



Daniel Khashabi, Mark Sammons, Ben Zhou, Tom Redman, **Christos Christodoulopoulos**, Vivek Srikumar, Nicholas Rizzolo, Lev Ratinov, Guanheng Luo, Quang Do, Chen-Tse Tsai Subhro Roy, Stephen Mayhew, Zhili Feng, John Wieting, Xiaodong Yu, Yangqiu Song, Shashank Gupta, Shyam Upadhyay, Naveen Arivazhagan, Qiang Ning, Shaoshi Ling, Dan Roth

## Components

- Supporting Natural Language Understanding applications requires preprocessing text at multiple, syntactic and semantic, levels.
- The process of managing and aggregating annotations is labor-intensive and error prone, requiring significant engineering.
- It is essential to building software frameworks for easy access to a wide range of NLP annotators and for straightforward use.

The diagram illustrates the CogComp-NLP architecture. It features a central 'Pipeline' box on the right, which is connected to a 'Core-Utilities' box on the left. The 'Pipeline' box contains three main components: 'Chunker', 'POS', and 'Comma-SRL', with an ellipsis indicating additional modules. The 'Core-Utilities' box is connected to 'Similarity Utilities' (top left) and 'Corpus-Utilities' (bottom left). The 'Edison' box is connected to 'Core-Utilities' and 'Comma-SRL'. The 'Annotator Modules' box is connected to 'Core-Utilities' and 'Comma-SRL'. The 'Pipeline' box is connected to 'Core-Utilities' and 'Comma-SRL'.

## Edison

A framework that extracts features to be used by machine learning algorithms. It enables users to define feature extraction functions that take as input the *Views* and *Constituents* created by *Annotators*.

Fundamental data-structures and operators;  
hence many of the other modules depend on it:

- SQL-like operations on *Text-Annotation*
- Experiment utilities & statistical significance
- String pattern-matching algorithms
- Utilities for reading and writing annotations.

NLP corpus readers that populate *Text-Annotation* objects. A few important datasets supported:

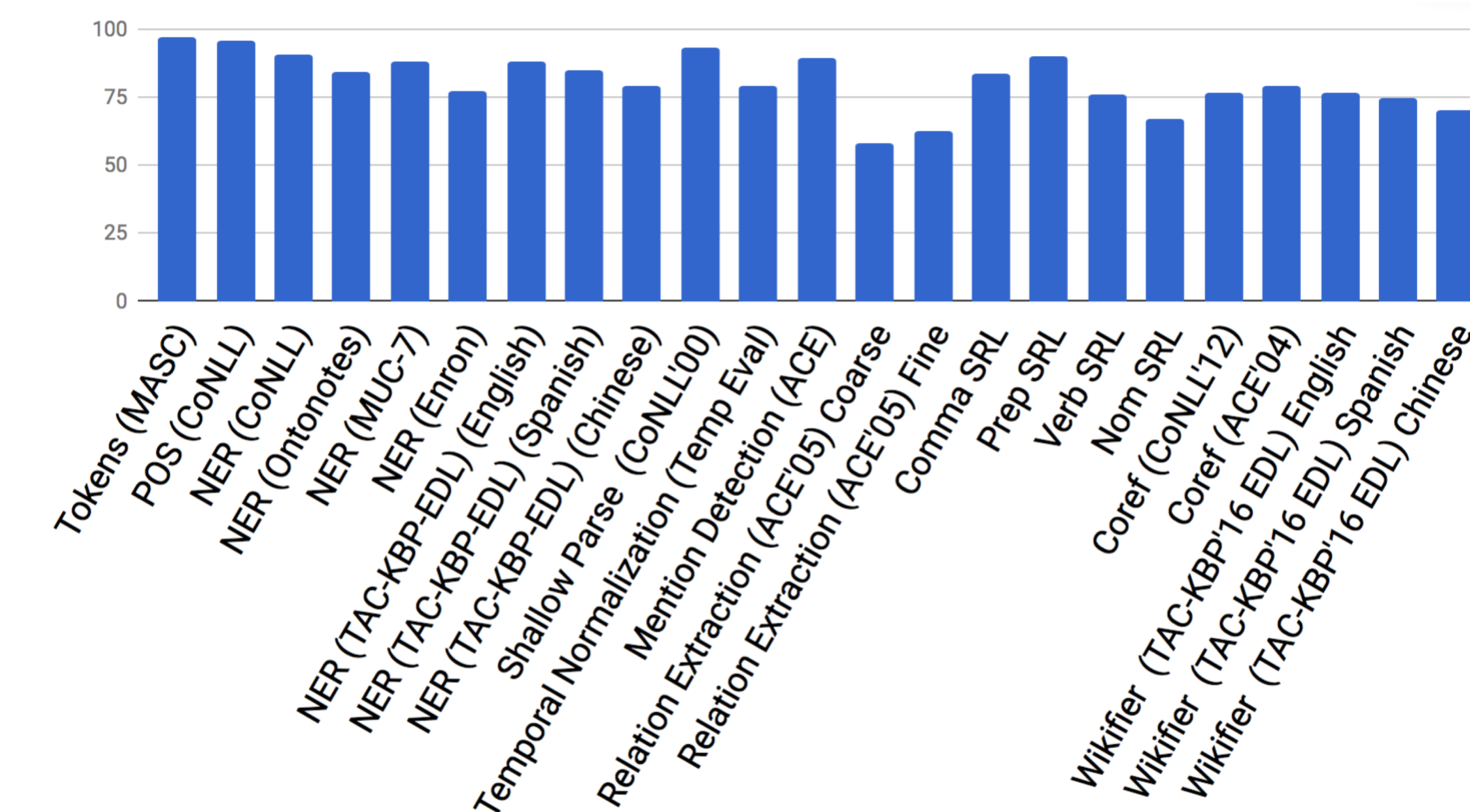
- CoNLL (shallow parsing/chunking)
- PennTreebank (constituency parsing)
- ACE 2004/2005 (NER)
- Ontonotes 5.0

## Similarity Utilities

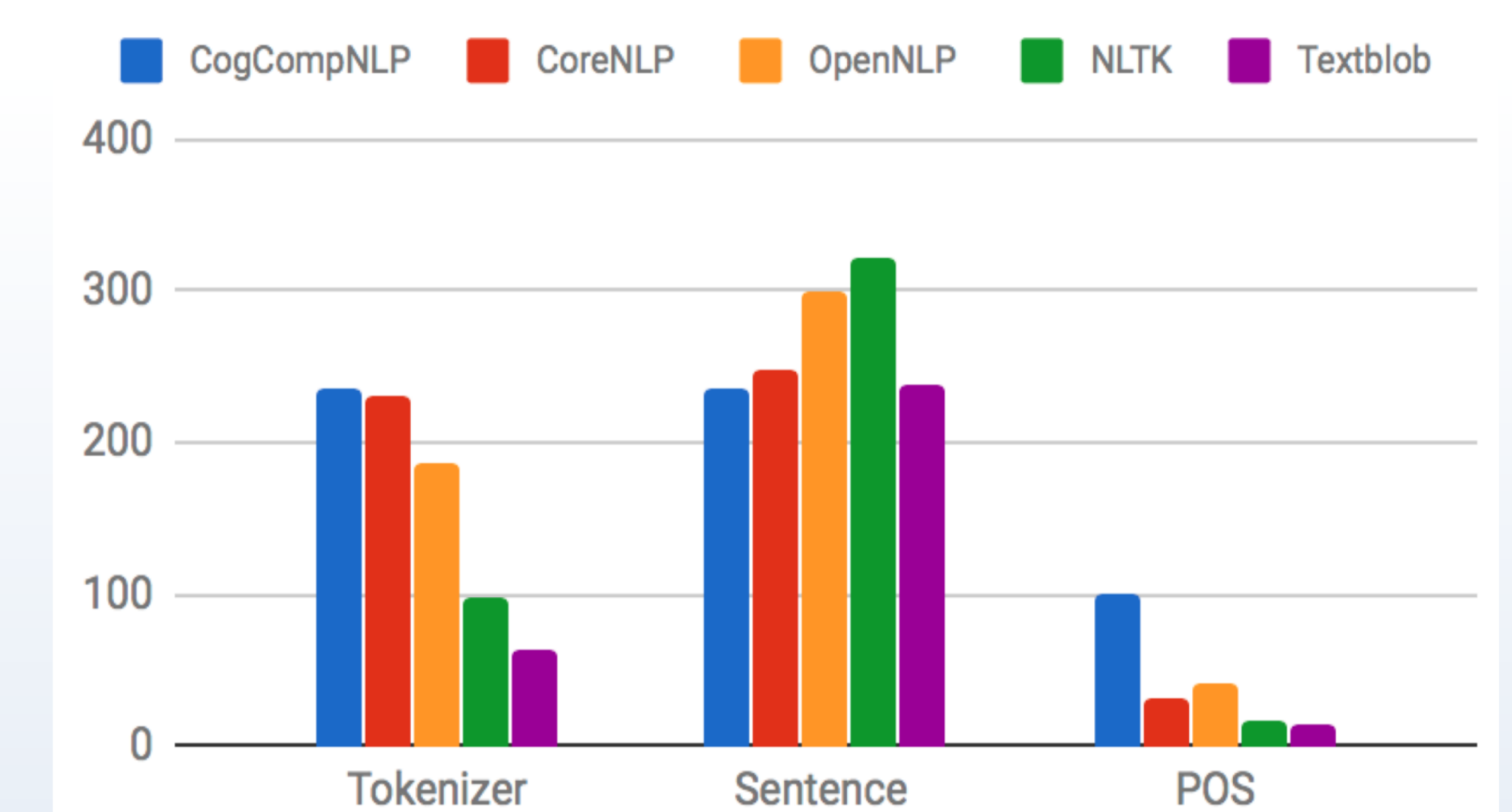
For calculating semantic similarity between words (e.g. Word2Vec, ESA, etc), phrases, and entities.

## Quantitative Evaluation

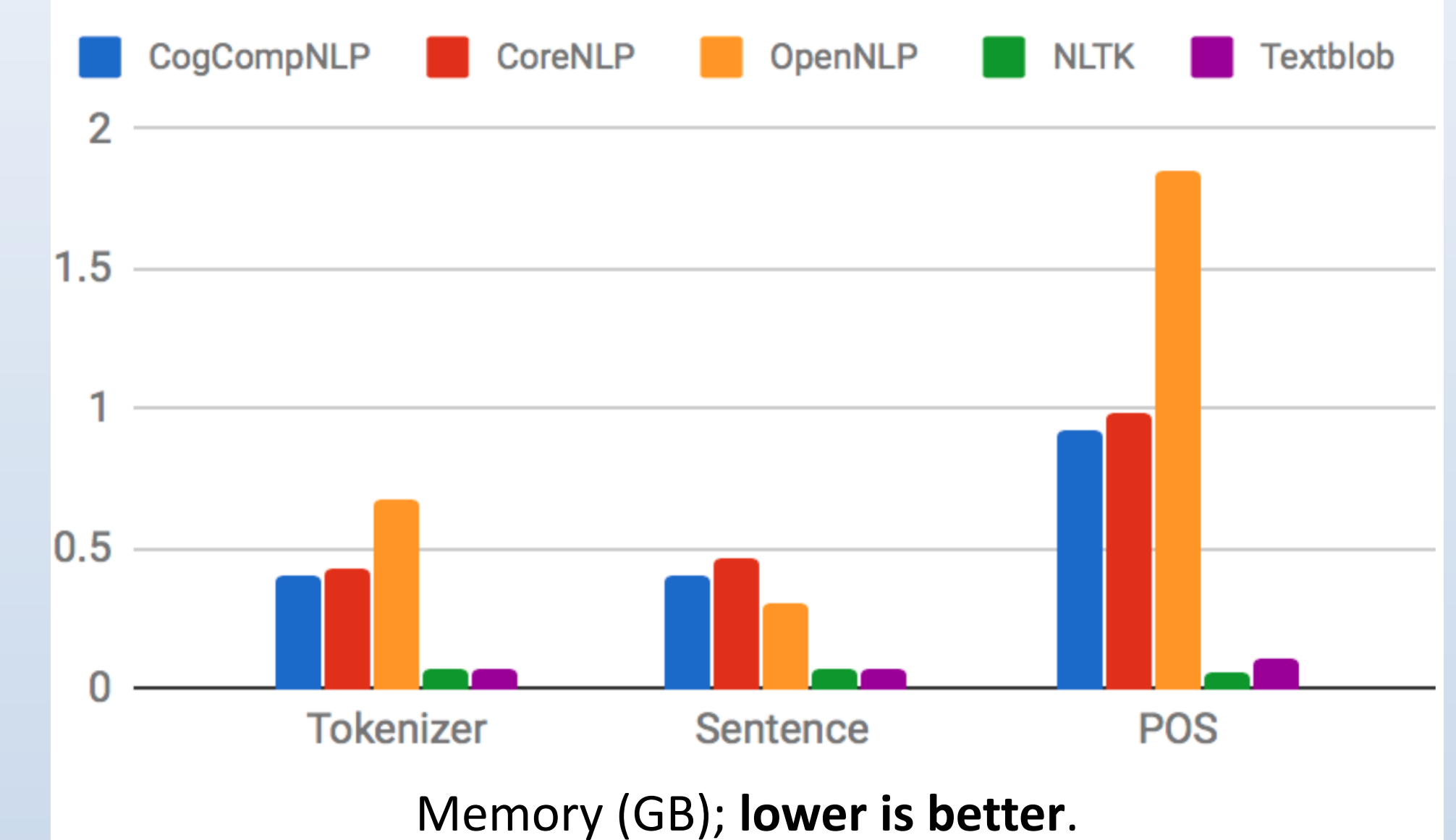
A qualitative assessment of the major components show that they have state-of-the-art quality or very close to the best existing results.



Speed and memory comparison between major NLP pipelines:



Speed: Wall clock time (thousand tokens / second); **higher is better**.



## Acknowledgements

Partly based on research sponsored by

- DARPA under agreement number FA8750-13-2- 0008.
- Contract HR0011-15-2-0025 with the US Defense Advanced Research Projects Agency (DARPA).
- Allen Institute for Artificial Intelligence ([allenai.org](http://allenai.org))
- Google
- NSF grant BCS-1348522; and by NIH grant R01-HD054448.

**TextAnnotation**

**Raw text:** John Smith bought the car.

**Tokens:** {0:John, 1:Smith, 2:bought, 3:the, 4:car, 5:..}

Views		
Name: SENTENCE	Constituents: {...}	Relations: {...}
Name: POS	Constituents: {...}	Relations: {...}

and other views....

**Java code:** <https://github.com/CogComp/cogcomp-nlp>

**Python code:** <https://github.com/CogComp/cogcomp-nlpy>

```
Java
// 'ta' is a partially annotated text
TextAnnotation ta = ...
AnnotatorService pipeline = PipelineFactory.buildPipeline(ViewNames.POS,
                                                         ViewNames.NER_CONLL);
TextAnnotation augTa = pipeline.annotateTextAnnotation(ta);
System.out.println(augTa.getView(ViewNames.POS).getConstituents());
// (NNP Pierre) (NNP Vinken) (, ,) (CD 61) (NNS years)...
```

```
Python
from ccg_nlp import remote_pipeline
pipeline = remote_pipeline.RemotePipeline()
text = "Hello, how are you. I am doing fine"
ta = pipeline.doc(text)
print(ta.get_pos)
# (UH Hello) (, ,) (WRB how) (VBP are) ...
```

**Link to demos:** <http://nlp.cogcomp.org>

