

Overview

Linguine aims to allow language science students to explore automated language analyses such as syntactic or semantic analyses.

This project built upon prior senior projects that created a user interface and web API. This year's project focused on linguistic analysis functionalities as well as system stability and performance.

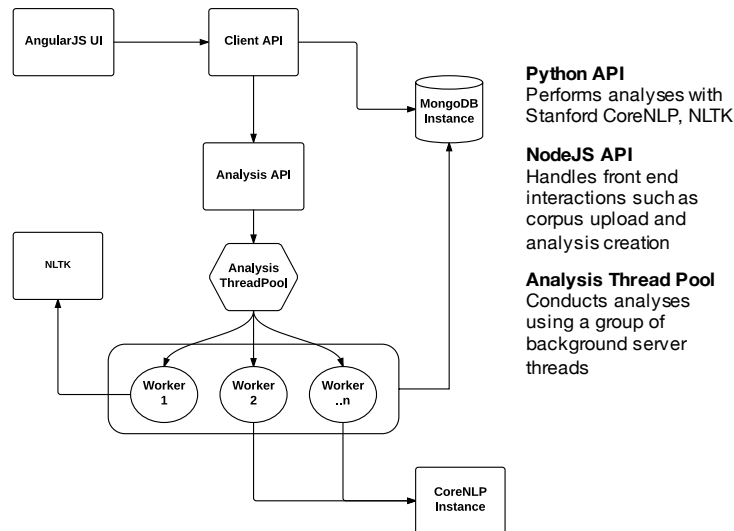
Contributions

- Implemented 5 additional analysis types
- Implemented visualizations for all analysis types
- Made multiple usability improvements to UI
- Implemented concurrency on Python backend
- Integrated Stanford CoreNLP and Illinois Curator
- Enabled users to work while analysis is processing

Technologies



Architecture

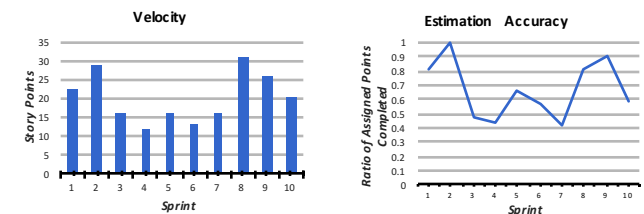


Methodology

We followed a Scrum methodology, using two week sprints. Several of our team members had experience using Scrum and we were likely to do much of the development independently, so it was the most appropriate choice.

In Fall, we met twice a week and remotely on the weekends. In Spring, we met in person three times a week to increase productivity.

We began each sprint with a sprint planning meeting to assign story points to each of our user stories. We used burndown charts to track our velocity and estimation accuracy.



Analyses & Visualizations

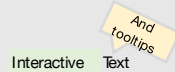
Term Frequency Analysis

Compute word frequencies in a text



Named Entity Recognition

Identify words by classes such as organization, place, or time expression



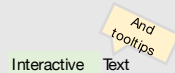
Coreference Resolution

Locate expressions that refer to the same entity in a text



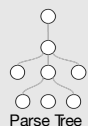
Relation Extraction

Find relationship triples between words



Parsing & Part of Speech Tagging

Construct a dependency parse tree and label words by part of speech



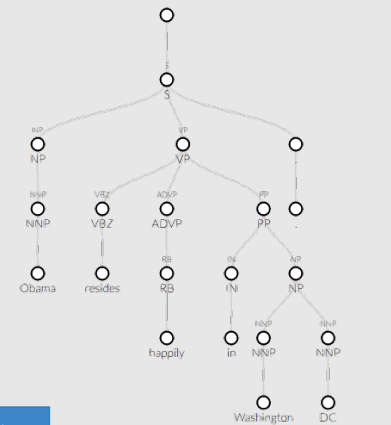
Sentiment Analysis

Estimate the sentiment of a text along with its sentences and tokens



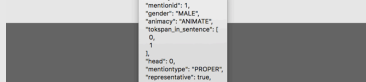
Obama resides happily in Washington DC.

Obama: PERSON



Select Entity

Obama resides happily in Washington DC.



Linguine
Corpora
Analyses
Documentation

Create Analysis
Search by corpus name, analysis type...

2015 News Article - Sentiment (Stanford CoreNLP) INCOMPLETE - Refresh to update

Analysis: nlp-sentiment
Cleanups:
Tokenizer:
Time Created: 4/19/2016 10:25:40 AM
ETA: 10:28:10 AM

2015 News Article - Term Frequency Analysis

Analysis: wordcloudop
Cleanups: stem_porter removecapsnp removepunct
Tokenizer: word_tokenize_treebank
Time Created: 4/19/2016 10:25:07 AM

Lessons Learned

- Consider a distributed computing model when working with computationally expensive operations
- As many system critical bugs as possible should be flushed out prior to user testing
- When inheriting a project, set aside enough time for potential rework or bug fixing
- Sponsor time is valuable - meetings need to focus on what the sponsor needs