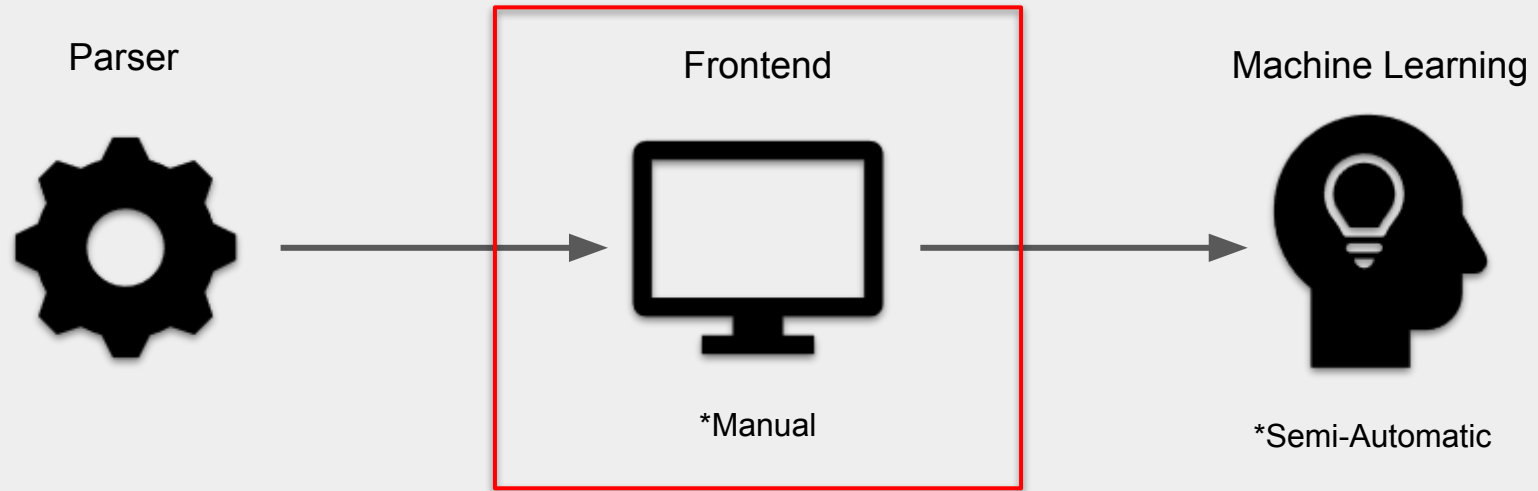# Natural Language Processing and Taxonomy Creation Tool

Brandon Christensen, Kadir Nour, Riley Hunter

The Boeing Company: Don Brancato & Rocky Bhatt
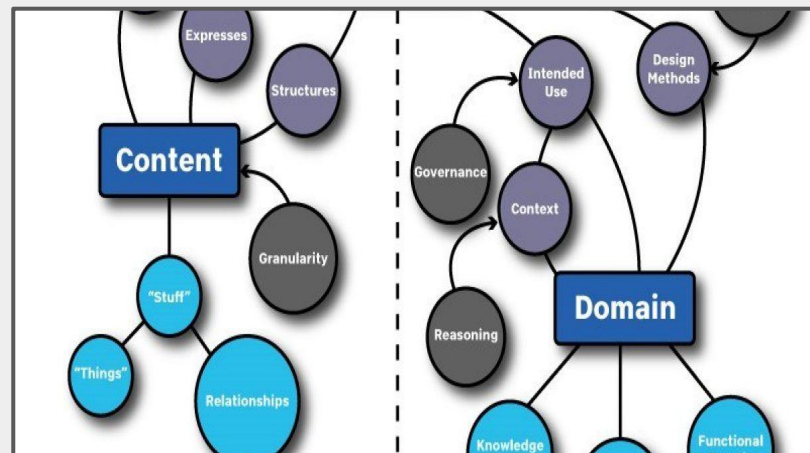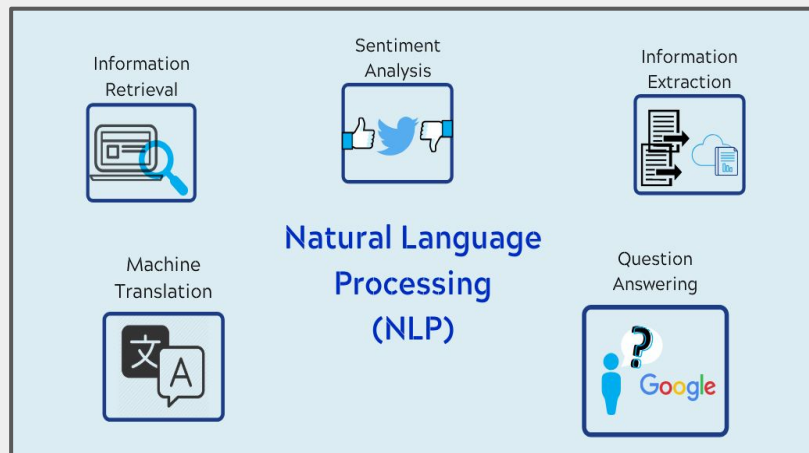
# Scope

Parser

Frontend

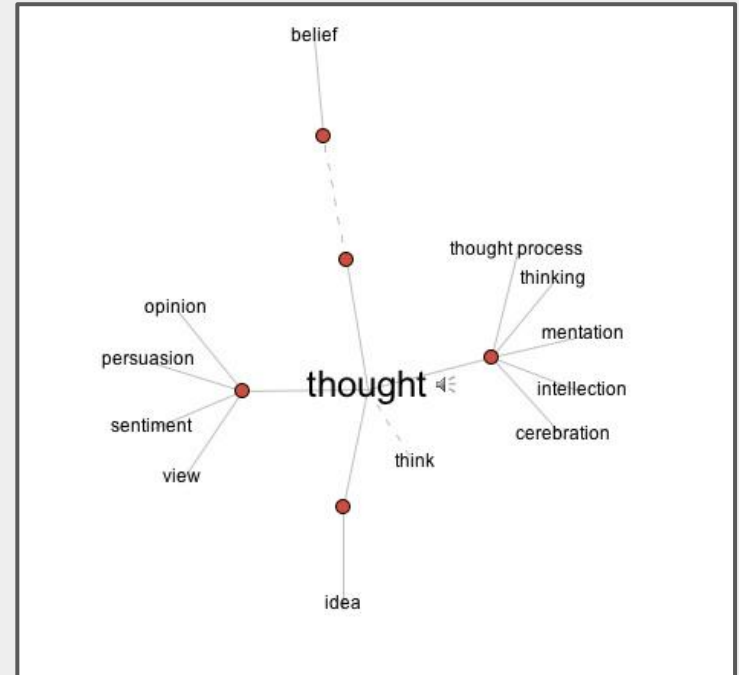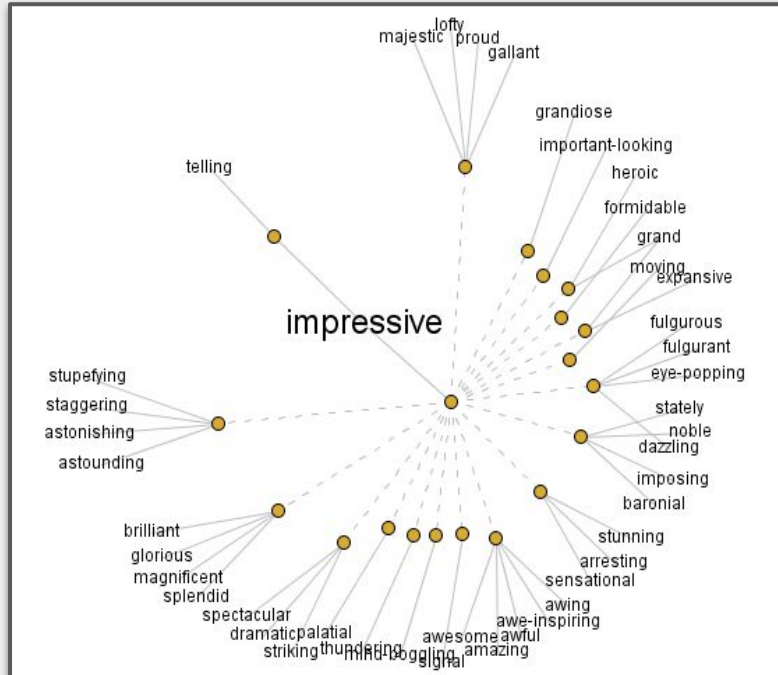*Manual

Machine Learning

*Semi-Automatic

# Motivation

**Goal:**

Create a state-of-the-art automatic taxonomy service that can parse and generate taxonomies for any constrained vertical or corpus based on semantic rule based system.
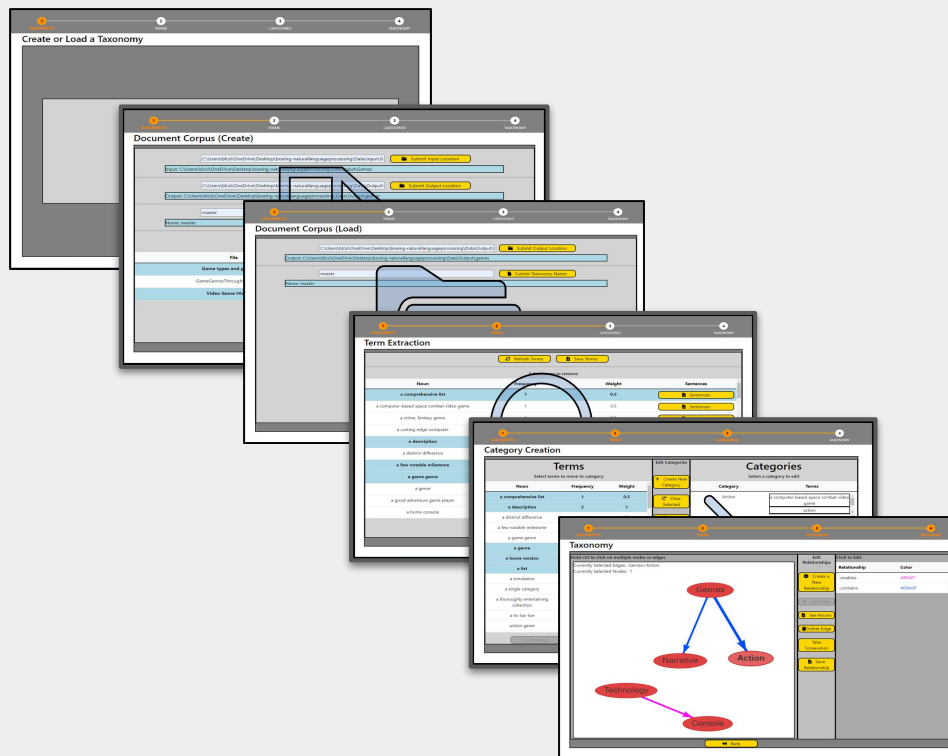
# Motivation

## Generalized Text

# Motivation

**Baseline Security**

```
Modified Security Configuration Name: Security Baseline - Windows 10 - November 2021 - Original
Modified Settings Count: 274
Original Security Configuration Name: Security Baseline - Windows 10 - December 2020 - IT Department
Original Settings Count: 273
****************************************************
TOTAL CHANGES():
ADDED IN Security Baseline - Windows 10 - November 2021 - Original :
windows10EndpointProtectionConfiguration defenderAllowScanScriptsLoadedInInternetExplorer
****************************************************
SETTING: windows10GeneralConfiguration passwordExpirationDays has differentiating values!!
****************************************************
SETTING: windows10GeneralConfiguration passwordBlockSimple has differentiating values!!
CONFIGURED IN Security Baseline - Windows 10 - November 2021 - Original : True
CONFIGURED IN Security Baseline - Windows 10 - December 2020 - IT Department : False
****************************************************
```
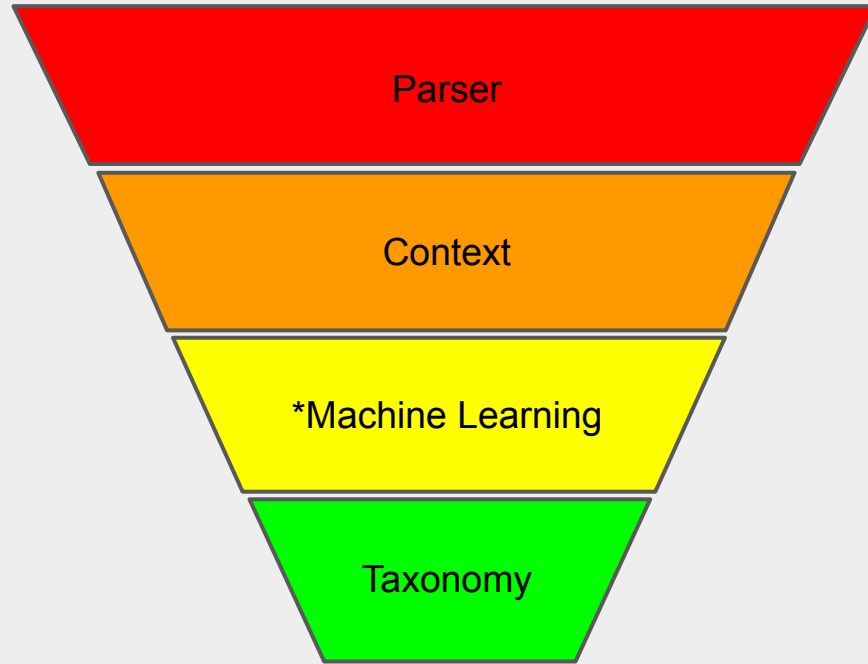
# Motivation

**Problem:**

1. Create a interactive and intuitive frontend.

2. Incorporate context.

3. Expert-driven manipulation of taxonomies.

4. Data saving.

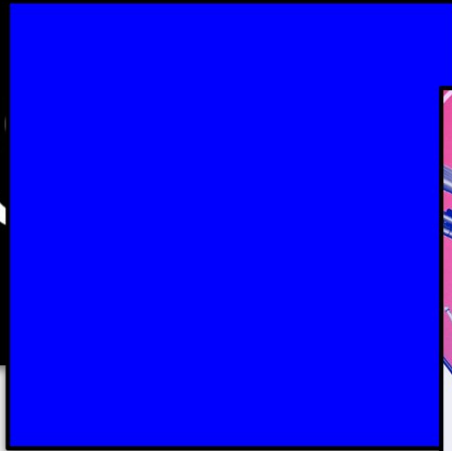# Solution

# Solution

## Why Do We Need Parsers?

"The quick brown fox jumps over the lazy dog" is a pangram—a sentence that contains all the letters of the alphabet. The phrase is commonly used for touch-typing practice, testing typewriters and computer keyboards, displaying examples of fonts, and other applications involving text where the use of all letters in the alphabet is desired.
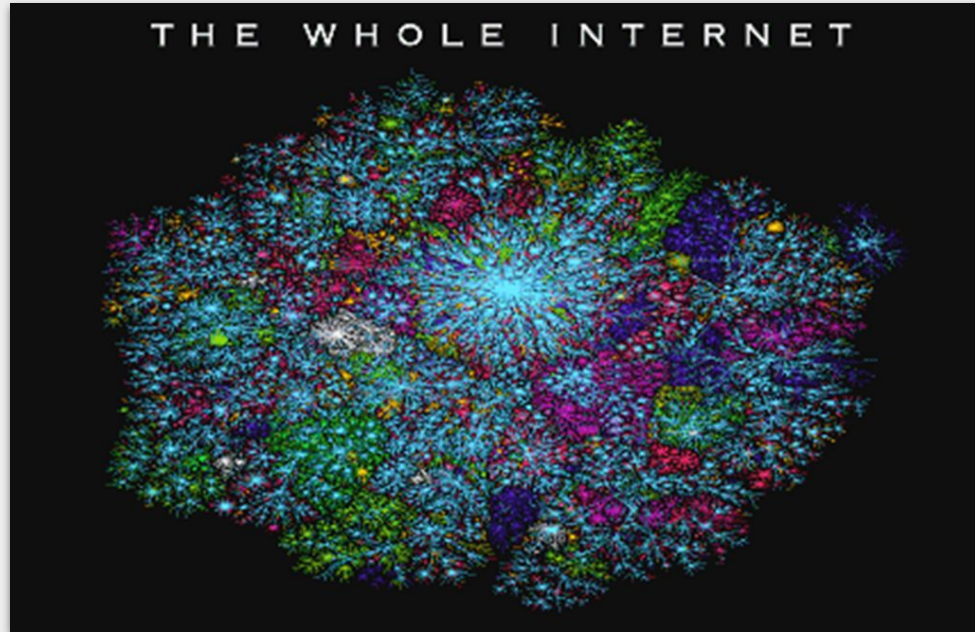
"The quick brown ==fox== jumps over the lazy ==dog==" is a ==pangram==—a ==sentence== that contains all the ==letters== of the ==alphabet==. The ==phrase== is commonly used for ==touch-typing== ==practice==, testing ==typewriters== and ==computer== ==keyboards==, displaying ==examples== of ==fonts==, and other ==applications== involving ==text== where the use of all ==letters== in the ==alphabet== is desired.

Solution
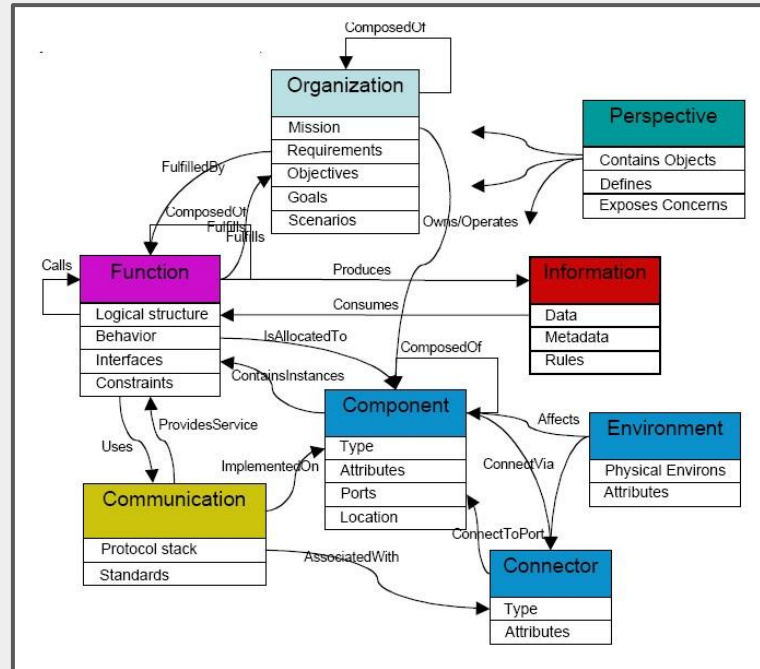
## **Why Do We Need Context?**

Solution

## Why Do We Need Machine Learning?

# Solution



**Why Do We Need A Taxonomy?**

# Solution



Parser

"The quick brown fox jumps over the lazy dog" is a pangram—a sentence that contains all the letters of the alphabet. The phrase is commonly used for touch-typing practice, testing typewriters and computer keyboards, displaying examples of fonts, and other applications involving text where the use of all letters in the alphabet is desired.
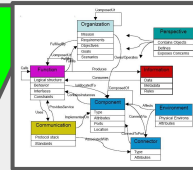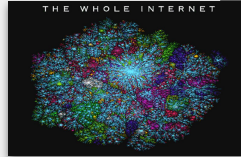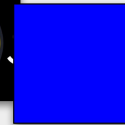
"The quick brown fox jumps over the lazy dog" is a pangram—a sentence that contains all the letters of the alphabet. The phrase is commonly used for touch-typing practice, testing typewriters and computer keyboards, displaying examples of fonts, and other applications involving text where the use of all letters in the alphabet is desired.

Context

THE WHOLE INTERNET

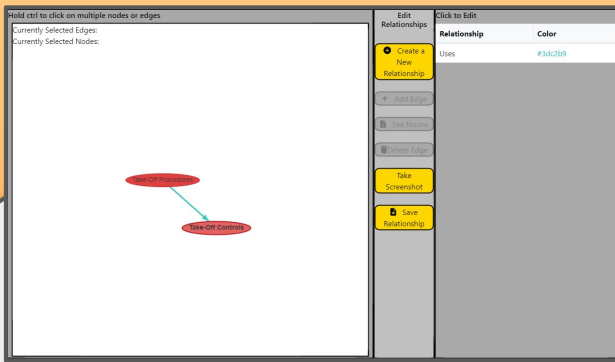*Machine Learning

Taxonomy

**Frontend**

# Use Case

## John at Boeing

John is a test pilot for Boeing. He would like to use the 737-flight manual as a reference, but he does not want to read through the entire 200-page document to find the information he needs.

1.  Running our tool first finds key terms from the document (the vocabulary).
2.  He then organizes the vocabulary into categories (as an expert himself, he knows which categories to create).
3.  One of the categories he creates is "Take-Off Procedures", and another is "Take-Off Controls". He finishes by creating a "uses" relationships between "Take-Off Procedures" and "Take-Off Controls".

Now, if John wants to find information on take-off procedures, he can search for that category. If he wants to find information on take-off controls, then he would check the related categories and find the "Take-Off Controls" category.
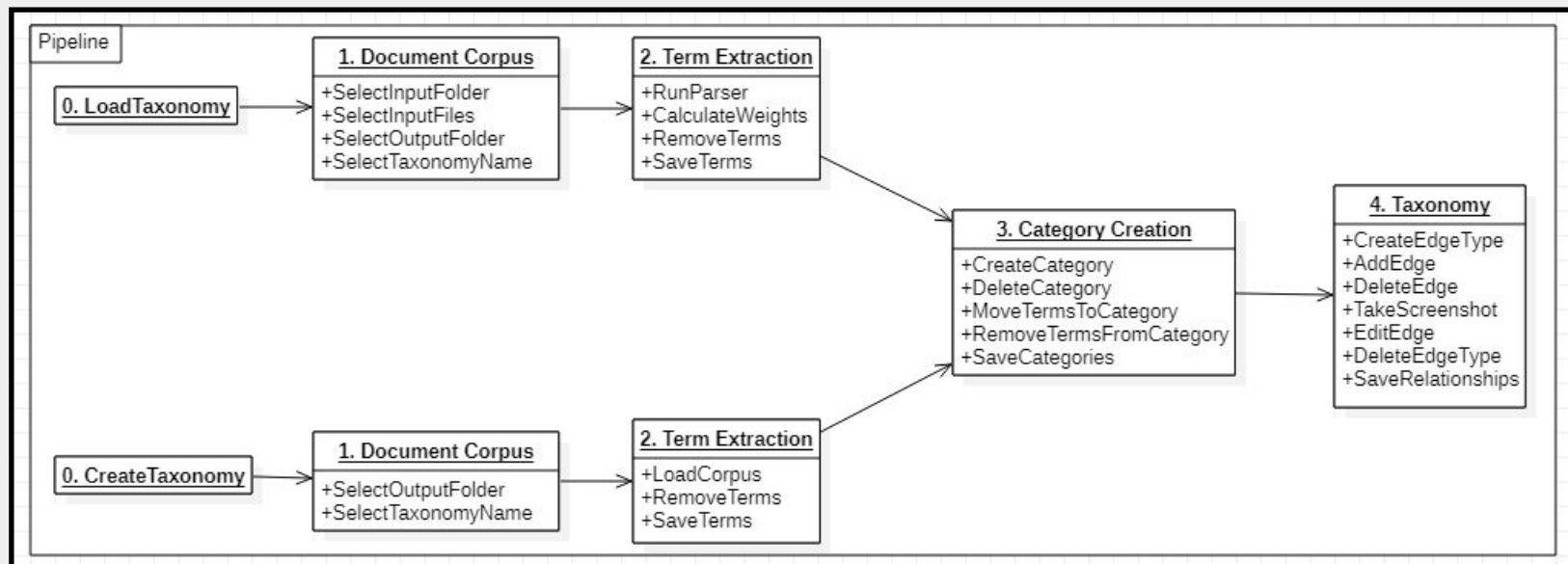
# John at Boeing

John is a test pilot for Boeing. He would like to use the 737-flight manual as a reference, but he does not want to read through the entire 200-page document to find the information he needs.

1. Running our tool first finds key terms from the document (the vocabulary).
2. He then organizes the vocabulary into categories (as an expert himself, he knows which categories to create).
3. One of the categories he creates is "Take-Off Procedures", and another is "Take-Off Controls". He finishes by creating a "uses" relationships between "Take-Off Procedures" and "Take-Off Controls".

Now, if John wants to find information on take-off procedures, he can search for that category. If he wants to find information on take-off controls, then he would check the related categories and find the "Take-Off Controls" category.

# Software Design



**Pipeline Architecture**

Pipeline

0. LoadTaxonomy

**1. Document Corpus**
+SelectInputFolder
+SelectInputFiles
+SelectOutputFolder
+SelectTaxonomyName

**2. Term Extraction**
+RunParser
+CalculateWeights
+RemoveTerms
+SaveTerms

**3. Category Creation**
+CreateCategory
+DeleteCategory
+MoveTermsToCategory
+RemoveTermsFromCategory
+SaveCategories

**4. Taxonomy**
+CreateEdgeType
+AddEdge
+DeleteEdge
+TakeScreenshot
+EditEdge
+DeleteEdgeType
+SaveRelationships

0. CreateTaxonomy

**1. Document Corpus**
+SelectOutputFolder
+SelectTaxonomyName

**2. Term Extraction**
+LoadCorpus
+RemoveTerms
+SaveTerms

# Tools and Technology

## Dev/Communications:

| | |
|---|---|
| Discord/Webex | Communications |
| GitHub | Version Control |
| Git Bash | Shell Environment |
| VS Code | Main Editor |
| Docker* | Containerization |

## Middle/ Backend:

| | |
|---|---|
| Flask | Front-to-back |
| SpaCy | NLP library |
| Pdfplummer | PDF library |

## Frontend:

| | |
|---|---|
| React | Frontend Framework |
| Bootstrap | Frontend Framework |
| Font Awesome | Buttons |
| Download | Download Screenshots |
| Html2Canvas | Screenshots |
| React-vis | Graph visualization |
| React-step-progress-bar | Progress Bar |

## Languages:

Python   JavaScript   HTML   CSS

# Testing

1. Unit and integration tests with Pytest. (Text extraction accuracy, Performance, etc.)

2. End-to-End tests with TestCafé.(Adding or deleting categories, Displaying terms in tables, etc.)

3. Automated CI/CD pipeline using GitHub

# Next Steps

- Machine Learning
  - Automation

- Frontend Tweeks
  - Search functionality

- Parser Accuracy
  - Optical Search

# Thanks!



**Thank you Don, Rocky, and AJ!**