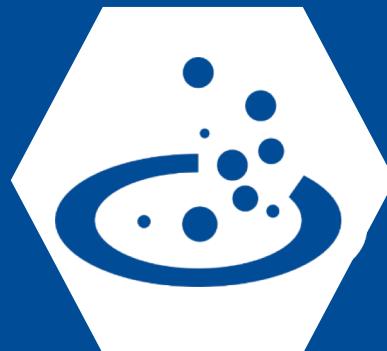




# Houston Data Analytics: Intro to Natural Language Processing (NLP)

Sense Corp Presentation

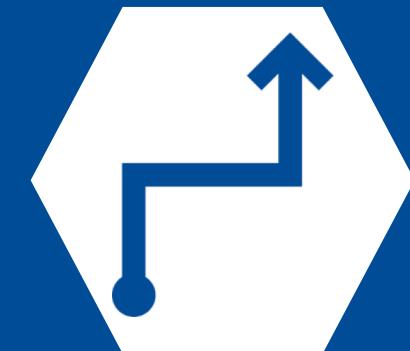
**We turn data into actionable insights and transform organizations for the digital era.**



**Data**



**Digital**



**Transformation**

# SENSE CORP

Powering Insight-Driven Organizations



Enterprise  
Data Strategy



Data  
Visualization



Data  
Governance



Advanced  
Analytics &  
Data Science



MDM



Big Data



Data  
Engineering

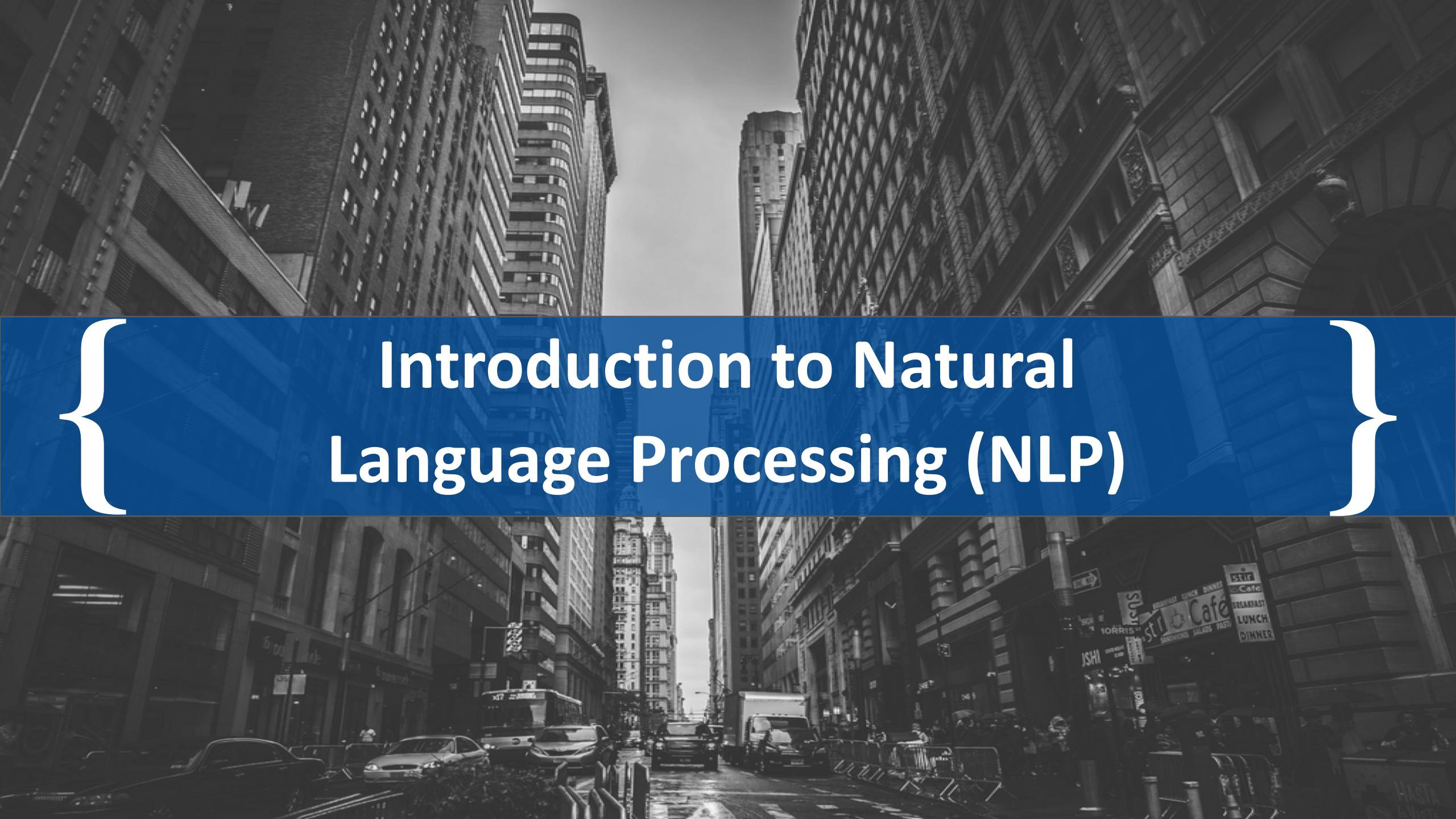


Data  
Security

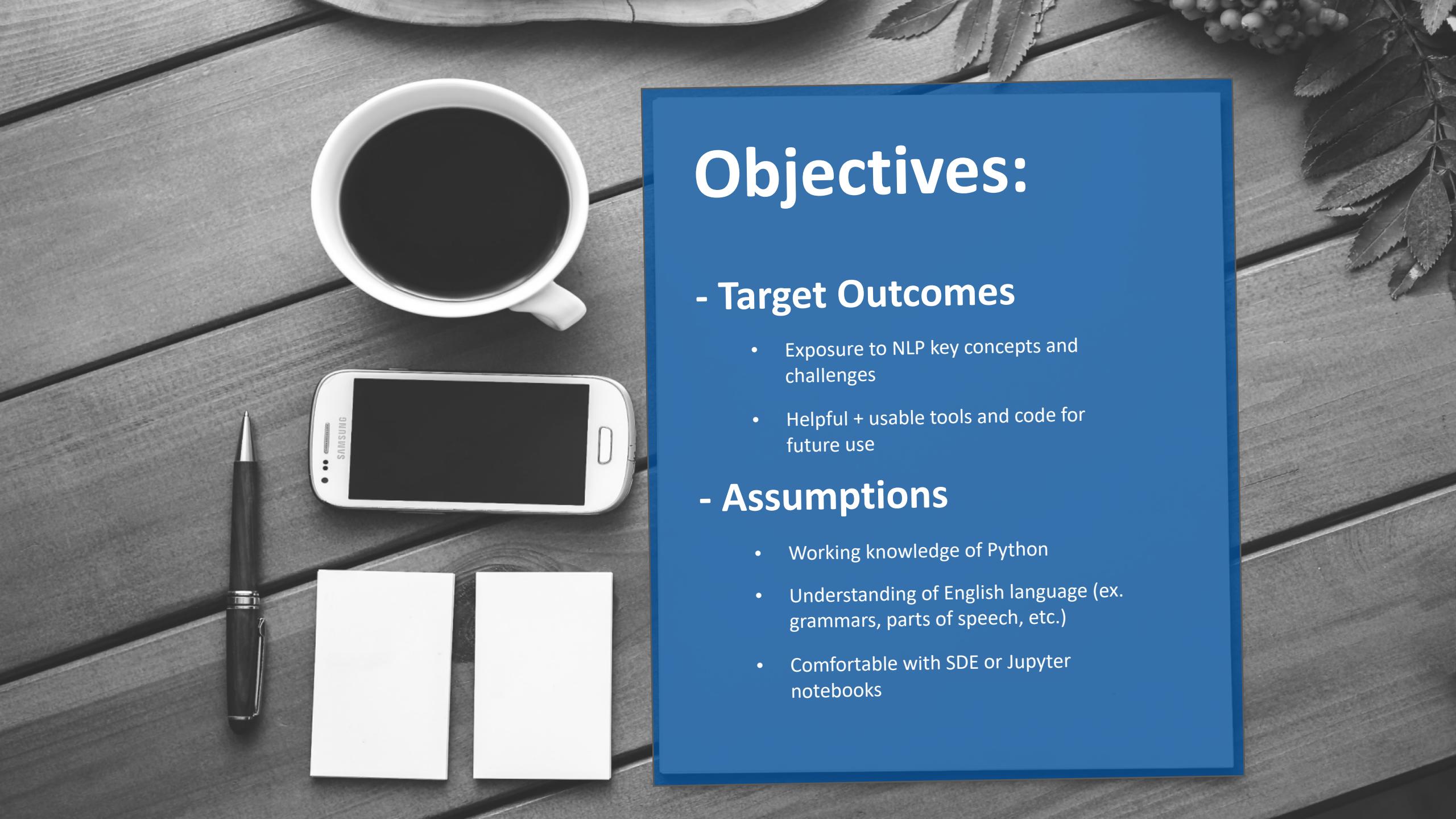
... we turn data into actionable  
insights and transform  
organizations for the digital era.

# Sense Corp has delivered exceptional results to hundreds of the world's largest global organizations and complex government agencies.

PUBLIC SECTOR	TELECOM/MEDIA	HEALTHCARE	FINANCIAL SERVICES	ENERGY	OTHER
 <b>TEXAS</b> Health and Human Services      	     	       	      	       	      



# Introduction to Natural Language Processing (NLP)



# Objectives:

## - Target Outcomes

- Exposure to NLP key concepts and challenges
- Helpful + usable tools and code for future use

## - Assumptions

- Working knowledge of Python
- Understanding of English language (ex. grammars, parts of speech, etc.)
- Comfortable with SDE or Jupyter notebooks

# What is NLP?

---

- **Wiki:** Natural language processing (NLP) is a field of computer science, artificial intelligence, and computational linguistics concerned with the interactions between computers and human (natural) languages.
- Identify the structure and meaning of words, sentences, texts and conversations
- Deep understanding of broad language
- NLP is all around us

# What are some applications?

- *Machine translation*
- *Dialog systems*
- *Sentiment analysis*
- *Text classification*
- *Question answering*
- *Digital assistants*
- *Information retrieval*

## Example: Language Comprehension

**Christopher Robin** is alive and well. **He** is the same person that you read about in the book, **Winnie the Pooh**. **As a boy**, **Chris** lived in a pretty home called **Cotchfield Farm**. When **Chris** was three years old, **his father** wrote a poem about **him**. The poem was printed in a magazine for others to read. **Mr. Robin** then wrote a book

Q: who wrote Winnie the Pooh?

Q: where is Chris lived?

# Why is NLP hard?

## News Headlines

- Ban on Nude Dancing on Governor's Desk
- Iraqi Head Seeks Arms
- Juvenile Court to Try Shooting Defendant
- Teacher Strikes Idle Kids
- Stolen Painting Found by Tree
- Local High School Dropouts Cut in Half
- Red Tape Holds Up New Bridges
- Clinton Wins on Budget, but More Lies Ahead
- Hospitals Are Sued by 7 Foot Doctors
- Kids Make Nutritious Snacks

**46 NEWS**  
WGCL TV ATLANTA

[cbs46.com](http://cbs46.com)  
Noon | 4:00 | 6:00 | 11:00

TOP STORIES  
BETTER MORNINGS  
WEATHER  
TRAFFIC  
46 INVESTIGATES  
RESTAURANT REPORT CARD  
CONSUMER  
SPORTS  
LOCAL HEADLINES  
ENTERTAINMENT  
CARTOON.TV  
CITY  
PROGRAM

[Email](#) [Print](#)

[A](#) [A](#) [A](#) Text Size

### Minister Accused Of Having 8 Wives In Jail

May 21, 2007 06:49 AM

**ATLANTA (AP)** -- A traveling minister who served two years in prison on bigamy charges has been jailed again for allegedly trying to marry more women.

Bishop Anthony Owens, 35, formerly of Duluth, Ga., is in a Gwinnett County jail after at least four women claimed he proposed to them after being released from prison in November 2005. Officials also say there is no evidence he divorced the eight wives he had married before going to prison.

A judge will decide whether he should go back to prison. Owens, who turned himself into the jail April 30, declined to be interviewed.

But his new fiancees aren't keeping quiet.

Betty Dixon, 38, met him last March in a casino near Memphis.

"Walker," the nurse said. "He told me God had sent him to help."

# Semantic Disambiguation

## Sentence

I ate spaghetti with meatballs.

I ate spaghetti with salad.

I ate spaghetti with abandon.

I ate spaghetti with a fork.

I ate spaghetti with a friend.

## Relation

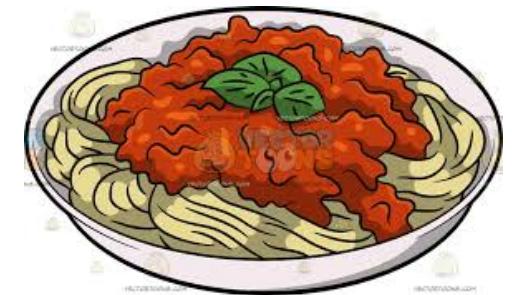
(ingredient of spaghetti)

(side dish of spaghetti)

(manner of eating)

(instrument of eating)

(accompanier of eating)



# NLP: Where do we start?

---

Like any other data analytics/science project: **Exploration!**

Code workshop: NLTK and Genesis (Bible)

- Lexical Dispersion Plot
- Basic statistics
  - Concordances, similarities, dispersion plots
  - Corpus Size (Vocabulary and lexical diversity)
  - Frequency Distributions
  - Collocations
  - Overall statistics (number of words, average word length, words per sentence)
  - Sentiment analysis

# (Larger) Example: The News

There are over 2M news articles published every day...

What is news?



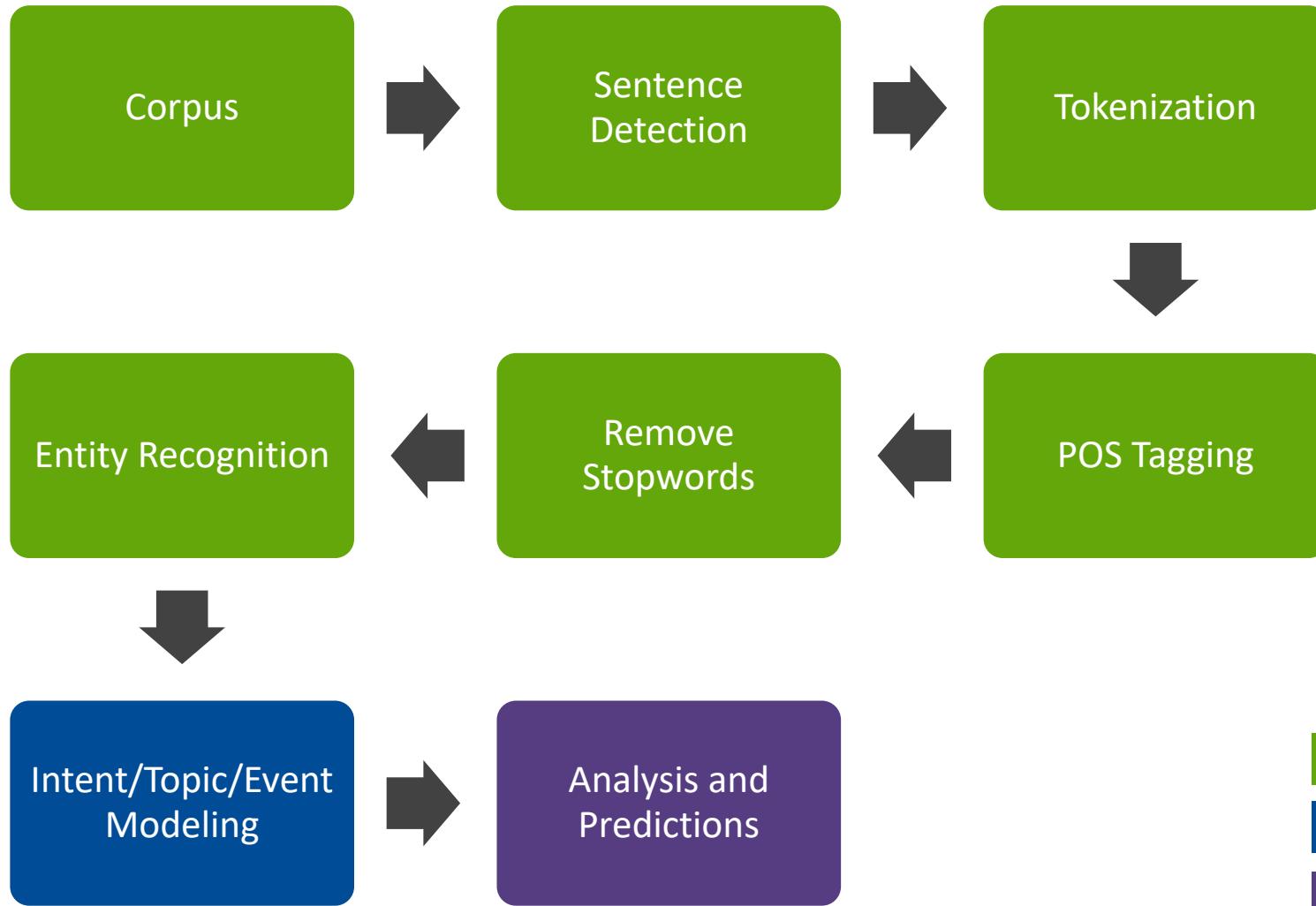
# News is whatever you want it to be

---



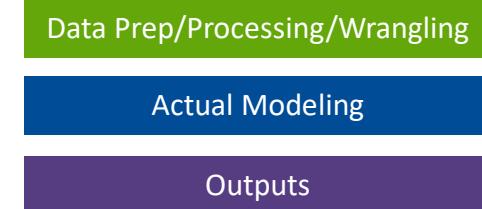
**Let's classify it (with NLP)**

# Walkthrough: NLP Pipeline

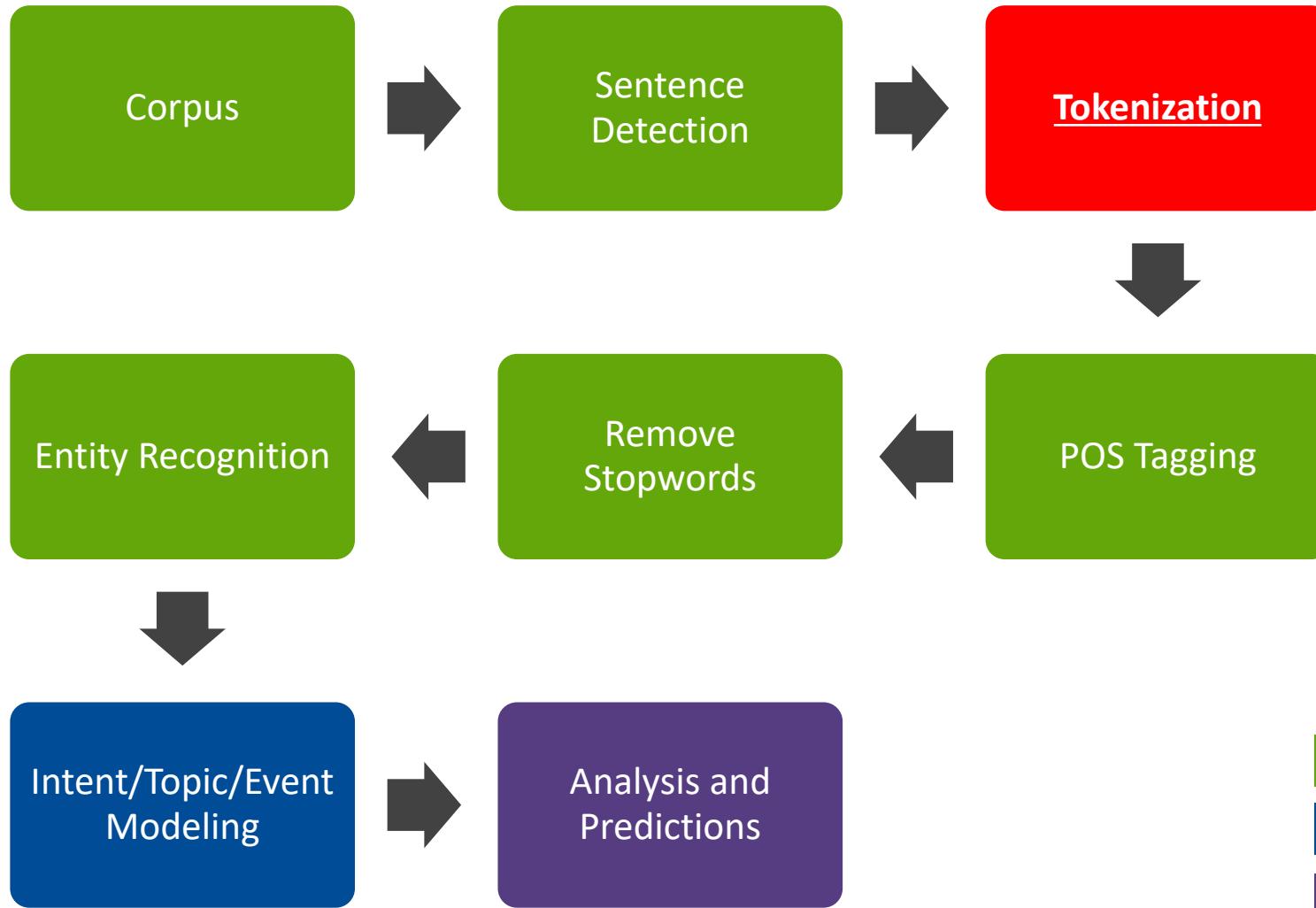


- General Framework (analogous to data processing pipeline)
- Data processing vs. modeling
- Extracting Signal vs. Removing Noise

## --LEGEND--



# Walkthrough: NLP Pipeline

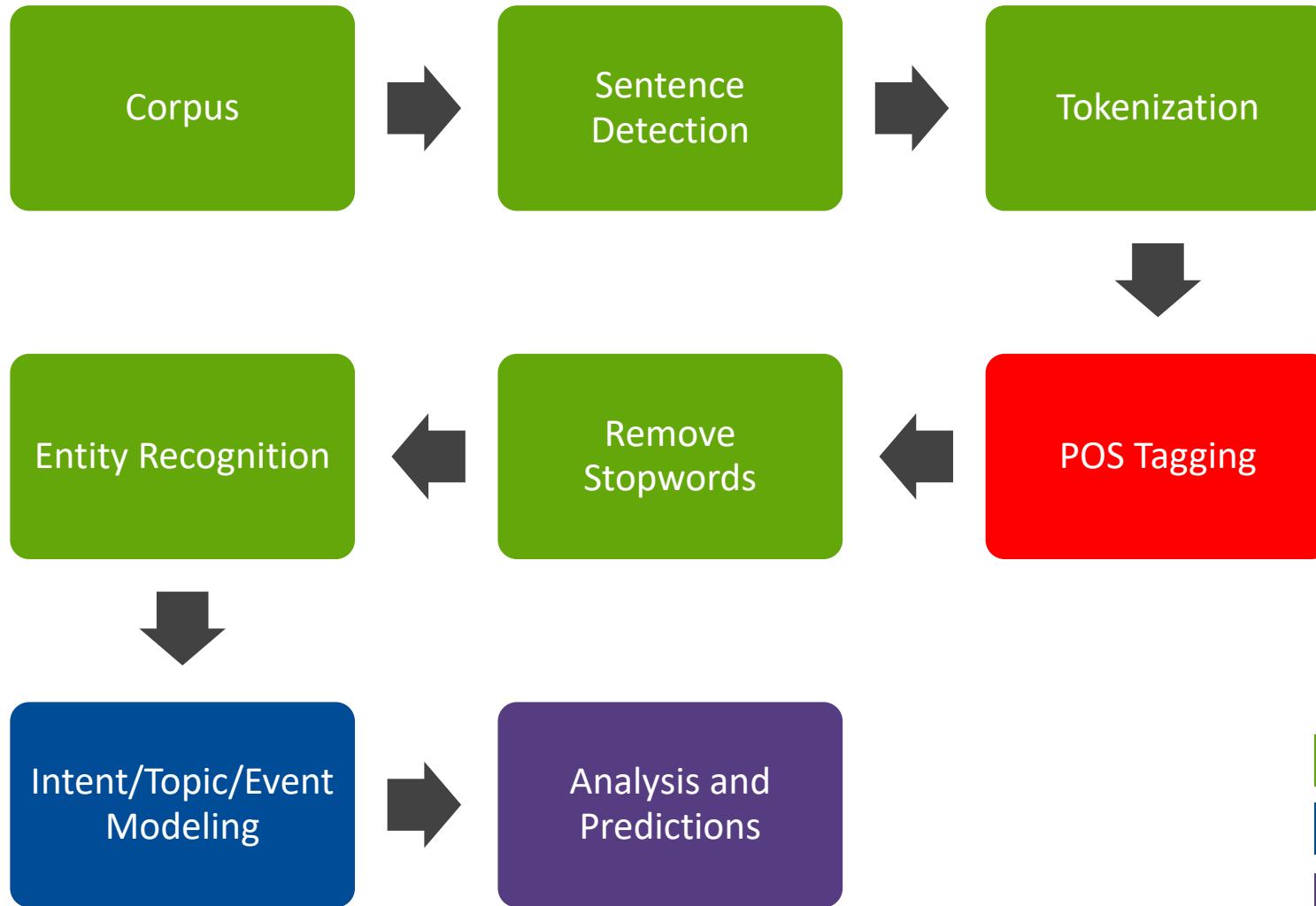


- Tokenization is 'chopping up' a text into pieces
- Tokens are also called 'grams' and can be designated to be any length

## --LEGEND--

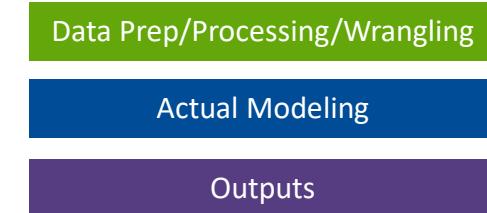
Data Prep/Processing/Wrangling
Actual Modeling
Outputs

# Walkthrough: NLP Pipeline

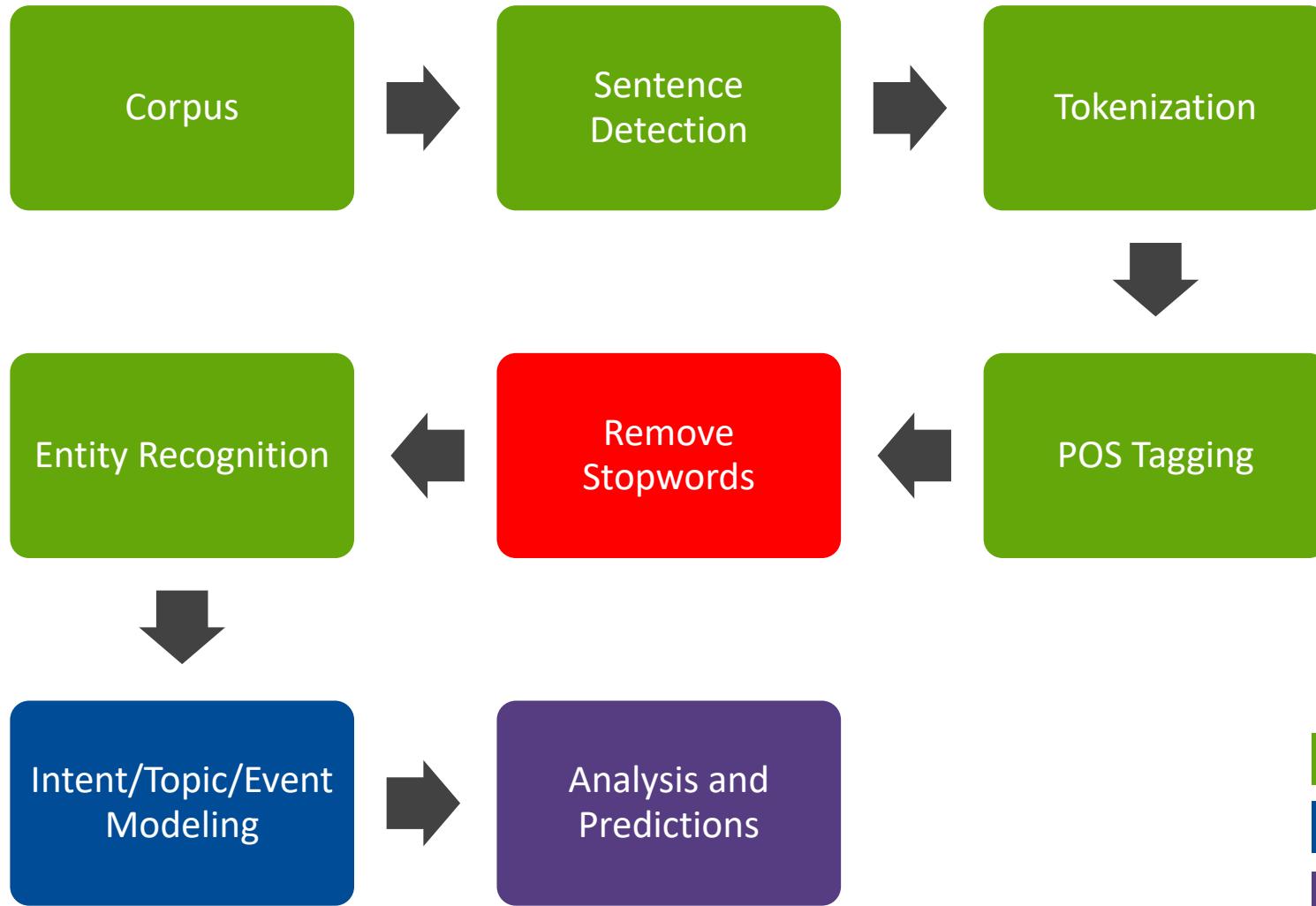


- POS Tagging, or “Part-of-speech tagging,” is identifying the part of speech of each token such as noun, adjective, etc.
- POS Tagging can be complex! The POS Tag is dependent on the definition of the word AND the usage.

## --LEGEND--



# Walkthrough: NLP Pipeline

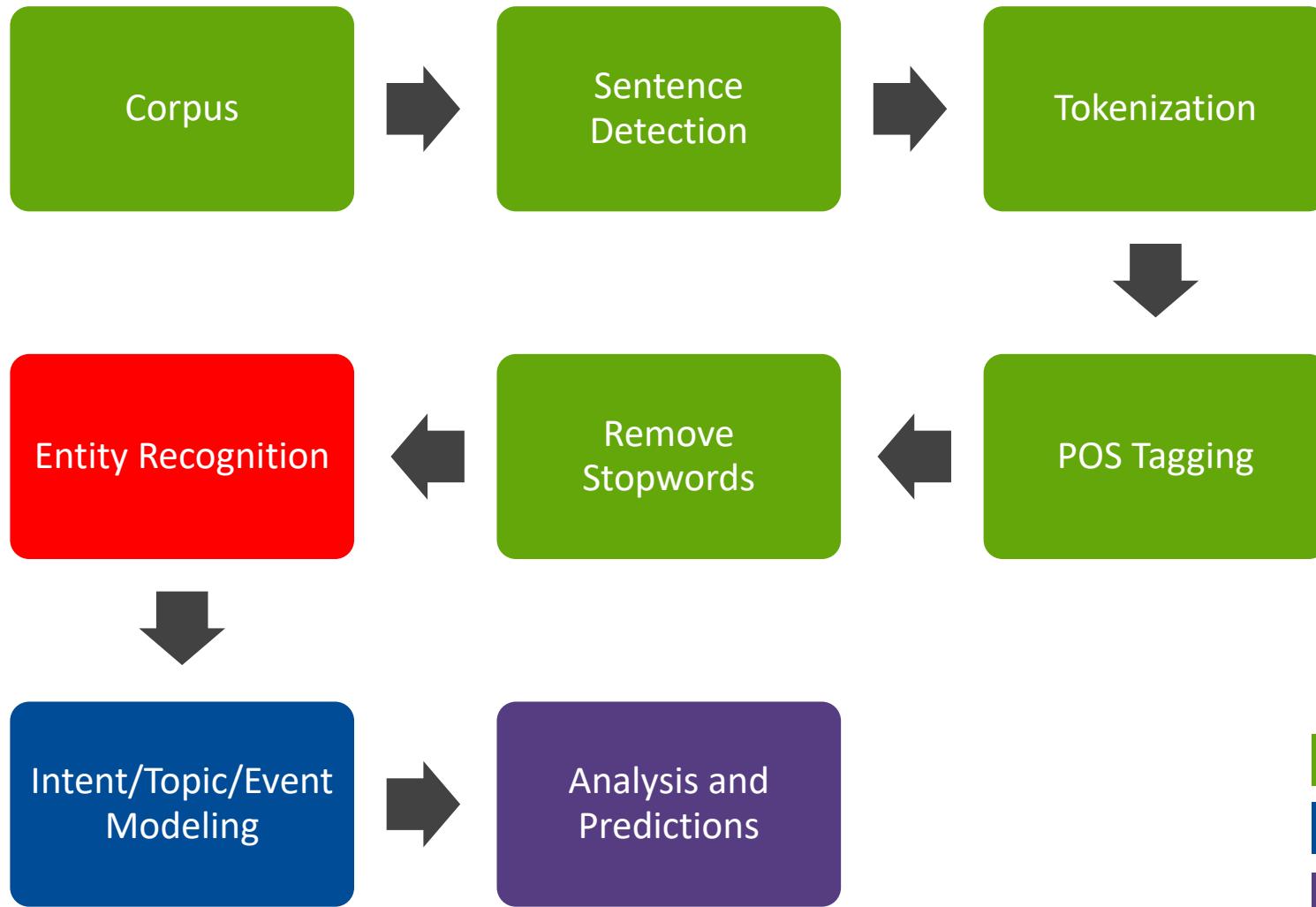


- Stopwords are extremely common words that have very little value, like 'the', 'is', etc.
- Because they are so common, they add significant 'noise' to our analysis

## --LEGEND--

Data Prep/Processing/Wrangling
Actual Modeling
Outputs

# Walkthrough: NLP Pipeline

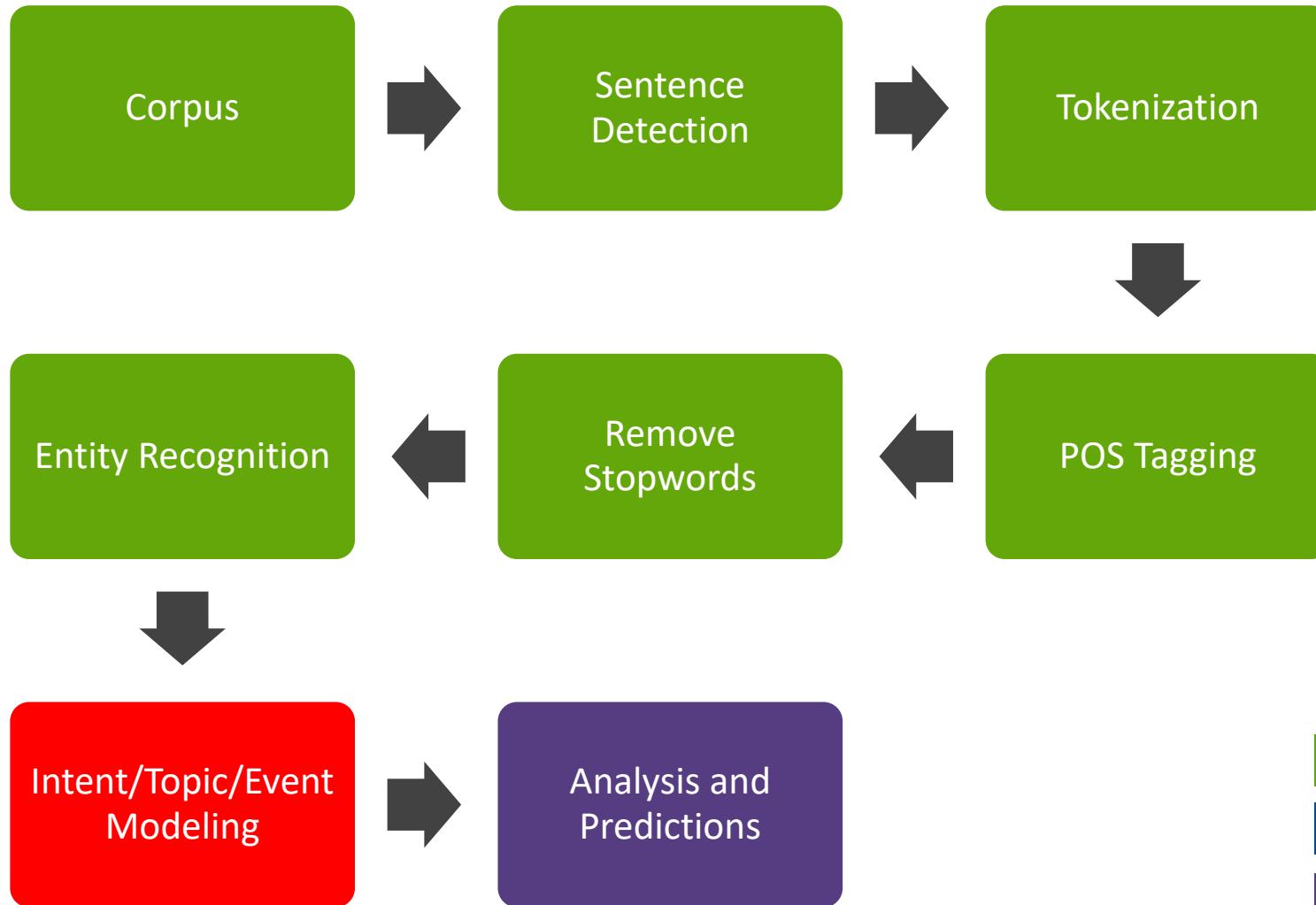


- Entity recognition is the task of identifying key nouns such as persons, places, etc.
- Entities can be identified by pre-defined dictionaries and/or grammars (POS tags).

## --LEGEND--

Data Prep/Processing/Wrangling
Actual Modeling
Outputs

# Walkthrough: NLP Pipeline



- For this example, we will use latent semantic analysis (LSA).
- LSA is a matrix decomposition technique that aims to map terms to a hidden topic layer. It is based on single value decomposition (SVD).

## --LEGEND--

Data Prep/Processing/Wrangling
Actual Modeling
Outputs

# Useful Tools

---

- Helpful Python packages
  - NLTK
    - Wordnet
  - GenSim
  - Textblob
  - SpaCy
  - Polyglot
  - Scikit-Learn
- 3<sup>rd</sup> Party Tools
  - Google Cloud NLP (Ex. Knowledge Graph)
  - Stanford Core NLP



# Takeaways



NLP have practical applications, but none do a great job in an open-ended domain



Sentences are understood through grammar, parsing and lexicons



Choosing a good interpretation of a sentence requires evidence from many sources



Most interesting NLP comes in connected discourse rather than in isolated sentences

THANKS!

**Justin J. Nguyen**

[jnguyen@sensecorp.com](mailto:jnguyen@sensecorp.com)

**Pratish Kanani**

[pkanani@sensecorp.com](mailto:pkanani@sensecorp.com)