# Tale Scanner : Natural Language Processing

## Matteo DeMilia, Kyrstn Hall, Christopher Mead, and Tiera Smith
## Professor: Dr. Jing Hou
## Sponsor: Dr. Darin Woolpert of the University of Vermont

California State University SAN MARCOS

The UNIVERSITY of VERMONT

## Background

- Linguistic analysis is slow and tedious

- Speech language pathologists(SLP) would benefit from efficient analysis tools

- Natural language processing(NLP) can speed up the analysis process exponentially

- The future of linguistic analysis lies in computed automation and artificial intelligence(AI)
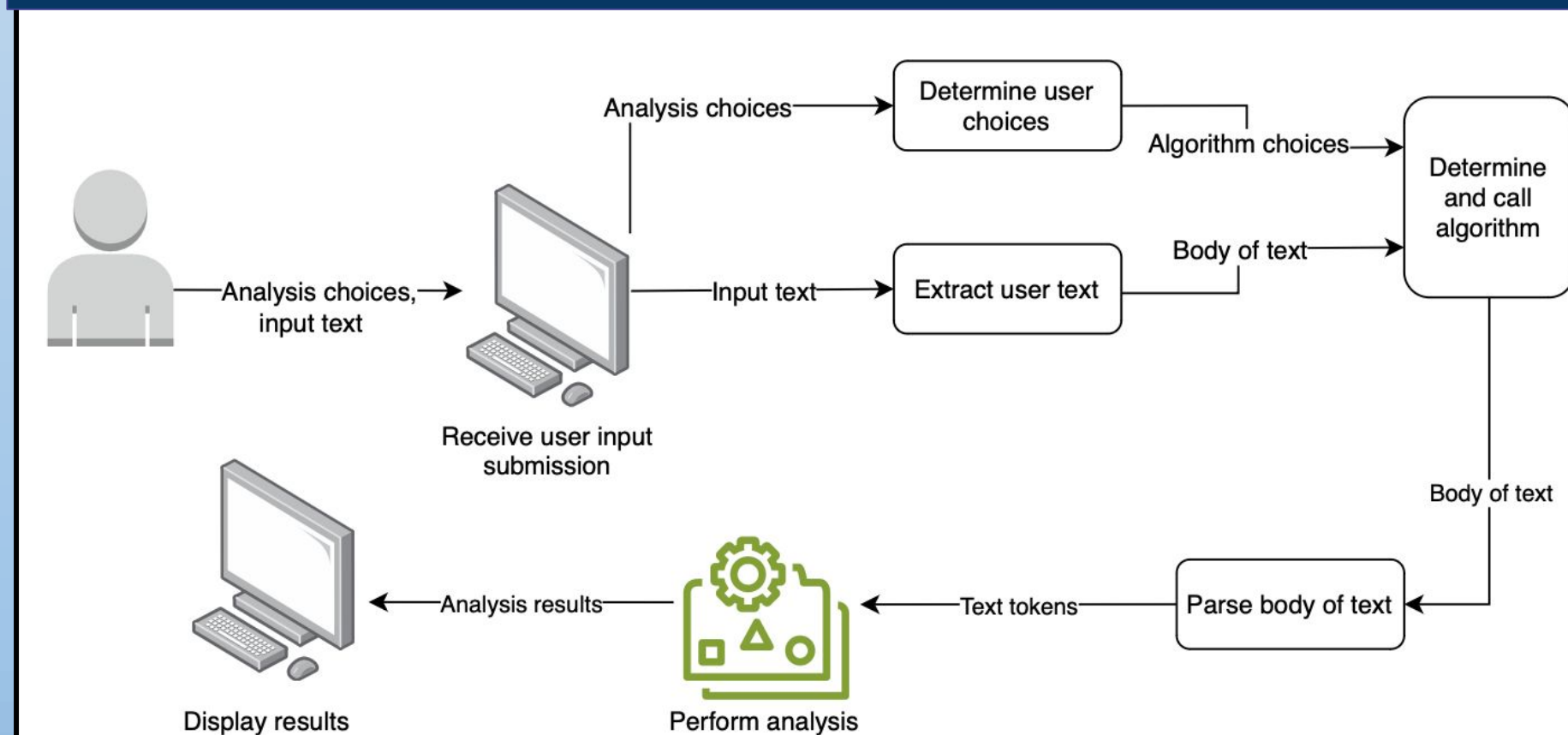
## Why Choose Tale Scanner?

- Tale Scanner is much more efficient than manual linguistic analysis

- Tale Scanner provides consistency, unlike general AI

- Choose from numerous analysis with results provided in <1 min

- Tale Scanner provides detailed description and background on the analysis process of each algorithm

## System Requirements

Goals of the Tale Scanner NLP System:
- Provide a simple user interface supporting:
  - Input text-box for 200-600 words
  - Analysis choice selection
  - Seamless execution and navigation
- Provide Text Analysis of:
  - Total number of words and unique words
  - Type token ratio
  - Total number of clauses
  - Total number of subordinate clauses
  - Syntactic subordination index
  - Total number of verb errors
  - Verb error ratio
  - Total number of morphemes

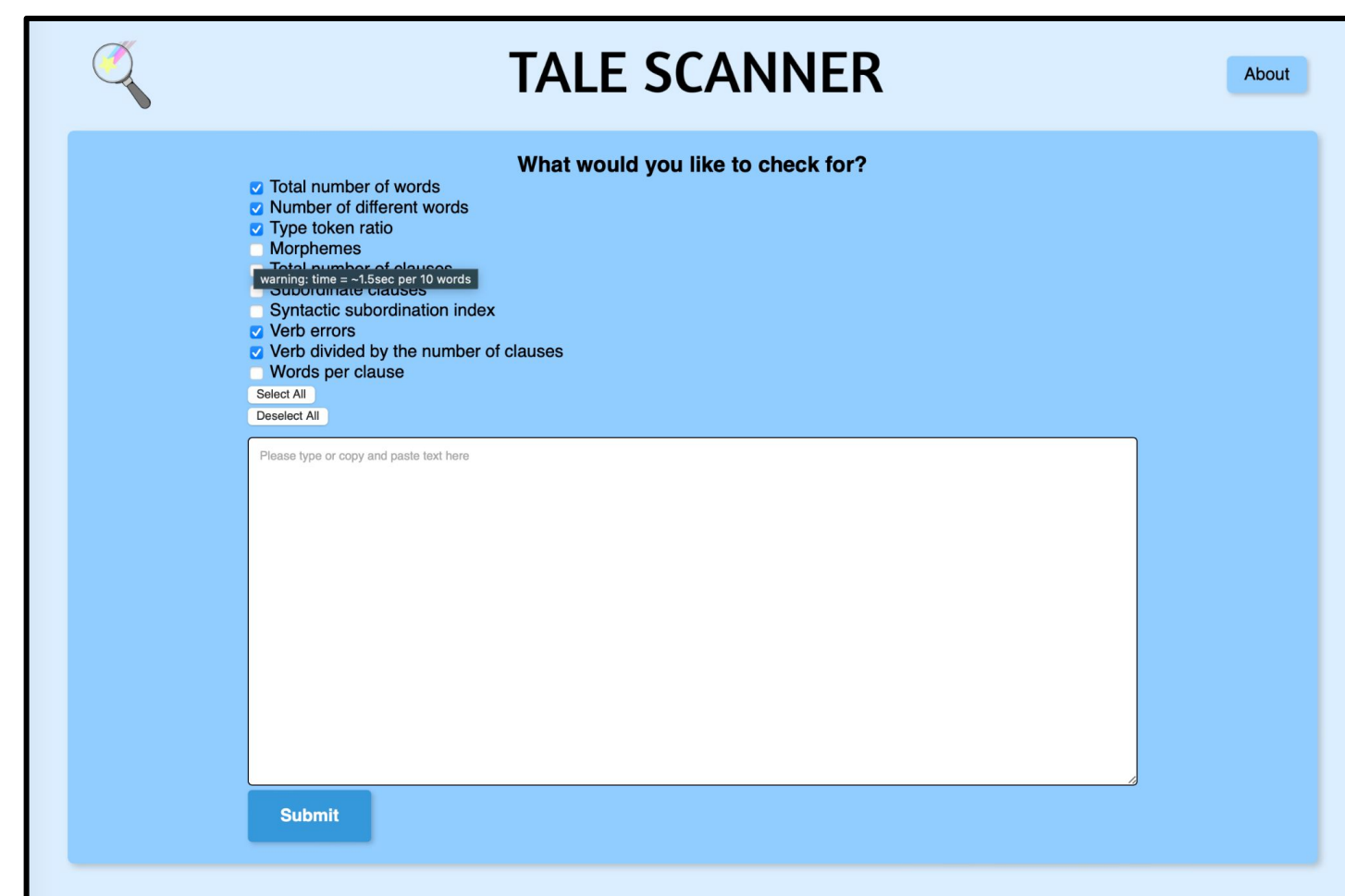## System Architecture and Design



- Simple, linear design; reduces dependencies and error cases

- Designed for efficiency and stability

- Two inputs:
  - Body of text
  - Analysis choices

## Tech Stack

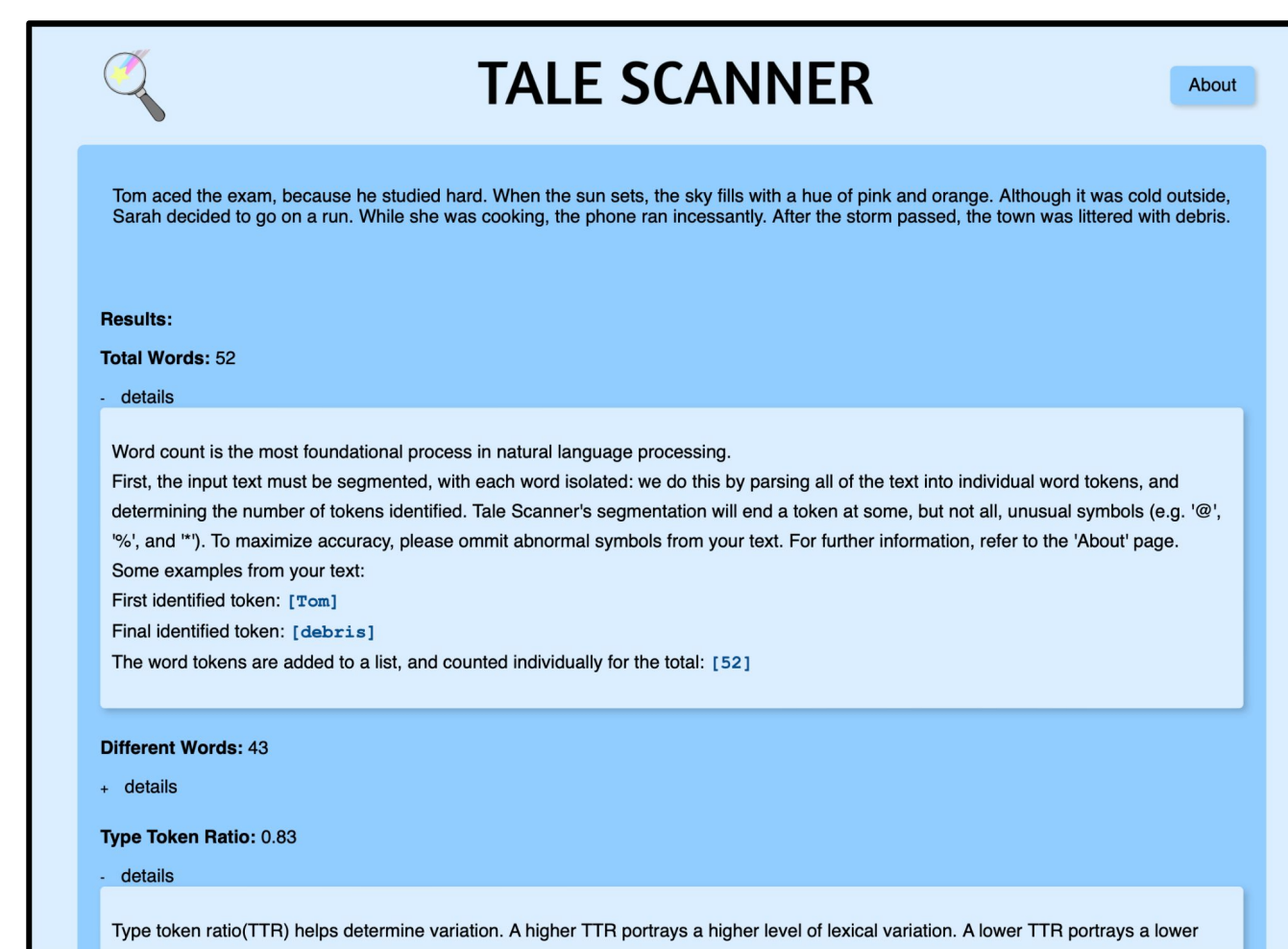HTML   CSS   JavaScript

Python   Flask   AWS Amplify
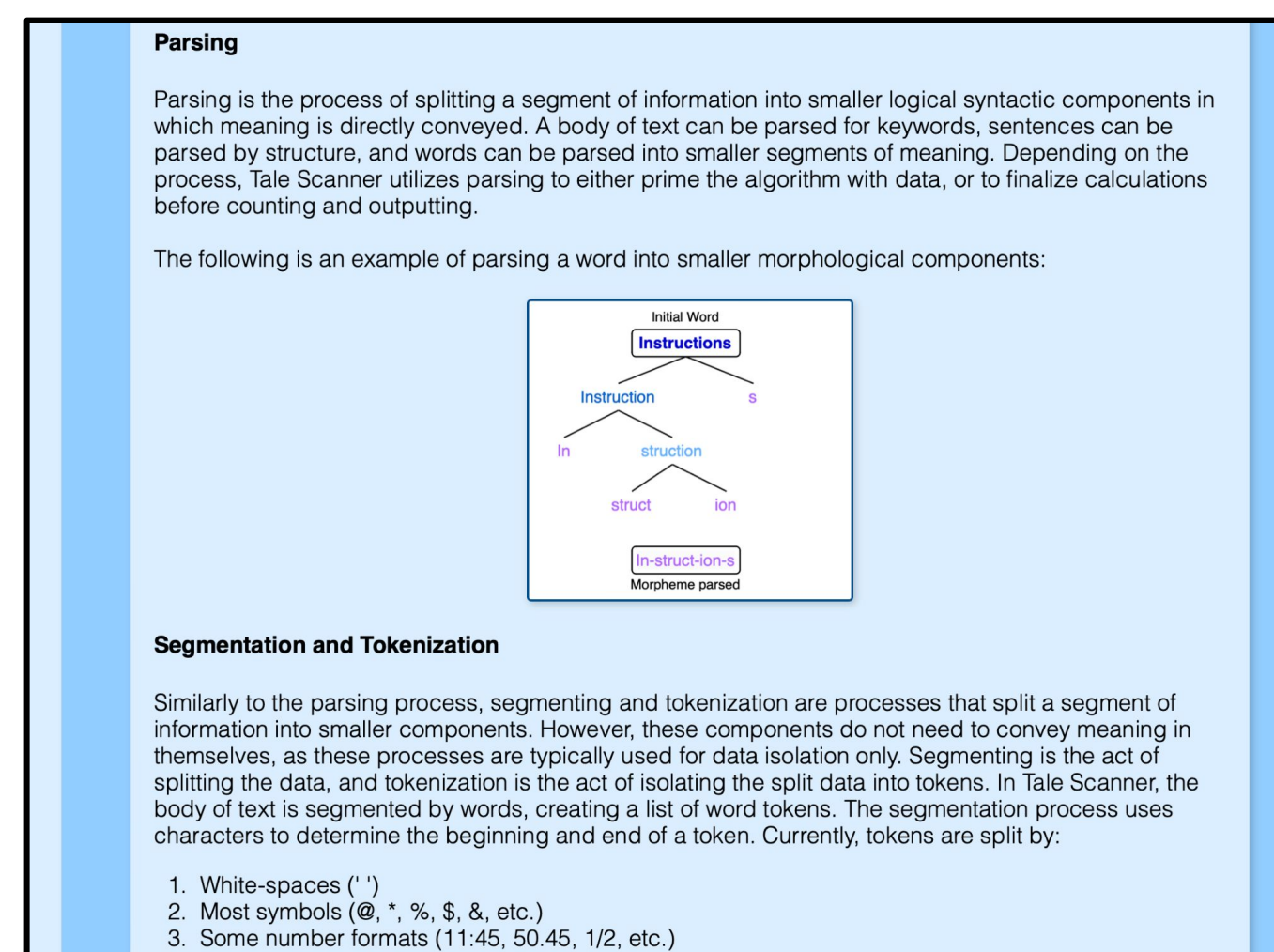
## Tale Scanner's User Interface



- Simple and clean UI
- Spinner submit button
- Error message for input errors
- Saves previously chosen analysis boxes
- Dynamic warnings are provided for a few analysis choices

- Results page loads with all chosen analysis results in drop down menu
- Drop down features include:
  - Dynamic data in blue text
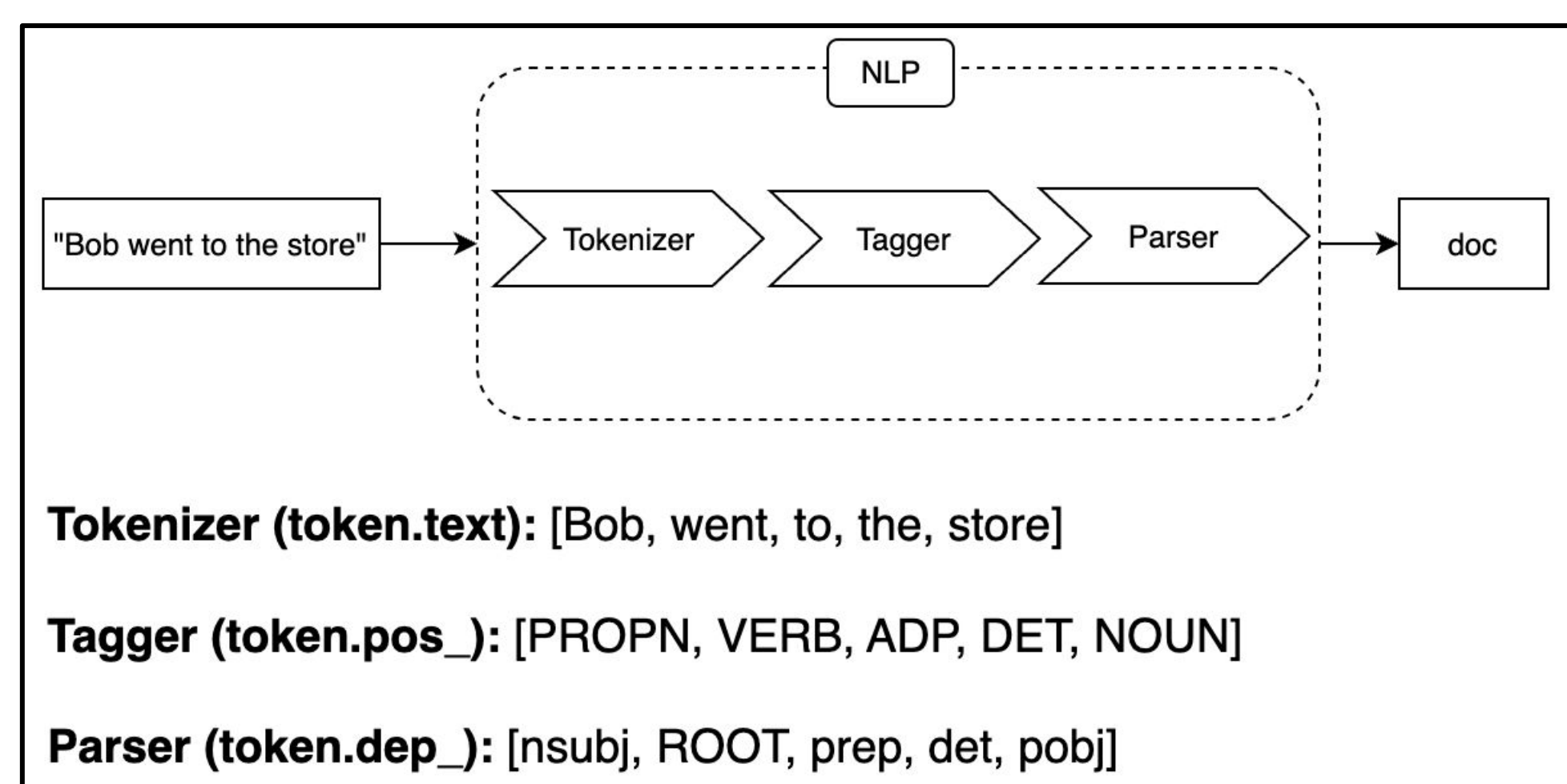  - Insight on the analysis algorithm



- About page provide important information regarding:
  - Error parameters
  - Algorithm details
  - Acknowledgment of anomalies
  - Foundational functionality

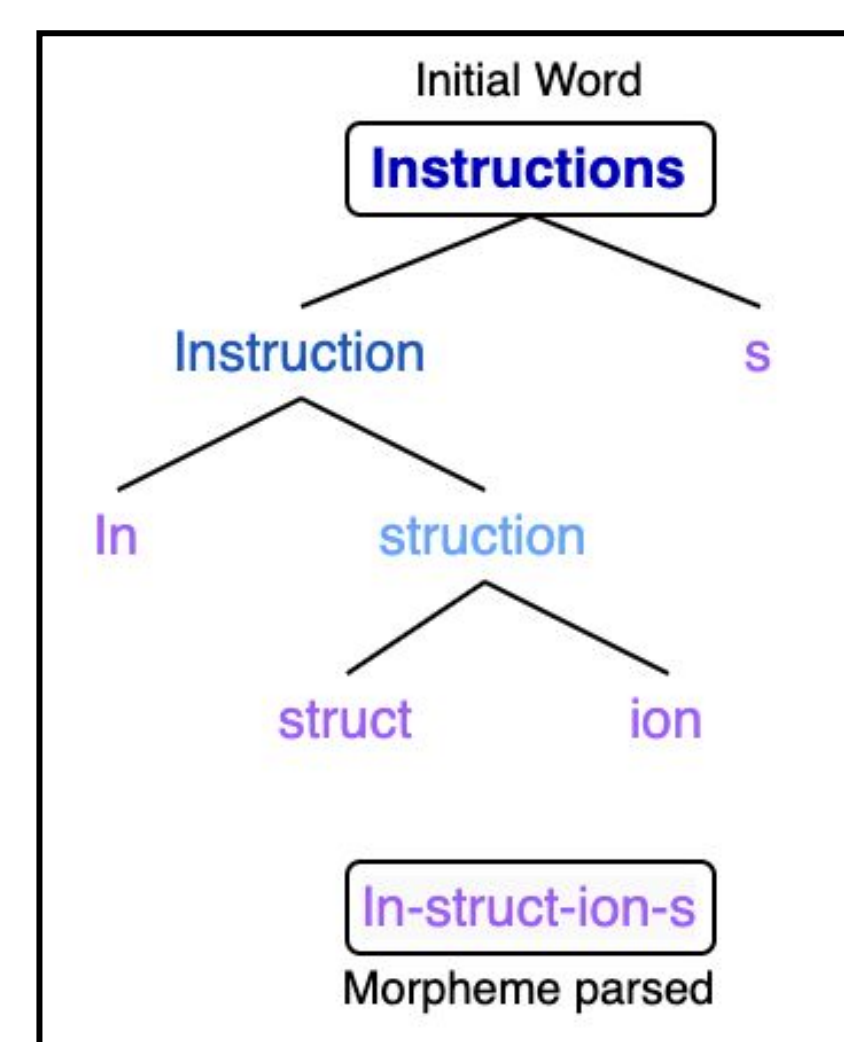## Overview of Tale Scanner Analysis (NLP)

### spaCy
- Open-source library for natural language processing (NLP) in Python
- Features include tokenization, parts of speech tagging and dependency parsing. Helped accomplish analysis for clauses and morphemes



**Tokenizer (token.text):** [Bob, went, to, the, store]

**Tagger (token.pos_):** [PROPN, VERB, ADP, DET, NOUN]

**Parser (token.dep_):** [nsubj, ROOT, prep, det, pobj]

### Morphemes
- Morphemes Python library that works in conjunction with the MorphoLex-en database

- Searches database and returns the information. Within this object is the word split up into it's morphology which includes its root, free, and bound aspects if they exist with in that word



Initial Word

Instructions

Instruction   s

In   struction

struct   ion

In-struct-ion-s

Morpheme parsed

### Gramformer
- Python library that highlights and detects grammar errors like spelling and punctuation via generative models
- Parses the user's input and then detects grammar errors, which has been learned by pre-trained and through fine tuning

## Case Testing and Analysis Gaps

- Case testing was foundational to analysis optimization

- Control cases were generated to isolate and identify exception cases

- Controlled test-runs and extensive system tests were recorded to keep track of consistency

- Warnings and parameters identified through case testing



- Values recorded per test-run:
  - Reload number
  - Analysis type
  - Expected word count, expected analysis output
  - Total run-time of test-run
  - Actual output
  - Respective label code for input text

## Future Enhancements

- Train large learning models for the remaining analysis types
  - Word choice errors
  - Story grammar elements
  - Cohesion

- Optimize each analysis type to handle all exception cases

- Implement multi-language analysis for necessary functions

- Handle different file input types (pdf, wordx, etc.)

## Summary

- Tale Scanner was designed with the intention of aiding speech language pathologists conduct statistical analysis more efficiently

- Ten of the original thirteen requirements were successfully implemented

- Tale Scanner provides a level of consistency in analysis that is often not represented by generalized artificial intelligence

- The system is designed and built with independent stability; changes to included libraries and technologies will not affect the build

- Case testing is foundational to optimizing natural language processing, and was fundamental to Tale Scanner's success