# The NLP Pipeline for Legal Text Processing
## Short Notes & Visual Mind Maps

December 2, 2025

## Overview

The NLP Pipeline for legal text processing consists of 6 sequential stages that transform raw legal documents into structured, analyzable data. Each stage processes the output of the previous stage, building progressively refined representations of the text.

**Pipeline Stages:**

1. Sentence Segmentation

2. Tokenization

3. Stop-Words Removal

4. Lemmatization

5. Part-of-Speech (POS) Tagging

6. Named Entity Recognition (NER)

## 1. Sentence Segmentation

**Definition:** Breaking raw text into individual sentences as complete units.
**Why it Matters in Legal Context:**

- Legal documents contain complex sentence structures with multiple clauses

- Abbreviations (e.g., "U.S.A.", "Dr.", "Inc.") can break naive algorithms

- Citations (e.g., "42 U.S.C. $1983") often appear mid-sentence

**Challenges:**

- **Abbreviations:** "The plaintiff, John Smith, Jr., filed suit." (multiple periods)

- **Decimal Numbers:** "The fine was $5.25 per unit." (period != sentence end)

- **Ellipsis:** "He said... then left." (three periods)

- **URLs & Sections:** "See regulation at 18 U.S.C. $242 for details."

**Techniques:**

- **Rule-Based:** Regex patterns with exception lists for common abbreviations

- **Statistical:** PUNKT (punkt) sentence tokenizer in NLTK

- **Deep Learning:** Transformer-based models (e.g., spaCy)

**Example:**

```
Input: "The␣plaintiff␣filed␣suit␣on␣Jan.␣15,␣2023.␣She␣requested␣$50,000␣
   in␣damages."
Output:
  - Sentence 1: "The␣plaintiff␣filed␣suit␣on␣Jan.␣15,␣2023."
  - Sentence 2: "She␣requested␣$50,000␣in␣damages."
```

# 2. Tokenization

**Definition:** Splitting sentences into individual tokens (words, punctuation, numbers).
**Why it Matters in Legal Context:**

- Legal terminology includes hyphenated compounds ("mother-in-law", "cross-examination")

- Possessives require careful handling: "Smith's attorney" vs. "attorneys"

- Currency and numbers must be recognized: "$50,000", "$242", "18 U.S.C."

**Challenges:**

- **Contractions:** "don't" → {"do", "n't"} or {"don't"}?

- **Hyphenated Terms:** Keep "well-being" intact or split?

- **Special Characters:** Currency symbols, section markers ($)

- **Punctuation:** Keep attached or separate?

**Techniques:**

- **Whitespace:** Simple but loses punctuation

- **Regex:** Custom patterns for legal tokens

- **spaCy:** Advanced tokenizer with linguistic awareness

**Example:**

```
Input: "The␣defendant's␣attorney␣filed␣a␣motion."
Tokens: ["The", "defendant", "'s", "attorney",
        "filed", "a", "motion", "."]
```

# 3. Stop-Words Removal

**Definition:** Removing high-frequency words with minimal semantic value.
**Why it Matters in Legal Context:**

- Reduces noise and focuses on key legal terms

- Improves downstream processing efficiency

- Common legal stop-words: "the", "a", "and", "or", "is", "was"

  **CRITICAL WARNING for Legal Texts:**

- **NEVER remove negations:** "not", "no", "neither"

    – "guilty" vs. "not guilty" = complete opposite meaning!
    – "liable" vs. "not liable" = legal outcome reversal

- **NEVER remove modal verbs:** "shall", "may", "must", "can"

    – "The defendant shall pay" (obligation) vs. "The defendant may pay" (option)

**Techniques:**

- **Standard Stop-word Lists:** NLTK, spaCy (generic English)

- **Domain-Specific Lists:** Custom list excluding legal negations

- **TF-IDF Filtering:** Remove tokens with high document frequency

**Example:**

```
Input: "The␣plaintiff␣is␣not␣liable␣for␣damages."
Standard removal: "plaintiff␣liable␣damages"
Legal-safe removal: "plaintiff␣NOT␣liable␣damages"
```

# 4. Lemmatization

**Definition:** Reducing words to their base/dictionary form (lemma).
**Why it Matters in Legal Context:**

- Groups related forms: "filed", "files", "filing" → "file"

- Standardizes legal terminology variations

- Improves matching of similar legal concepts

**Lemmatization vs. Stemming:**

| Word | Lemmatization | Stemming |
|------|---------------|----------|
| filed | file (valid word) | fil |
| attorney | attorney | attorney |
| liabilities | liability | liabil |
| studying | study | studi |

**Key Difference:** Lemmatization produces valid dictionary words; stemming produces stems (may not be valid).
**Techniques:**

- **Dictionary-Based:** spaCy, NLTK WordNetLemmatizer (accuracy ∼95%)

- **Morphological Rules:** Custom rules for legal suffixes

- **Context-Aware:** Modern transformers (BERT, RoBERTa)

**Example:**

```
words = ["filed", "attorney", "liabilities", "worse"]
lemmas = ["file", "attorney", "liability", "bad"]
```

# 5. Part-of-Speech (POS) Tagging

**Definition:** Assigning grammatical tags (noun, verb, adjective, etc.) to each token.
**Why it Matters in Legal Context:**

- **Intent Detection:** Modal verbs ("shall", "may", "must") indicate legal obligations

- **Key Information Extraction:** Nouns = entities, Verbs = actions

- **Sentiment Analysis:** Adverbs indicate intent ("deliberately", "knowingly")

**Common Legal POS Tags:**

| Tag | Example | Legal Significance |
|-----|---------|--------------------|
| NN | defendant, court | Entities |
| VB | filed, testified | Actions |
| MD | shall, may, must | Obligations |
| RB | deliberately, knowingly | Intent/Adverbial |
| JJ | liable, criminal | State/Property |

**Techniques:**

- **Rule-Based:** Manual POS tagging rules

- **Statistical:** Hidden Markov Models (HMM)

- **Deep Learning:** BiLSTM, Transformers (BERT, RoBERTa)

**Example:**

```
Sentence: "The␣defendant␣shall␣pay␣damages."
POS Tags:
  - The: DET (determiner)
  - defendant: NN (noun)
  - shall: MD (modal verb) <- Legal obligation!
  - pay: VB (verb)
  - damages: NN (noun)
```

# 6. Named Entity Recognition (NER)

**Definition:** Identifying and classifying named entities (persons, organizations, locations, etc.).
**Why it Matters in Legal Context:**

- Extracts key actors: parties (plaintiff, defendant, judge), law firms

- Identifies laws and regulations: "42 U.S.C. $1983", "Roe v. Wade"

- Captures critical details: dates, amounts, locations, verdicts

**Legal-Specific Entity Types:**

| Entity Type | Example | Format |
|---|---|---|
| PERSON | John Smith, Judge Davis | Names |
| ORG | Smith & Associates LLC | Entities |
| GPE | New York, USA | Geographic |
| LAW | 42 U.S.C. $1983 | Statute/Code |
| DATE | January 15, 2023 | Temporal |
| MONEY | $50,000 | Monetary |
| VERDICT | guilty, acquitted | Outcomes |

**Techniques:**

- **Rule-Based:** Regex patterns (works well for legal citations)

- **Statistical:** CRF (Conditional Random Fields)

- **Deep Learning:** BiLSTM-CRF, Transformers (BERT-based models)

**Example:**

```
Text: "Smith␣v.␣Jones,␣42␣U.S.C.␣§1983␣(2020)"
Entities:
  - Smith: PERSON
  - Jones: PERSON
  - 42 U.S.C. §1983: LAW
  - 2020: DATE
```

# Complete Pipeline Execution

**Input:** *"Smith v. Jones: The plaintiff was not liable for damages under 42 U.S.C. $1983."*
   **Stage 1 - Sentence Segmentation:**

```
Sentence 1: "Smith v. Jones: The plaintiff was not liable..."
Sentence 2: "...for damages under 42 U.S.C. $1983."
```

   **Stage 2 - Tokenization:**

```
["Smith", "v", ".", "Jones", ":", "The", "plaintiff",
 "was", "not", "liable", "for", "damages", "under",
 "42 U.S.C. \$1983", "."]
```

   **Stage 3 - Stop-Words Removal (Domain-Safe):**

```
["Smith", "v", "Jones", "plaintiff", "NOT", "liable",
 "damages", "42 U.S.C. \$1983"]
(Preserved: NOT = critical negation)
```

**Stage 4 - Lemmatization:**

```
["Smith", "v", "Jones", "plaintiff", "not", "liable",
 "damage", "42 U.S.C. \$1983"]
```

**Stage 5 - POS Tagging:**

```
Smith/NN v/NN Jones/NN plaintiff/NN not/RB liable/JJ
damage/NN 42 U.S.C. \$1983/NNP
```

**Stage 6 - Named Entity Recognition:**

```
PERSON: Smith, Jones
LAW: 42 U.S.C. \$1983
VERDICT: not liable
```

# Key Takeaways for Legal NLP

1. **Negations are Sacred:** Never remove "not", "no", "neither", "cannot" from legal texts.

2. **Modal Verbs Matter:** "Shall" (obligation), "may" (permission), "must" (requirement) have distinct meanings.

3. **Citations are Atomic:** Preserve legal citations intact (42 U.S.C. $1983, Roe v. Wade).

4. **Context is King:** Legal language is highly context-dependent.

5. **Domain-Specific Tools:** Build custom stop-word lists for legal text.

6. **Pipeline Order Matters:** Errors compound downstream. Quality matters at each stage.

7. **Modern Transformers Win:** BERT-based models outperform traditional pipelines.

# Quick Reference: Legal Abbreviations

| Abbr. | Full Form | Example |
|---|---|---|
| U.S.C. | United States Code | 42 U.S.C. $1983 |
| C.F.R. | Code of Federal Regulations | 29 C.F.R. $1910 |
| v. | versus | Smith v. Jones |
| J. | Judge/Justice | Judge Smith |
| Esq. | Esquire | John Smith, Esq. |
| et al. | and others | Smith et al. v. State |

# Implementation Tips

- Use spaCy for production NLP (fast, accurate)

- Use NLTK for educational/research purposes

- Combine regex with statistical methods for citations

- Build custom domain-specific stop-word lists

- Test on actual legal domain data

- Preserve POS tags throughout pipeline

- Use ensemble methods for NER