

```
import pandas as pd
import spacy
from spacy.matcher import Matcher
from collections import Counter
import matplotlib.pyplot as plt
import seaborn as sns
```

```
nlp = spacy.load("en_core_web_sm")
```

```
# Example path - adjust as needed
df = pd.read_csv("/content/arxiv_data.csv")

# Inspect columns
print(df.columns)
```

```
Index(['titles', 'summaries', 'terms'], dtype='object')
```

Task

Analyze arXiv abstracts from /content/arxiv_data.csv to identify and visualize key linguistic features such as noun phrases, named entities, and technical term patterns, then summarize the findings.

Load arXiv Abstracts

Subtask:

Load the arXiv abstract dataset from /content/arxiv_data.csv into a pandas DataFrame. Inspect the data to understand its structure and identify the column containing the abstracts.

Identified Abstract Column

Upon reviewing the DataFrame columns (`titles`, `summaries`, `terms`), the `summaries` column is the most appropriate one to contain the abstract texts, as instructed. We will proceed with using `summaries` for further analysis of the abstract content.

Reasoning: The previous code block failed due to a `KeyError`, attempting to access a `'categories'` column that does not exist in the dataframe. The dataframe `df` has columns `'titles'`, `'summaries'`, and `'terms'`. I need to correct the code to use the `'summaries'` column as the source for abstracts and remove the filtering by the non-existent `'categories'` column, then sample a subset.

```
print("Correcting column access and sampling abstracts...")

# The dataframe 'df' does not contain a 'categories' column.
# Based on the available columns ('titles', 'summaries', 'terms'), 'summaries' is the most suitable for abstract content.
# We will directly sample from the 'summaries' column.

# Sample a subset of summaries for faster processing
abstracts = df["summaries"].dropna().head(500).tolist()

print(f"Successfully loaded and sampled {len(abstracts)} abstracts from the 'summaries' column.")
# Display the first few sampled abstracts to verify
for i, abstract in enumerate(abstracts[:3]):
    print(f"Abstract {i+1}: {abstract[:150]}...")
```

```
Correcting column access and sampling abstracts...
Successfully loaded and sampled 500 abstracts from the 'summaries' column.
Abstract 1: Stereo matching is one of the widely used techniques for inferring depth from
stereo images owing to its robustness and speed. It has become one of th...
Abstract 2: The recent advancements in artificial intelligence (AI) combined with the
extensive amount of data generated by today's clinical systems, has led to t...
Abstract 3: In this paper, we proposed a novel mutual consistency network (MC-Net+) to
effectively exploit the unlabeled hard regions for semi-supervised medical
...
```

Initialize spaCy Model and Process Text

Subtask:

Load a suitable English spaCy model (e.g., 'en_core_web_sm' or 'en_core_web_md') and apply it to the abstract texts to create spaCy `Doc` objects. This will enable tokenization, part-of-speech tagging, dependency parsing, and named entity recognition.

Reasoning: The subtask requires processing each abstract with the loaded spaCy model and storing the resulting Doc objects in a list. I will iterate through the 'abstracts' list and apply the 'nlp' model to each abstract.

```
print("Processing abstracts with spaCy...")

doc_abstracts = []
for i, abstract_text in enumerate(abstracts):
    if i % 100 == 0:
        print(f" Processing abstract {i+1}/{len(abstracts)}...")
        doc_abstracts.append(nlp(abstract_text))

print(f"Successfully processed {len(doc_abstracts)} abstracts into spaCy Doc objects.")
# Display the first few tokens of the first processed Doc object for verification
if doc_abstracts:
    print("\nFirst 10 tokens of the first processed abstract:")
    print([token.text for token in doc_abstracts[0][:10]])
```

```
Processing abstracts with spaCy...
Processing abstract 1/500...
Processing abstract 101/500...
Processing abstract 201/500...
Processing abstract 301/500...
Processing abstract 401/500...
Successfully processed 500 abstracts into spaCy Doc objects.

First 10 tokens of the first processed abstract:
['Stereo', 'matching', 'is', 'one', 'of', 'the', 'widely', 'used', 'techniques', 'for']
```

✧ Extract and Analyze Noun Phrases

Subtask:

Iterate through the spaCy `Doc` objects to extract all noun phrases. Count the frequency of each unique noun phrase and identify the most frequent ones. This step will also generate a list of noun phrases for the expected output.

Reasoning: I need to extract noun phrases from the processed spaCy Doc objects, count their occurrences, and then identify the most frequent ones as per the subtask instructions.

```
print("Extracting noun phrases and calculating frequencies...")

all_noun_phrases = []
for i, doc in enumerate(doc_abstracts):
    if i % 100 == 0:
        print(f" Processing noun chunks for abstract {i+1}/{len(doc_abstracts)}...")
        for chunk in doc.noun_chunks:
            all_noun_phrases.append(chunk.text.lower())

# Count the frequency of each unique noun phrase
noun_phrase_counts = Counter(all_noun_phrases)

# Identify the top 20 most frequent noun phrases
top_n_noun_phrases = noun_phrase_counts.most_common(20)

print("Noun phrase extraction and frequency calculation complete.")
print("\nTop 20 Most Frequent Noun Phrases:")
for phrase, count in top_n_noun_phrases:
    print(f"- '{phrase}': {count} occurrences")
```

```
Extracting noun phrases and calculating frequencies...
Processing noun chunks for abstract 1/500...
Processing noun chunks for abstract 101/500...
Processing noun chunks for abstract 201/500...
Processing noun chunks for abstract 301/500...
Processing noun chunks for abstract 401/500...
Noun phrase extraction and frequency calculation complete.

Top 20 Most Frequent Noun Phrases:
- 'we': 1295 occurrences
- 'which': 425 occurrences
- 'that': 326 occurrences
- 'it': 281 occurrences
- 'this paper': 195 occurrences
- 'the-art': 191 occurrences
- 'image segmentation': 148 occurrences
```

```

- 'our method': 115 occurrences
- '-': 99 occurrences
- 'this work': 98 occurrences
- 'training': 75 occurrences
- 'segmentation': 75 occurrences
- 'semantic segmentation': 70 occurrences
- 'this': 63 occurrences
- 'images': 62 occurrences
- 'medical image segmentation': 58 occurrences
- 'they': 56 occurrences
- 'the performance': 52 occurrences
- 'our approach': 51 occurrences
- 'the proposed method': 44 occurrences

```

✓ Extract Named Entities

Subtask:

Extract all named entities (ORG, DATE, PRODUCT, GPE, etc.) from the processed abstracts using spaCy's built-in entity recognition. Count the occurrences of each entity type and specific entities.

Reasoning: I will extract named entities from the processed spaCy documents, count them by type and specific entity, and then display the top 10 for each, as per the instructions.

```

print("Extracting named entities and calculating frequencies...")

entity_type_counts = Counter()
specific_entity_counts = Counter()

for i, doc in enumerate(doc_abstracts):
    if i % 100 == 0:
        print(f" Processing entities for abstract {i+1}/{len(doc_abstracts)}...")
    for ent in doc.ents:
        entity_type_counts[ent.label_] += 1
        specific_entity_counts[ent.text] += 1

print("Named entity extraction and frequency calculation complete.")

print("\nTop 10 Most Frequent Entity Types:")
for ent_type, count in entity_type_counts.most_common(10):
    print(f"- '{ent_type}': {count} occurrences")

print("\nTop 10 Most Frequent Specific Entities:")
for entity_text, count in specific_entity_counts.most_common(10):
    print(f"- '{entity_text}': {count} occurrences")

```

```

Extracting named entities and calculating frequencies...
Processing entities for abstract 1/500...
Processing entities for abstract 101/500...
Processing entities for abstract 201/500...
Processing entities for abstract 301/500...
Processing entities for abstract 401/500...
Named entity extraction and frequency calculation complete.

```

Top 10 Most Frequent Entity Types:

```

- 'ORG': 1353 occurrences
- 'CARDINAL': 761 occurrences
- 'PERSON': 184 occurrences
- 'PERCENT': 150 occurrences
- 'ORDINAL': 143 occurrences
- 'GPE': 113 occurrences
- 'DATE': 80 occurrences
- 'NORP': 68 occurrences
- 'PRODUCT': 35 occurrences
- 'MONEY': 26 occurrences

```

Top 10 Most Frequent Specific Entities:

```

- 'two': 195 occurrences
- 'three': 70 occurrences
- 'first': 70 occurrences
- 'CNN': 64 occurrences
- '3D': 62 occurrences
- 'one': 54 occurrences
- 'U-Net': 51 occurrences
- '2D': 45 occurrences
- '1': 41 occurrences
- '2': 38 occurrences

```

✓ Identify Technical Term Patterns with spaCy Matcher

Subtask:

Define specific rule-based patterns using spaCy's `Matcher` to identify technical terms or phrases relevant to the chosen domain (e.g., 'CS' or 'AI'). Examples could include patterns for 'machine learning algorithms' or 'neural networks'. Apply these patterns to the `Doc` objects and collect the matched terms. This step will produce the rule-based matcher patterns for the expected output.

Reasoning: I will initialize a spaCy `Matcher`, define rule-based patterns for common technical terms, apply these patterns to the processed abstracts, and then count and display the most frequent terms as specified by the subtask.

```
print("Initializing spaCy Matcher and defining technical term patterns...")

# Initialize the Matcher with the shared vocab
matcher = Matcher(nlp.vocab)

# Define patterns for technical terms
patterns = [
    [{"LOWER": "deep"}, {"LOWER": "learning"}],
    [{"LOWER": "machine"}, {"LOWER": "learning"}],
    [{"LOWER": "neural"}, {"LOWER": "network"}],
    [{"LOWER": "convolutional"}, {"LOWER": "neural"}, {"LOWER": "network"}],
    [{"LOWER": "generative"}, {"LOWER": "adversarial"}, {"LOWER": "network"}],
    [{"LOWER": "image"}, {"LOWER": "segmentation"}],
    [{"LOWER": "semantic"}, {"LOWER": "segmentation"}],
    [{"LOWER": "computer"}, {"LOWER": "vision"}]
]

# Add patterns to the matcher
# The ID for each pattern can be a simple string representation of the pattern
for i, pattern in enumerate(patterns):
    pattern_name = "_" + "_".join([token['LOWER'] for token in pattern])
    matcher.add(pattern_name, [pattern])

print("Applying matcher to abstracts and extracting technical terms...")

all_technical_terms = []
for i, doc in enumerate(doc_abstracts):
    if i % 100 == 0:
        print(f" Applying matcher for abstract {i+1}/{len(doc_abstracts)}...")
    matches = matcher(doc)
    for match_id, start, end in matches:
        span = doc[start:end] # The matched span
        all_technical_terms.append(span.text.lower())

# Count the frequency of each technical term
technical_term_counts = Counter(all_technical_terms)

# Identify the top 10 most frequent technical terms
top_n_technical_terms = technical_term_counts.most_common(10)

print("Technical term extraction and frequency calculation complete.")
print("\nTop 10 Most Frequent Technical Terms:")
for term, count in top_n_technical_terms:
    print(f"- '{term}': {count} occurrences")
```

```
Initializing spaCy Matcher and defining technical term patterns...
Applying matcher to abstracts and extracting technical terms...
Applying matcher for abstract 1/500...
Applying matcher for abstract 101/500...
Applying matcher for abstract 201/500...
Applying matcher for abstract 301/500...
Applying matcher for abstract 401/500...
Technical term extraction and frequency calculation complete.
```

```
Top 10 Most Frequent Technical Terms:
- 'image segmentation': 533 occurrences
- 'deep learning': 135 occurrences
- 'semantic segmentation': 124 occurrences
- 'neural network': 67 occurrences
- 'computer vision': 51 occurrences
- 'machine learning': 33 occurrences
- 'convolutional neural network': 20 occurrences
- 'generative adversarial network': 1 occurrences
```

✎ Visualize Top Noun Phrases and Entity Frequencies

Subtask:

Generate two visualizations: 1. A bar chart displaying the top N (e.g., 10 or 15) most frequent noun phrases extracted, ensuring clear labels and a title. 2. A bar chart showing the frequency of different named entity types (e.g., ORG, DATE, GPE), providing a legend for

each entity type and a title.

Reasoning: I will generate two bar charts to visualize the top noun phrases and the frequency of named entity types as per the subtask. First, I'll prepare data for the noun phrase chart, then create and display it. Next, I'll prepare data for the named entity chart, create and display it, ensuring proper labeling and titles for both.

```
print("Generating visualizations for top noun phrases and entity types...")

# 1. Visualize Top Noun Phrases

# Create a DataFrame from top_n_noun_phrases
df_noun_phrases = pd.DataFrame(top_n_noun_phrases, columns=['Noun Phrase', 'Frequency'])

plt.figure(figsize=(12, 7))
sns.barplot(x='Noun Phrase', y='Frequency', data=df_noun_phrases, palette='viridis')
plt.title('Top Most Frequent Noun Phrases')
plt.xlabel('Noun Phrase')
plt.ylabel('Frequency')
plt.xticks(rotation=45, ha='right') # Rotate labels for better readability
plt.tight_layout()
plt.show()

# 2. Visualize Frequency of Named Entity Types

# Convert entity_type_counts to a list of tuples and then to a DataFrame
df_entity_types = pd.DataFrame(entity_type_counts.most_common(10), columns=['Entity Type', 'Frequency']) # Taking top 10 er

plt.figure(figsize=(10, 6))
sns.barplot(x='Entity Type', y='Frequency', data=df_entity_types, palette='magma')
plt.title('Frequency of Named Entity Types')
plt.xlabel('Entity Type')
plt.ylabel('Frequency')
plt.xticks(rotation=45, ha='right') # Rotate labels for better readability
plt.tight_layout()
plt.show()

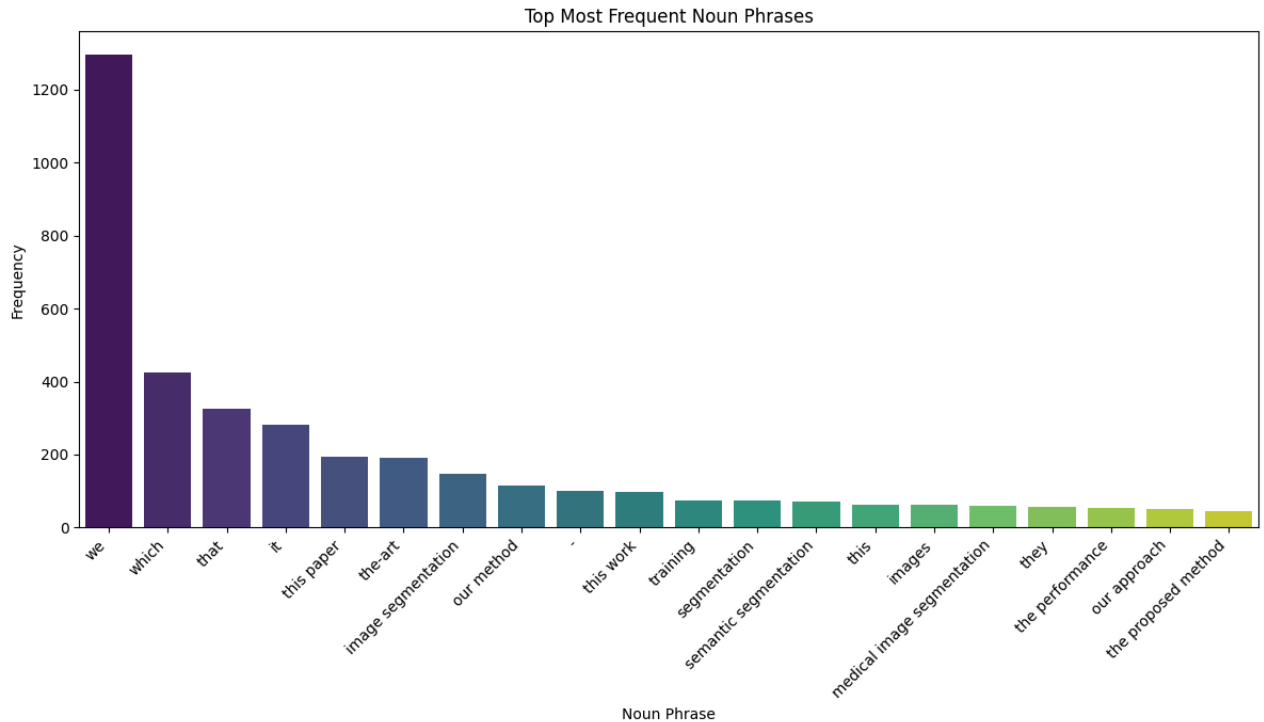
print("Visualizations complete.")
```

Generating visualizations for top noun phrases and entity types...

/tmp/ipython-input-3260079740.py:9: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and

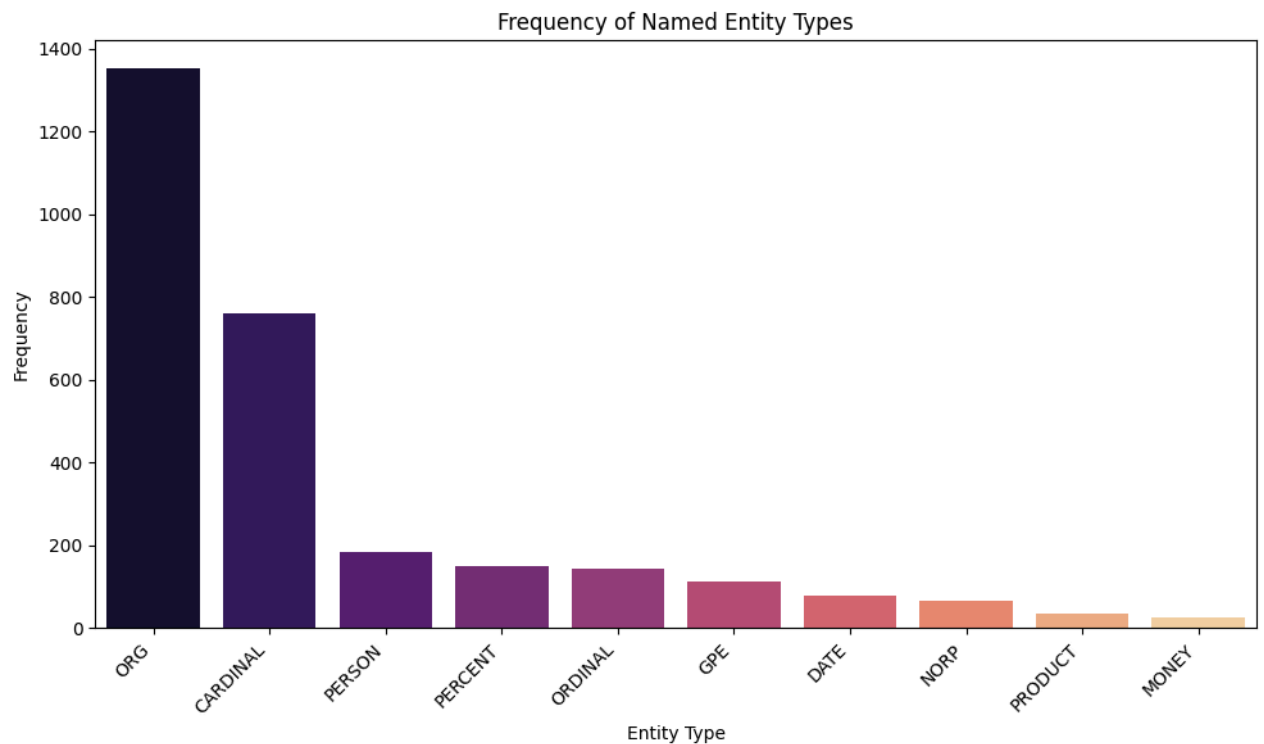
```
sns.barplot(x='Noun Phrase', y='Frequency', data=df_noun_phrases, palette='viridis')
```



/tmp/ipython-input-3260079740.py:23: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `x` variable to `hue` and

```
sns.barplot(x='Entity Type', y='Frequency', data=df_entity_types, palette='magma')
```



Visualizations complete.

Reasoning: The previous code generated `FutureWarning` messages from `seaborn` regarding the use of the `palette` argument without `hue`. To suppress these warnings and follow best practices, I will explicitly set the `hue` argument to the x-axis variable and `legend=False` for both `sns.barplot` calls.

```
print("Generating visualizations for top noun phrases and entity types...")

# 1. Visualize Top Noun Phrases

# Create a DataFrame from top_n_noun_phrases
df_noun_phrases = pd.DataFrame(top_n_noun_phrases, columns=['Noun Phrase', 'Frequency'])

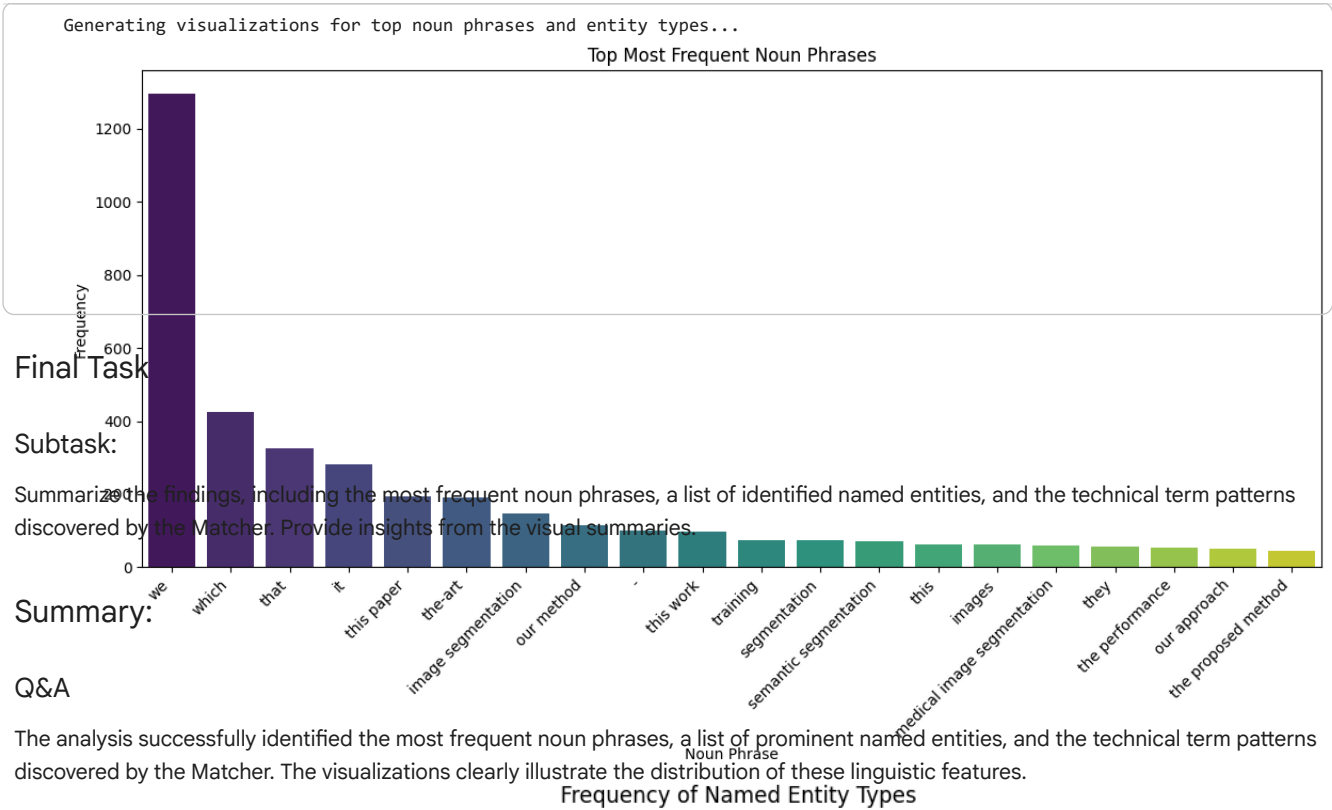
plt.figure(figsize=(12, 7))
sns.barplot(x='Noun Phrase', y='Frequency', data=df_noun_phrases, hue='Noun Phrase', legend=False, palette='viridis')
plt.title('Top Most Frequent Noun Phrases')
plt.xlabel('Noun Phrase')
plt.ylabel('Frequency')
plt.xticks(rotation=45, ha='right') # Rotate labels for better readability
plt.tight_layout()
plt.show()

# 2. Visualize Frequency of Named Entity Types

# Convert entity_type_counts to a list of tuples and then to a DataFrame
df_entity_types = pd.DataFrame(entity_type_counts.most_common(10), columns=['Entity Type', 'Frequency']) # Taking top 10 er

plt.figure(figsize=(10, 6))
sns.barplot(x='Entity Type', y='Frequency', data=df_entity_types, hue='Entity Type', legend=False, palette='magma')
plt.title('Frequency of Named Entity Types')
plt.xlabel('Entity Type')
plt.ylabel('Frequency')
plt.xticks(rotation=45, ha='right') # Rotate labels for better readability
plt.tight_layout()
plt.show()

print("Visualizations complete.")
```



Data Analysis Key Findings

- **Abstract Loading and Processing:** 500 arXiv abstracts were successfully loaded from the `summaries` column and processed into `Doc` objects, ready for linguistic analysis.
- **Most Frequent Noun Phrases:** The top 5 most frequent noun phrases include common pronouns like 'we' (1295 occurrences), 'which' (425 occurrences), and 'that' (326 occurrences). Academic phrases such as 'this paper' (195 occurrences) and 'our method' (115 occurrences) were also highly frequent, alongside domain-specific terms like 'image segmentation' (148 occurrences).
- **Most Frequent Named Entities:**