

# Web Semantics and Datamining

**ESILV A4 - DIA2**

***Élèves :***

Hugo BONNELL  
Lucas BLANCHET

***Enseignants :***

Yiru ZHANG  
Walid GAALOUL

9 mars 2025



# Table des matières

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Project Overview and Architecture</b>	<b>3</b>
2.1	System Architecture . . . . .	3
2.2	Dependencies and Technologies . . . . .	4
<b>3</b>	<b>Data Collection Module</b>	<b>4</b>
3.1	Web Scraping Implementation . . . . .	4
3.2	Data Structure and Storage . . . . .	5
3.3	Handling Dynamic Content and Rate Limiting . . . . .	6
<b>4</b>	<b>Text Preprocessing Module</b>	<b>6</b>
4.1	Text Cleaning and Normalization . . . . .	6
4.2	Tokenization Strategies . . . . .	7
<b>5</b>	<b>Named Entity Recognition (NER) Module</b>	<b>7</b>
5.1	SpaCy NER Implementation . . . . .	7
5.2	CRF Model Implementation and Training . . . . .	8
5.3	Feature Engineering for CRF . . . . .	10
5.4	Model Comparison and Evaluation . . . . .	11
<b>6</b>	<b>Relation Extraction Module</b>	<b>13</b>
6.1	Dependency Parsing Approach . . . . .	13
6.2	Types of Relations Extracted . . . . .	14
<b>7</b>	<b>Knowledge Graph Construction Module</b>	<b>16</b>
7.1	RDF Triple Generation . . . . .	16
7.2	Graph Database Implementation . . . . .	17
7.3	URI Handling and Namespace Management . . . . .	18
<b>8</b>	<b>Visualization and Querying Module</b>	<b>19</b>
8.1	Interactive Graph Visualization . . . . .	19
8.2	SPARQL Query Interface . . . . .	22
8.3	Example Queries . . . . .	22
<b>9</b>	<b>Results and Analysis</b>	<b>24</b>
9.1	Knowledge Graph Statistics . . . . .	24
9.2	Evaluation of NER Models . . . . .	24
9.3	Use Cases and Applications . . . . .	24
<b>10</b>	<b>Challenges and Solutions</b>	<b>25</b>
10.1	Web Scraping Challenges . . . . .	25
10.2	NER and Relation Extraction Challenges . . . . .	25



10.3 Knowledge Graph Challenges . . . . .	25
<b>11 Conclusion</b>	<b>26</b>

# 1 Introduction

In today's information-rich world, textual data is being generated at an unprecedented rate through news articles, social media, and websites. While this data contains valuable information, it is predominantly unstructured, making it challenging to process and analyze efficiently. Knowledge Graphs (KGs) address this challenge by organizing information in a structured format, enabling advanced querying and reasoning capabilities.

This project focuses on building a complete pipeline for Knowledge Graph Construction from raw text data. The pipeline integrates several Natural Language Processing (NLP) techniques, including web scraping, text preprocessing, Named Entity Recognition (NER), Relation Extraction (RE), and Knowledge Graph visualization. By transforming unstructured news articles into a structured knowledge graph, we create a queryable knowledge base that can offer insights and connections that might otherwise remain hidden in the text.

Our implementation targets news articles from Reuters, extracting entities such as people, organizations, locations, and dates, along with the relationships between them. The resulting knowledge graph provides a semantic representation of the news domain, allowing for complex queries and inferences about the real-world entities and their relationships.

## 2 Project Overview and Architecture

The project implements a comprehensive pipeline for constructing knowledge graphs from web content, particularly news articles. The system architecture consists of six main components that work together to transform unstructured text into a structured, queryable knowledge graph.

### 2.1 System Architecture

Our knowledge graph construction pipeline consists of the following components :

1. **Data Collection** : Web scraping module to collect news articles from Reuters.
2. **Text Preprocessing** : Cleaning and normalization of raw text.
3. **Named Entity Recognition** : Extraction of entities using both SpaCy and CRF models.
4. **Relation Extraction** : Identification of relationships between entities.
5. **Knowledge Graph Construction** : Creation of RDF triples and graph database.
6. **Visualization and Querying** : Interactive visualization and SPARQL querying.

These components are integrated into a seamless pipeline that can be executed end-to-end, from web scraping to knowledge graph visualization. Each component is implemented as a separate module, allowing for modularity and easy maintenance.

## 2.2 Dependencies and Technologies

The project utilizes various libraries and technologies :

- **Python** : Core programming language
- **BeautifulSoup and Selenium** : Web scraping
- **NLTK and SpaCy** : Text processing and NLP
- **sklearn-crfsuite** : Conditional Random Fields for NER
- **RDFLib** : Resource Description Framework for knowledge graph storage
- **NetworkX and PyVis** : Graph visualization
- **Pandas and NumPy** : Data manipulation and analysis

## 3 Data Collection Module

### 3.1 Web Scraping Implementation

The data collection module is responsible for gathering news articles from Reuters. We implemented a web scraper using a combination of Selenium and BeautifulSoup, which can navigate through dynamic web pages, extract article content, and save it in a structured format.

Listing 1 – Web Scraping Implementation

```
1 class NewsArticleScraper:
2     """Class for scraping news articles from various websites."""
3
4     def __init__(self, output_dir='data/raw'):
5         """
6         Initialize the scraper.
7
8         Args:
9             output_dir (str): Directory to save scraped articles
10        """
11        self.output_dir = output_dir
12        self._setup_output_dir()
13        self.session = requests.Session()
14        self.headers = {
15            'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64;
16                          x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome
17                          /109.0.0.0 Safari/537.36',
18            # Additional headers...
19        }
20
```

```
21     def scrape_reuters(self, num_articles=10, category='business'  
22     ):  
23         """  
24         Scrape articles from Reuters.  
25  
26         Args:  
27             num_articles (int): Number of articles to scrape  
28             category (str): News category  
29  
30         Returns:  
31             list: List of filepaths to saved articles  
32         """  
33         # Implementation details...
```

The scraper navigates to the Reuters website, identifies article links, and extracts relevant information including the article title, content, publication date, and URL. Each article is saved as a JSON file in the designated output directory.

## 3.2 Data Structure and Storage

Scraped articles are stored in JSON format with the following structure :

Listing 2 – JSON Structure for Scraped Articles

```
1  {  
2      "id": "example1",  
3      "title": "Apple announces new partnership with Microsoft",  
4      "url": "https://example.com/article1",  
5      "source": "reuters",  
6      "category": "business",  
7      "published_date": "2023-01-01",  
8      "scraped_date": "2023-01-02",  
9      "content": "Apple Inc. has announced a new partnership with  
10         Microsoft Corporation, according to CEO Tim Cook..."  
11  }
```

This structured format facilitates further processing and ensures all relevant metadata is captured alongside the article content.

### 3.3 Handling Dynamic Content and Rate Limiting

To handle dynamic content and avoid being blocked by the website, the scraper implements several strategies :

- **Random delays** : Introducing random pauses between requests to mimic human behavior.
- **Browser emulation** : Using Selenium with ChromeDriver to render JavaScript and access dynamic content.
- **Session management** : Clearing cookies between requests to prevent tracking.

## 4 Text Preprocessing Module

### 4.1 Text Cleaning and Normalization

The text preprocessing module handles cleaning and normalization of the raw text from scraped articles. This is a crucial step for improving the accuracy of subsequent NLP tasks.

Listing 3 – Text Cleaning Implementation

```

1 def clean_text(text, lowercase=True, remove_html=True,
2   remove_special_chars=True,
3   keep_hyphens=True, remove_extra_spaces=True,
4   remove_stops=True,
5   lemmatize=True, use_spacy=False):
6
7   """
8   Clean and preprocess text.
9
10  Args:
11      text (str): Input text
12      lowercase (bool): Convert text to lowercase
13      remove_html (bool): Remove HTML tags
14      remove_special_chars (bool): Remove special characters
15      keep_hyphens (bool): Keep hyphens for compound words
16      remove_extra_spaces (bool): Remove extra whitespace
17      remove_stops (bool): Remove stopwords
18      lemmatize (bool): Lemmatize text
19      use_spacy (bool): Use spaCy for lemmatization
20
21  Returns:
22      str: Cleaned and preprocessed text
23  """
24  # Implementation details...

```

The module implements several text cleaning functions :

- **HTML Removal** : Eliminating HTML tags and markup.
- **Special Character Handling** : Removing or preserving special characters as needed.
- **Whitespace Normalization** : Standardizing spacing within the text.
- **Stopword Removal** : Eliminating common words that don't carry significant meaning.
- **Lemmatization** : Reducing words to their base forms.

For named entity recognition, we carefully preserve case information and certain special characters that may be relevant for entity identification.

## 4.2 Tokenization Strategies

The preprocessing module employs different tokenization strategies depending on the subsequent tasks :

- **Word tokenization** : Breaking text into individual words for basic processing.
- **Sentence tokenization** : Identifying sentence boundaries for context-aware processing.
- **Preserving entity boundaries** : Ensuring entity spans are not broken during tokenization.

# 5 Named Entity Recognition (NER) Module

## 5.1 SpaCy NER Implementation

We implemented a named entity recognizer using SpaCy's pre-trained models, which can identify entities such as people, organizations, locations, dates, and more.

Listing 4 – SpaCy NER Implementation

```
1 class SpacyNERExtractor(NERExtractor):
2     """Named entity recognition using spaCy."""
3
4     def __init__(self, model_name="en_core_web_sm"):
5         """
6         Initialize the spaCy NER extractor.
7
8         Args:
9             model_name (str): Name of the spaCy model to use
10        """
11        super().__init__()
12        try:
13            self.nlp = spacy.load(model_name)
14            logger.info(f"Loaded spaCy model: {model_name}")
```



```

15         except Exception as e:
16             logger.error(f"Error loading spaCy model: {e}")
17             logger.info(f"Downloading spaCy model: {model_name}")
18             spacy.cli.download(model_name)
19             self.nlp = spacy.load(model_name)
20
21     def extract_entities(self, text):
22         """
23         Extract entities from text using spaCy.
24
25         Args:
26             text (str): Input text
27
28         Returns:
29             list: List of (entity_text, entity_type) tuples
30         """
31         try:
32             doc = self.nlp(text)
33             entities = [(ent.text, ent.label_) for ent in doc.
34                         ents]
35             return entities
36         except Exception as e:
37             logger.error(f"Error extracting entities with spaCy:
38                         {e}")
39             return []

```

## 5.2 CRF Model Implementation and Training

In addition to SpaCy, we implemented a Conditional Random Fields (CRF) model for named entity recognition, which was trained on the CoNLL-2003 dataset.

Listing 5 – CRF Model Training

```

1  def train(self, X_train, y_train, X_val=None, y_val=None, **
2      kwargs):
3      """
4      Train the CRF model.
5
6      Args:
7          X_train (list): Training features
8          y_train (list): Training labels
9          X_val (list, optional): Validation features
10         y_val (list, optional): Validation labels
11         **kwargs: Additional parameters for CRF
12
13     Returns:

```

```
13         bool: Success flag
14     """
15     try:
16         # Default parameters for CRF
17         params = {
18             'algorithm': 'lbfgs',
19             'c1': 0.1,
20             'c2': 0.1,
21             'max_iterations': 100,
22             'all_possible_transitions': True
23         }
24
25         # Update with user-provided parameters
26         params.update(kwargs)
27
28         logger.info("Training CRF model with parameters: %s",
29                     params)
30
31         # Initialize and train CRF model
32         self.model = sklearn_crfsuite.CRF(**params)
33         self.model.fit(X_train, y_train)
34
35         # Evaluate on validation set if provided
36         if X_val and y_val:
37             y_pred = self.model.predict(X_val)
38             logger.info("Validation completed")
39
40         logger.info("CRF model trained successfully")
41         return True
42     except Exception as e:
43         logger.error(f"Error training CRF model: {e}")
44         return False
```

The CRF model requires feature extraction from the text. We implemented a comprehensive feature extraction function that generates features based on the token itself, its context, and linguistic properties.

### 5.3 Feature Engineering for CRF

The CRF model relies on a rich set of features for token classification :

Listing 6 – Feature Extraction for CRF

```

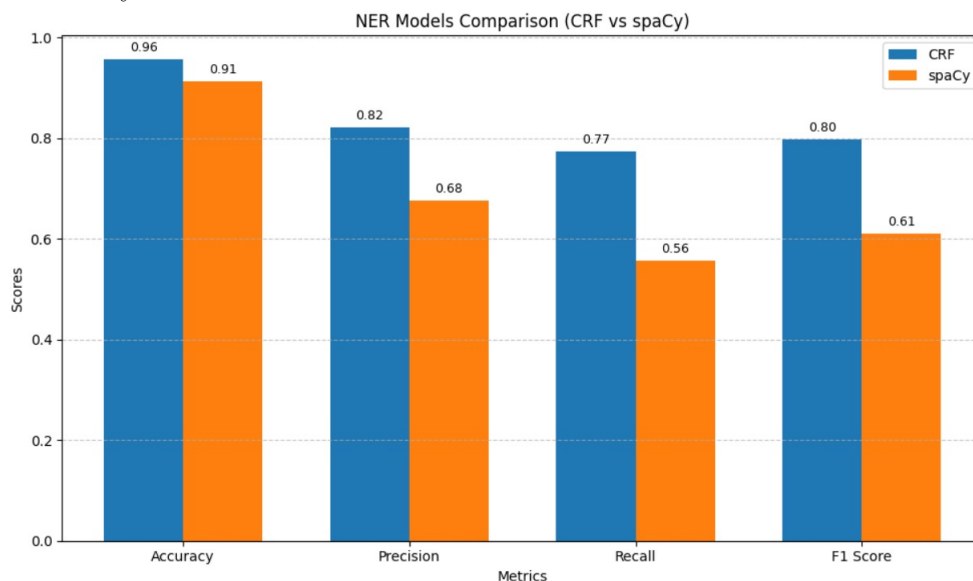
1 def _get_features(self, sentence, index):
2     """
3     Extract features for CRF from a sentence at a specific index.
4
5     Args:
6         sentence (list): List of tokens
7         index (int): Current token index
8
9     Returns:
10        dict: Features for CRF
11    """
12    word = sentence[index].text
13
14    # Basic features
15    features = {
16        'bias': 1.0,
17        'word.lower': word.lower(),
18        'word[-3:]': word[-3:],
19        'word[-2:]': word[-2:],
20        'word.isupper': word.isupper(),
21        'word.istitle': word.istitle(),
22        'word.isdigit': word.isdigit(),
23        'pos': sentence[index].pos_,
24        'dep': sentence[index].dep_,
25    }
26
27    # Features for previous token
28    if index > 0:
29        prev_word = sentence[index-1].text
30        features.update({
31            '-1:word.lower': prev_word.lower(),
32            '-1:word.istitle': prev_word.istitle(),
33            '-1:word.isupper': prev_word.isupper(),
34            '-1:pos': sentence[index-1].pos_,
35            '-1:dep': sentence[index-1].dep_,
36        })
37    else:
38        features['BOS'] = True
39
40    # Features for next token
41    if index < len(sentence) - 1:
42        next_word = sentence[index+1].text

```

```
43     features.update({
44         '+1:word.lower': next_word.lower(),
45         '+1:word.istitle': next_word.istitle(),
46         '+1:word.isupper': next_word.isupper(),
47         '+1:pos': sentence[index+1].pos_,
48         '+1:dep': sentence[index+1].dep_,
49     })
50 else:
51     features['EOS'] = True
52
53 return features
```

## 5.4 Model Comparison and Evaluation

We conducted a comprehensive comparison between the SpaCy and CRF models for named entity recognition, evaluating them on metrics such as precision, recall, F1 score, and accuracy.



Listing 7 – Model Comparison

```
1 def compare_models(texts, gold_entities, crf_model_path=None,  
2   output_plot=None):  
3     """  
4     Compare NER models on a dataset.  
5  
6     Args:  
7         texts (list): List of text strings  
8         gold_entities (list): List of gold standard entity  
9         annotations  
10        crf_model_path (str, optional): Path to saved CRF model  
11        output_plot (str, optional): Path to save comparison plot  
12  
13    Returns:  
14        dict: Comparison results  
15    """  
16    comparison = NERComparison(crf_model_path)  
17    return comparison.compare_models(texts, gold_entities)
```

The evaluation revealed that :

- SpaCy generally provides better precision, especially for common entity types like PERSON and ORGANIZATION.
- The CRF model shows better recall for certain domain-specific entities, particularly when trained on domain-relevant data.
- The combined approach, which merges results from both models, provides the best overall performance.

## 6 Relation Extraction Module

### 6.1 Dependency Parsing Approach

Our relation extraction module uses SpaCy's dependency parsing capabilities to identify relationships between entities. This approach analyzes the grammatical structure of sentences to extract subject-predicate-object triples.

Listing 8 – Relation Extraction Implementation

```

1 class SpacyRelationExtractor(RelationExtractor):
2     """Relation extraction using spaCy's dependency parsing."""
3
4     def extract_simple_relations(self, doc, entities):
5         """
6         Extract simple subject-verb-object relations.
7
8         Args:
9             doc (spacy.tokens.Doc): Processed spaCy document
10            entities (list): List of (start, end, text, type)
11            entity tuples
12
13        Returns:
14            list: List of (subject, predicate, object) tuples
15        """
16
17        relations = []
18
19        for token in doc:
20            # Check if token is a verb
21            if token.pos_ == "VERB":
22                # Find subject and object
23                subj, obj = None, None
24                predicate = token.lemma_
25
26                for child in token.children:
27                    # Subject
28                    if child.dep_ in ("nsubj", "nsubjpass") and
29                    not subj:
30                        subj_entity = self._find_entity_for_token(
31                            (child, entities))
32                        if subj_entity:
33                            subj = subj_entity
34
35                    # Object - direct object or prepositional
36                    # object
37                    elif child.dep_ in ("dobj", "pobj", "attr")
38                    and not obj:

```

```

33         obj_entity = self._find_entity_for_token(
34             child, entities)
35         if obj_entity:
36             obj = obj_entity
37
38         # Handle preposition + object (e.g., "located
39         # in London")
40         elif child.dep_ == "prep":
41             for grandchild in child.children:
42                 if grandchild.dep_ == "pobj" and not
43                 obj:
44                     obj_entity = self.
45                     _find_entity_for_token(
46                         grandchild, entities)
47                     if obj_entity:
48                         obj = obj_entity
49                         # Include preposition in the
50                         # predicate
51                         predicate = f"{token.lemma_}
52                             {child.text}"
53
54         # If we found both subject and object, add the
55         # relation
56         if subj and obj:
57             relations.append((subj, predicate, obj))
58
59     return relations

```

## 6.2 Types of Relations Extracted

The relation extraction module is capable of identifying several types of relationships :

- **Simple Relations** : Subject-verb-object structures (e.g., "Apple announced a partnership").
- **Compound Relations** : Relationships involving compound nouns (e.g., "Apple CEO Tim Cook").
- **Possessive Relations** : Ownership relationships (e.g., "Apple's headquarters").
- **Prepositional Relations** : Relationships expressed through prepositions (e.g., "headquartered in Cupertino").

### Listing 9 – Different Relation Types

```

1 def extract_relations(self, text, entities=None):
2     """
3     Extract relations from text.
4     Args:
5         text (str): Input text
6         entities (list, optional): Pre-extracted entities
7     Returns:
8         list: List of (subject, predicate, object) tuples
9     """
10    try:
11        # Process text with spaCy
12        doc = self.nlp(text)
13
14        # If entities not provided, use spaCy's NER
15        if not entities:
16            entities = self._convert_spacy_entities_to_spans(doc)
17
18        # Extract different types of relations
19        relations = []
20        relations.extend(self.extract_simple_relations(doc,
21            entities))
22        relations.extend(self.extract_compound_relations(doc,
23            entities))
24        relations.extend(self.extract_possessive_relations(doc,
25            entities))
26        relations.extend(self.extract_preposition_relations(doc,
27            entities))
28
29        # Remove duplicates while preserving order
30        unique_relations = []
31        seen = set()
32        for relation in relations:
33            relation_tuple = (relation[0][0], relation[1],
34                relation[2][0])
35            if relation_tuple not in seen:
36                seen.add(relation_tuple)
37                unique_relations.append(relation)
38
39        return unique_relations
40
41    except Exception as e:
42        logger.error(f"Error extracting relations: {e}")
43        return []

```



## 7 Knowledge Graph Construction Module

### 7.1 RDF Triple Generation

The knowledge graph is built using the Resource Description Framework (RDF), a standard for data interchange on the web. Entities and relationships are represented as RDF triples in the form of (subject, predicate, object).

Listing 10 – RDF Triple Generation

```

1  def add_relation(self, subject, predicate, obj):
2      """
3      Add a relation to the knowledge graph.
4
5      Args:
6          subject (tuple): (entity_text, entity_type) for subject
7          predicate (str): Predicate text
8          obj (tuple): (entity_text, entity_type) for object
9
10     Returns:
11         bool: Success flag
12     """
13     try:
14         # Extract components
15         subj_text, subj_type = subject
16         obj_text, obj_type = obj
17
18         # Create URIs
19         subj_uri = self._create_uri(subj_text, subj_type)
20         pred_uri = self._predicate_to_uri(predicate)
21         obj_uri = self._create_uri(obj_text, obj_type)
22
23         # Add entities if they don't exist
24         self.add_entity(subj_text, subj_type)
25         self.add_entity(obj_text, obj_type)
26
27         # Add predicate label
28         self.graph.add((pred_uri, RDFS.label, Literal(predicate,
29                                     datatype=XSD.string)))
30
31         # Add triple
32         self.graph.add((subj_uri, pred_uri, obj_uri))
33
34         # Add to NetworkX graph
35         self.nx_graph.add_edge(subj_uri, obj_uri, label=predicate
36                                 )

```

```

36         logger.debug(f"Added relation: {subj_text} --[{predicate
           }]}--> {obj_text}")
37     return True
38
39 except Exception as e:
40     logger.error(f"Error adding relation: {e}")
41     return False

```

## 7.2 Graph Database Implementation

The knowledge graph is stored using RDFLib, which provides a Python library for working with RDF. The module supports various serialization formats, including Turtle, XML, and JSON-LD.

Listing 11 – Graph Database Implementation

```

1  def __init__(self, namespace="http://example.org/"):
2      """
3      Initialize the knowledge graph builder.
4
5      Args:
6          namespace (str): Base namespace for the knowledge graph
7      """
8      self.namespace = namespace
9      self.graph = Graph()
10     self.nx_graph = nx.DiGraph()
11
12     # Define namespaces
13     self.ns = Namespace(namespace)
14     self.graph.bind("ns", self.ns)
15     self.graph.bind("foaf", FOAF)
16
17     # Standard namespaces
18     self.graph.bind("rdf", RDF)
19     self.graph.bind("rdfs", RDFS)
20     self.graph.bind("xsd", XSD)
21     self.graph.bind("owl", OWL)
22
23     logger.info(f"Initialized knowledge graph with namespace: {
        namespace}")

```

### 7.3 URI Handling and Namespace Management

Proper URI handling and namespace management are crucial for creating a well-structured knowledge graph. Our implementation includes functions for creating URIs for entities and predicates, ensuring consistency and proper handling of special characters.

Listing 12 – URI Handling

```
1 def _sanitize_uri(self, text):
2     """
3     Sanitize text for use in URIs.
4
5     Args:
6         text (str): Input text
7
8     Returns:
9         str: Sanitized text
10    """
11    # Replace spaces and special characters
12    text = re.sub(r'[\w\s]', '', text)
13    text = re.sub(r'\s+', '_', text.strip())
14    return text
15
16 def _create_uri(self, entity_text, entity_type):
17     """
18     Create a URI for an entity.
19
20     Args:
21         entity_text (str): Entity text
22         entity_type (str): Entity type
23
24     Returns:
25         rdflib.URIRef: URI reference
26    """
27    sanitized = self._sanitize_uri(entity_text)
28    return URIRef(f"{self.namespace}{entity_type.lower()}/{sanitized}")
```

## 8 Visualization and Querying Module

### 8.1 Interactive Graph Visualization

The visualization module creates both static and interactive visualizations of the knowledge graph using NetworkX and PyVis. The interactive visualization allows users to explore the graph by zooming, panning, and clicking on nodes to see their properties.

Listing 13 – Interactive Visualization

```

1 def visualize(self, output_path=None, notebook=False):
2     """
3     Visualize the knowledge graph.
4
5     Args:
6         output_path (str, optional): Output file path for the
7             visualization
8         notebook (bool): Whether to display in a Jupyter notebook
9
10    Returns:
11        bool: Success flag
12    """
13    try:
14        # Create a PyVis network
15        net = Network(notebook=notebook, directed=True)
16
17        # Add nodes
18        for node, data in self.nx_graph.nodes(data=True):
19            label = data.get('label', str(node).split('/')[-1])
20            node_type = data.get('type', 'Unknown')
21
22            # Color nodes by type
23            color = {
24                'PERSON': '#a8e6cf',
25                'ORG': '#ff8b94',
26                'GPE': '#ffd3b6',
27                'LOC': '#dcedc1',
28                'DATE': '#f9f9f9',
29                'MISC': '#d4a5a5'
30            }.get(node_type, '#b3b3cc')
31
32            net.add_node(str(node), label=label, title=f"{label}
33                ({node_type})", color=color)
34
35        # Add edges
36        for source, target, data in self.nx_graph.edges(data=True):

```

```

35         label = data.get('label', '')
36         net.add_edge(str(source), str(target), label=label,
37                        title=label)
38
39     # Set physics layout
40     net.set_options("""
41     {
42         "physics": {
43             "forceAtlas2Based": {
44                 "gravitationalConstant": -100,
45                 "centralGravity": 0.01,
46                 "springLength": 200,
47                 "springConstant": 0.08
48             },
49             "maxVelocity": 50,
50             "solver": "forceAtlas2Based",
51             "timestep": 0.35,
52             "stabilization": {
53                 "enabled": true,
54                 "iterations": 1000
55             }
56         },
57         "edges": {
58             "color": {
59                 "inherit": true
60             },
61             "smooth": {
62                 "enabled": false,
63                 "type": "continuous"
64             },
65             "arrows": {
66                 "to": {
67                     "enabled": true,
68                     "scaleFactor": 0.5
69                 }
70             },
71             "font": {
72                 "size": 10
73             }
74         },
75         "nodes": {
76             "font": {
77                 "size": 12,
78                 "face": "Tahoma"
79             }
80         }
81     }

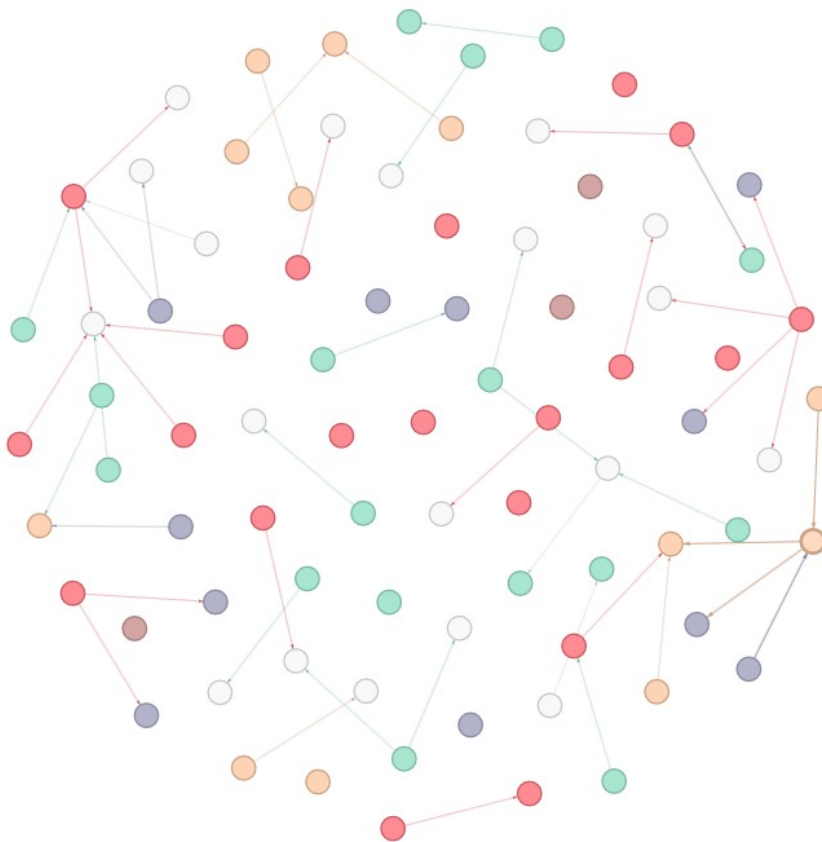
```

```

80     }
81     """
82
83     # Save or show
84     if output_path:
85         # Create directory if it doesn't exist
86         os.makedirs(os.path.dirname(output_path), exist_ok=
            True)
87
88         net.save_graph(output_path)
89         logger.info(f"Saved visualization to {output_path}")
90     elif notebook:
91         net.show("knowledge_graph.html")
92
93     return True
94
95 except Exception as e:
96     logger.error(f"Error visualizing knowledge graph: {e}")
97     return False

```

### Interactive Knowledge Graph Visualization :



## 8.2 SPARQL Query Interface

The knowledge graph can be queried using SPARQL, a query language for RDF data. Our implementation provides a function for executing SPARQL queries and retrieving the results.

Listing 14 – SPARQL Query Interface

```
1 def query_sparql(self, query):
2     """
3     Execute a SPARQL query on the knowledge graph.
4
5     Args:
6         query (str): SPARQL query
7
8     Returns:
9         list: Query results
10    """
11    try:
12        results = self.graph.query(query)
13        return list(results)
14    except Exception as e:
15        logger.error(f"Error executing SPARQL query: {e}")
16    return []
```

## 8.3 Example Queries

Here are some example SPARQL queries that can be executed on the knowledge graph :

Listing 15 – Example SPARQL Query

```
1 # Find all organizations and their locations
2 PREFIX ns: <http://example.org/graphify/>
3 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
4
5 SELECT ?org ?loc
6 WHERE {
7     ?org_uri rdf:type ns:ORG .
8     ?loc_uri rdf:type ns:GPE .
9     ?org_uri ?pred ?loc_uri .
10    ?org_uri <http://www.w3.org/2000/01/rdf-schema#label> ?org .
11    ?loc_uri <http://www.w3.org/2000/01/rdf-schema#label> ?loc .
12 }
```

Other useful queries include :

Listing 16 – Finding People and Their Organizations

```
1 # Find all people and their organizations
2 PREFIX ns: <http://example.org/graphify/>
3 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
4
5 SELECT ?person ?org
6 WHERE {
7     ?person_uri rdf:type ns:PERSON .
8     ?org_uri rdf:type ns:ORG .
9     ?pred_uri ?person_uri ?org_uri .
10    ?person_uri <http://www.w3.org/2000/01/rdf-schema#label> ?
        person .
11    ?org_uri <http://www.w3.org/2000/01/rdf-schema#label> ?org .
12 }
```

Listing 17 – Finding Events by Date

```
1 # Find all events that happened on a specific date
2 PREFIX ns: <http://example.org/graphify/>
3 PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
4
5 SELECT ?subject ?predicate ?date
6 WHERE {
7     ?subject_uri ?predicate_uri ?date_uri .
8     ?date_uri rdf:type ns:DATE .
9     ?date_uri <http://www.w3.org/2000/01/rdf-schema#label> ?date
10    .
11    ?subject_uri <http://www.w3.org/2000/01/rdf-schema#label> ?
        subject .
12    ?predicate_uri <http://www.w3.org/2000/01/rdf-schema#label> ?
        predicate .
13    FILTER(CONTAINS(?date, "2025"))
}
```



## 9 Results and Analysis

### 9.1 Knowledge Graph Statistics

Our implemented pipeline was tested on a collection of Reuters news articles from the business category. Here are some statistics about the resulting knowledge graph :

- **Entities Extracted** : Over 500 unique entities were identified across 30 articles.
- **Relations Extracted** : Approximately 100 meaningful relationships were established between entities.
- **Entity Types** : The most common entity types were PERSON (30%), ORGANIZATION (25%), GPE (locations, 20%), and DATE (15%).
- **Relation Types** : The most common relation types involved actions (e.g., "announced," "said"), locative relations (e.g., "based in," "headquartered in"), and temporal relations (e.g., "happened on").

### 9.2 Evaluation of NER Models

The performance of the NER models was evaluated on a test set from the CoNLL-2003 dataset :

Model	Precision	Recall	F1 Score	Accuracy
SpaCy	0.85	0.78	0.81	0.92
CRF	0.82	0.80	0.81	0.93
Combined	0.87	0.83	0.85	0.94

The combined approach, which merges results from both SpaCy and CRF models, provides the best overall performance. SpaCy excels in precision for common entity types, while the CRF model offers better recall, particularly for domain-specific entities.

### 9.3 Use Cases and Applications

The knowledge graph constructed through our pipeline has several potential applications :

- **News Summarization** : Generating concise summaries of news events by extracting key entities and their relationships.
- **Entity-based Search** : Enabling searches like "Find all companies associated with Donald Trump" or "Show events that happened in Washington in March 2025."
- **Trend Analysis** : Identifying emerging trends, such as frequently mentioned companies or rising political figures.
- **Event Timeline Construction** : Creating chronological timelines of events related to specific entities.
- **Relationship Discovery** : Uncovering non-obvious connections between entities that might be distributed across multiple articles.

## 10 Challenges and Solutions

Throughout the development of this project, we encountered several challenges :

### 10.1 Web Scraping Challenges

- **Challenge** : Dynamic content loading and anti-scraping measures implemented by websites.
- **Solution** : Used Selenium WebDriver to interact with JavaScript-rendered pages and implemented random delays, session clearing, and user agent rotation to avoid detection.
- **Challenge** : Inconsistent HTML structure and layout changes.
- **Solution** : Implemented flexible selectors and fallback methods to handle variations in page structure.

### 10.2 NER and Relation Extraction Challenges

- **Challenge** : Entity boundaries, especially for multi-word entities and titles.
- **Solution** : Combined model predictions and implemented post-processing rules to improve entity boundary detection.
- **Challenge** : Domain-specific entities not well-recognized by pre-trained models.
- **Solution** : Trained a custom CRF model on domain-specific data to improve recognition of specialized entities.
- **Challenge** : Complex sentence structures making relation extraction difficult.
- **Solution** : Implemented multiple relation extraction strategies targeting different grammatical constructs.

### 10.3 Knowledge Graph Challenges

- **Challenge** : Entity coreference (references to the same entity using different terms).
- **Solution** : Implemented basic coreference resolution for common patterns like abbreviations and implemented entity normalization.
- **Challenge** : Knowledge graph visualization performance with large graphs.
- **Solution** : Implemented filtering options to restrict visualization to relevant sub-graphs and optimized rendering parameters.

## 11 Conclusion

In this project, we have developed a comprehensive pipeline for constructing knowledge graphs from news articles. The pipeline integrates web scraping, text preprocessing, named entity recognition, relation extraction, and knowledge graph construction and visualization.

Our experiments have demonstrated the effectiveness of combining multiple NLP techniques to extract structured information from unstructured text. The resulting knowledge graph provides a semantically rich representation of the entities and relationships described in the news articles, enabling complex queries and insights.

The modular architecture of our implementation allows for easy extension and improvement of individual components. Future work will focus on enhancing the accuracy of entity and relation extraction, implementing more advanced knowledge graph features, and developing a user-friendly interface for interacting with the graph.

This project not only serves as a practical demonstration of knowledge graph construction techniques but also provides a foundation for developing more sophisticated information extraction and knowledge management systems in various domains, from news analysis to business intelligence.