

# Building End-to-End NLP Systems

## Lecture 6: NLP Pipelines

Harito

September 15, 2025

# Recap: Our NLP Components

So far, we've built several key components:

- **Tokenizer** (Lab 1): Breaks text into words.
- **Vectorizer** (Lab 2 3): Converts words into numerical vectors (Count, TF-IDF).
- **TextClassifier** (Lab 5): Learns to predict labels from these vectors.

Each component performs a specific task, but in a real-world application, they need to work together seamlessly.

# The Need for Pipelines

Imagine processing a new document for classification:

- ① Take the raw text.
- ② Pass it to the tokenizer.
- ③ Take the tokens and pass them to the vectorizer.
- ④ Take the vector and pass it to the classifier.
- ⑤ Get the final prediction.

This sequential process is exactly what an **NLP Pipeline** formalizes.

# What is an NLP Pipeline?

## Definition

An NLP pipeline is a sequence of processing steps applied to raw text to achieve a specific NLP task. The output of one component serves as the input for the next.

# Benefits of Using Pipelines

Pipelines are not just about convenience; they offer significant advantages:

- **Modularity:** Each step is a distinct, interchangeable component. You can easily swap out a 'SimpleTokenizer' for a 'RegexTokenizer'.
- **Reusability:** Once a pipeline is defined, it can be applied consistently to any new data.
- **Consistency:** Ensures that the exact same preprocessing and modeling steps are applied during both training and prediction.
- **Simplicity:** Reduces boilerplate code and makes complex workflows easier to manage and understand.
- **Deployment:** Simplifies deploying your NLP solution, as the entire process is encapsulated.

# Scikit-learn's Pipeline (Concept)

'scikit-learn' provides a powerful 'Pipeline' class that allows you to chain multiple estimators.

python from sklearn.pipeline import Pipeline from

sklearn.feature\_extraction.text import TfidfVectorizer from sklearn.linear\_model import

This is a conceptual example using sklearn's components

`text_clf = Pipeline([('tfidf', TfidfVectorizer()), ('clf', LogisticRegression()),])`

You can then fit and predict with the whole pipeline

`text_clf.fit(X_train, y_train)`  
`predicted = text_clf.predict(X_test)`

For our lab, we will build a custom pipeline using our own components.

Time for Lab 6!

## Objective:

- Implement a 'TextPipeline' class.
- Integrate your 'Tokenizer', 'Vectorizer', and 'TextClassifier'.
- Demonstrate end-to-end text processing.