Extraction

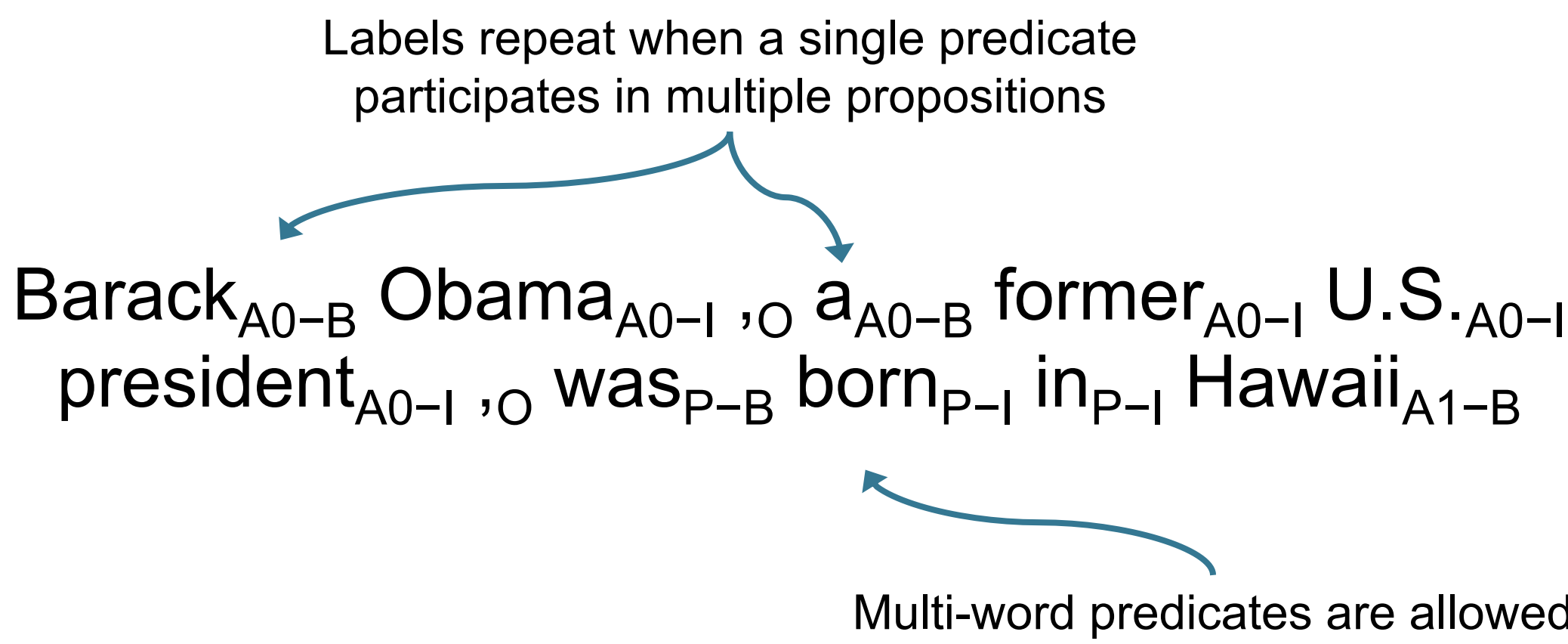
Aims to extract asserted propositions from unstructured text:

“Barack Obama, a former U.S president, **was born in** Hawaii.”

- 1. (Barack Obama; **was born in**; Hawaii)
- 2. (a former U.S. president; **was born in**; Hawaii)

BIO Encoding

Each tuple is encoded with respect to a single predicate, where argument labels indicate their position in the tuple.

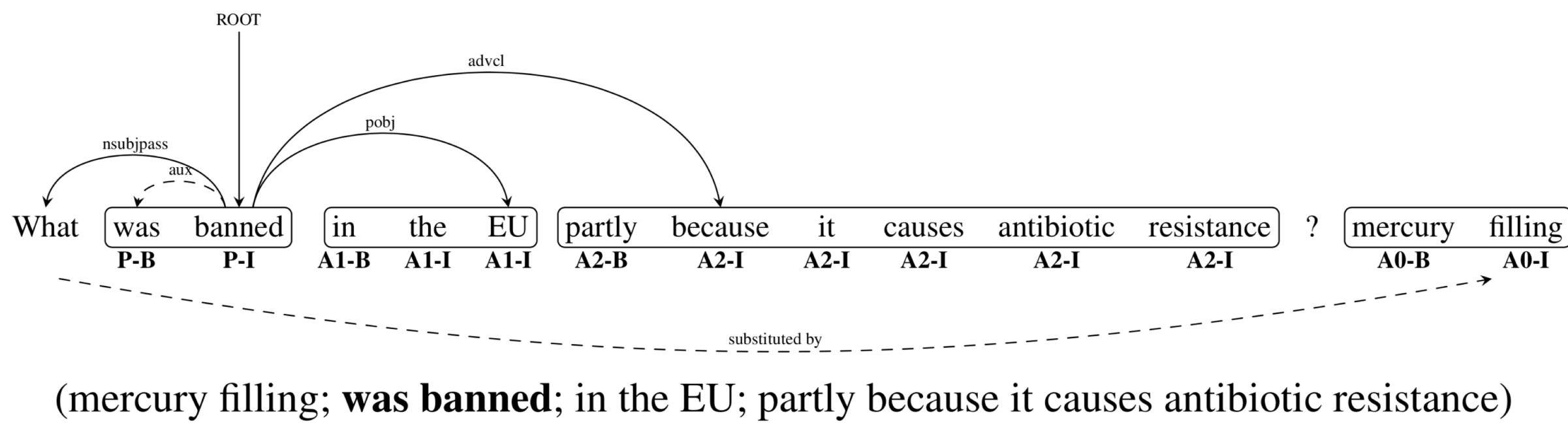


Training Data

We used the QA-SRL to Open IE conversion (**OIE2016**, Stanovsky and Dagan, 2016) to train our model. This consists of verbal propositions, automatically extracted from template QA-SRL annotations.

Augmenting with QAMR annotations

In addition, we converted the Question-Answer Meaning Representation bank (Michael et al, 2018 – **Come see our poster tomorrow!**), consisting of free-form question-answer format over a wide range of predicates. The conversion was achieved with heuristics over the QA parse tree.



Resulting Training Corpus

Dataset	Domain	#Sentences	#Tuples
OIE2016	News, Wiki	3200	5077
QAMR	Wikinews, wiki	3300	12952

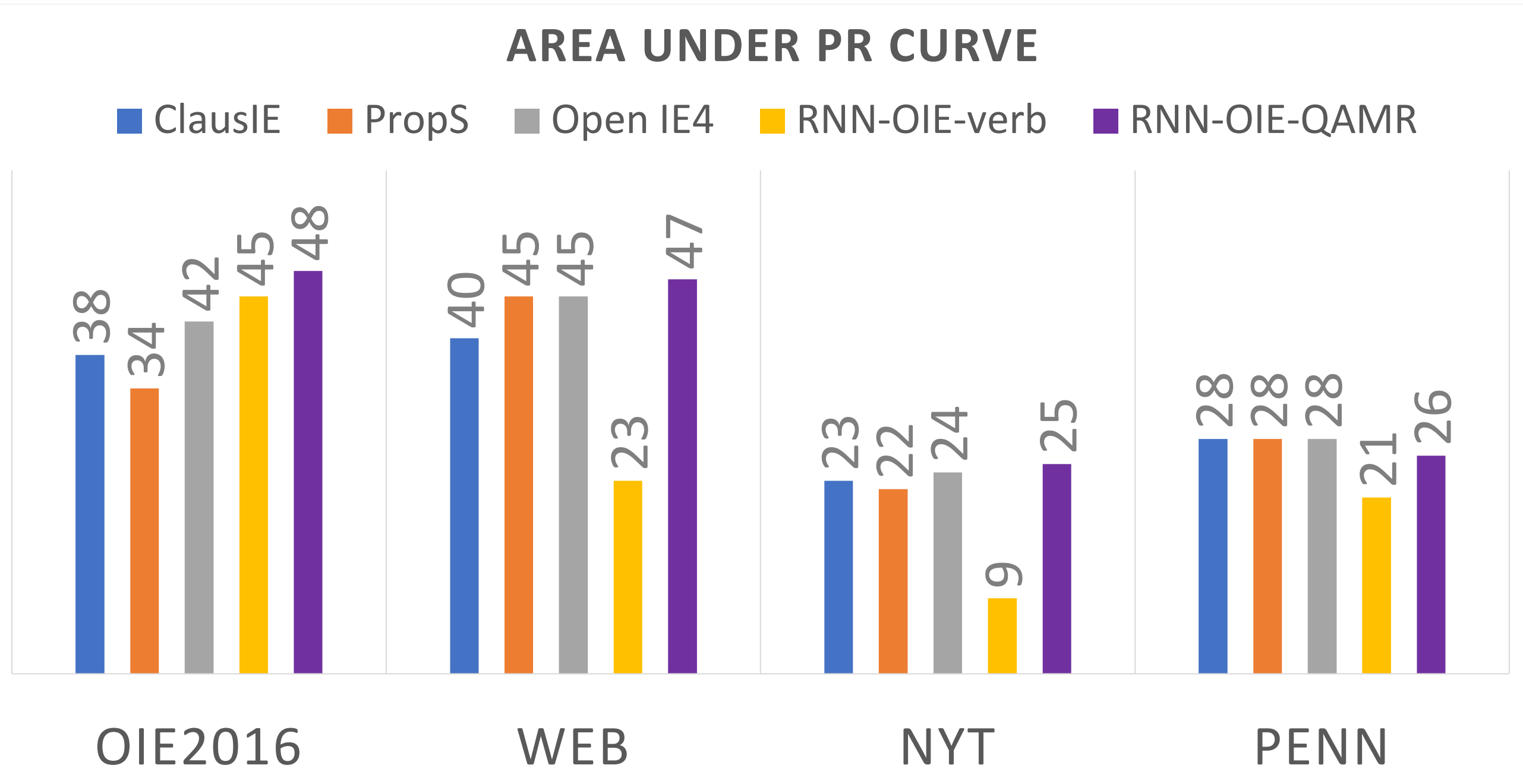
Test Data

We test our model on four publicly available Open IE corpora, following (Schneider et al., 2017).

Dataset	Domain	#Sentences	#Tuples
OIE2016	News, Wiki	3200	1729
WEB	News, Web	500	461
NYT	News, Wiki	222	222
PENN	Mixed	100	51

Evaluation

We compare RNN-OIE against top performing Open IE systems:



RNN-OIE performs competitively across all test sets, outperforming all other systems on the larger test sets. **QAMR improves performance**, especially on more diverse test sets.

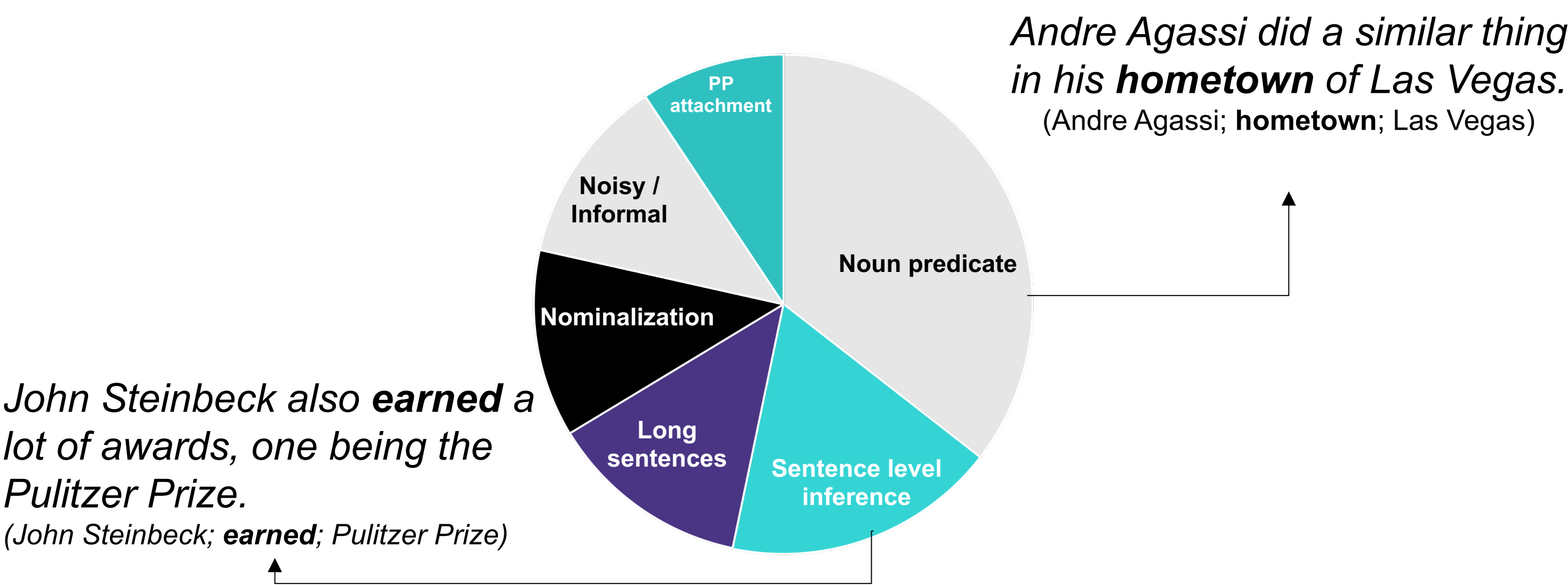
Run-time Analysis

Rnn-OIE is able to leverage GPU architecture to achieve a 10 times improvement over the previous fastest system (measured in sentences per second).

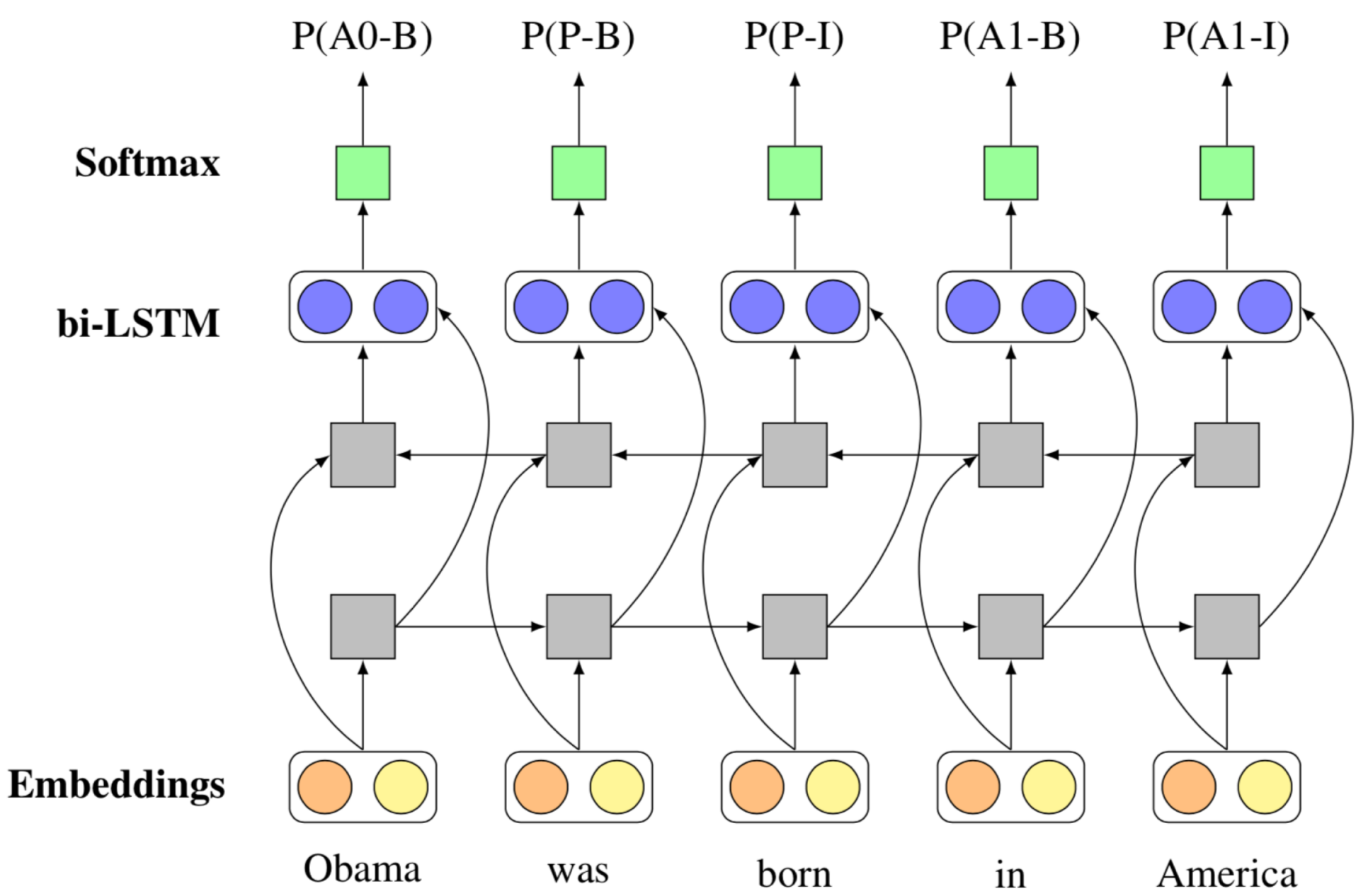
	ClausIE	PropS	Open IE4	RNN-OIE
CPU	4.07	4.59	15.38	13.51
GPU	---	---	---	149.25

Error Analysis

An analysis of 100 gold propositions which were missed by **all** systems (i.e., recall errors) reveals that they all struggle with noun relations, sentence-level inference and long or informal sentences.



RNN-OIE: Bi-LSTM Sequence Tagger Inspired by recent state of the art in Semantic Role Labelling (Zhou and Xu, 2015; He et al., 2017).



Features: Concatenated pretrained embeddings of current word and target predicate (identified by a verb POS).

Decoding: Ignores malformed spans - if an A0-I label is not preceded by A0-I or A0-B, we treat it as O.

Confidence: Estimated for an extraction E by $\prod_{l \in E} P(l)$