# Project: Coreference Resolution

**Team:** Clever Iguanas

Kaveri Gupta, Sagar Chaturvedi

THE UNIVERSITY OF UTAH

## Objective

In this project, we find antecedents for the anaphors annotated in the input documents. The input files contain anaphora tags and the output files contain the antecedent tag-id as a property of the anaphora tags.

## Methods

This project is a deterministic, rule-based engine for antecedent resolution which is built upon basic properties of the coreferences and the antecedents. E.g. Exact string matching, appositive matching, Head noun matching etc.

## External Tools

The system requires NLP pre-processing to extract word, phrase and sentence properties, like PoSTags, Named entities and Head Nouns, in order to determine the antecedents. To serve this purpose, Stanford CoreNLP library is used.

We have used Tokenizer, Sentence splitter, PoS Tagger, Lemmatizer, Named Entity Tagger, Parser, Dependency Parser and Gender Identifier in a Stanford CoreNLP pipeline.
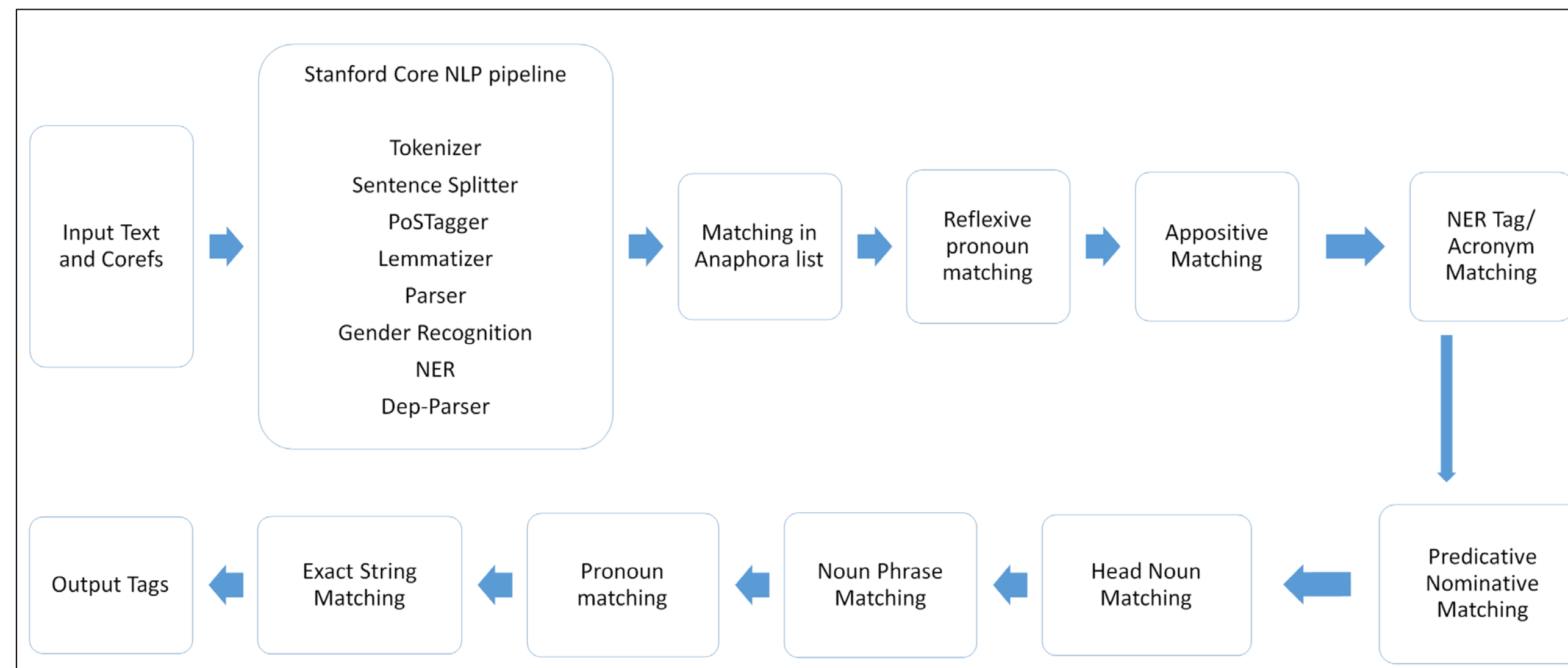
## System

The system is a sequence of sub-modules where each sub-module covers a subset of antecedents in additive manner. It consists of the following high level modules:

☐ **Preprocessing** module reads anaphora tags and text from the input XML files and executes the Stanford NLP pipeline.

☐ **Resolution engine** performs sequential deterministic pattern matching. It matches exact anaphors, reflexive pronouns, appositives, NER Tags and acronyms, predicate nominatives, head nouns, noun phrases, pronouns and exact strings in the given order. We also tried role appositives but it did not perform well as it was tagging wrong antecedents.

Following is the flow diagram of the final system:



## Team Member Contributions

Kaveri worked on the Preprocessing module and the matching of reflexive pronouns, pronoun matching, appositives and predicate nominatives. Rest was implemented by Sagar. The paper reading and the design were done together.

## Results

Following are the results achieved on Test Set #1. Total accuracy was 60.47%. The table mentions all the methods and the accuracies achieved by them when applied individually.

| Over All Accuracy | 60.47 |
|---|---|
| Anaphora Matching | 39.98 |
| Reflexive Pronoun Matching | 0.08 |
| Appositive Matching | 5.52 |
| Named Entity matching | 1.06 |
| Predicate Nominative Matching | 1.21 |
| Matching Head nouns in anaphora tags | 41.42 |
| Noun phrase matching | 35.98 |
| Matching all pronouns | 2.87 |
| Exact string matching outside anaphora tags | 26.46 |

On the Test Set #2, the system achieved 61.83% while on the Test Set #3 it got 61.49%.

## Analysis

While the simpler approaches like anaphora matching, head noun matching and exact string matching performed well, the complex approaches like Predicate nominatives, role appositives and reflexive pronouns did not perform as expected. Given more time, we would have tried adding synonym matching in noun phrases and checking context similarity around the anaphors and the antecedents.

## References

1. Stanford's Multi-Pass Sieve Coreference Resolution System - http://nlp.stanford.edu/pubs/conllst2011-coref.pdf.

2. Coreference Resolution: Current Trends and Future Directions - https://www.cs.cmu.edu/~jhclark/pubs/clark_gonzalez_coreference.pdf