

Named Entity Recognition (NER)

Mohammad Matar 219200060

07.11.2024

Inhaltsverzeichnis

| | | |
|----------|-------------------------------------------------------------|-----------|
| 1 | Grundlage von NER-Konzepte | 4 |
| 1.1 | Intro | 4 |
| 1.2 | Motivation und Zielsetzung | 4 |
| 1.3 | Definition der NER | 5 |
| 2 | Einsatz von Tools | 6 |
| 2.1 | Überblick | 6 |
| 2.2 | Auswahl der Tools | 10 |
| 2.3 | Einsatz des Spacy-Tools | 12 |
| 2.4 | Einsatz des Flair-Tools | 13 |
| 2.5 | Einsatz des weiter zu trainierenden Flair-Modells | 15 |
| 2.6 | Evaluation der eingesetzten Modelle | 21 |
| 3 | Zusammenfassung und Ausblick | 24 |
| 3.1 | Zusammenfassung | 25 |
| 3.2 | Ausblick | 25 |
| 4 | Python-Codes zum Einsatz der Tools | 27 |
| 4.1 | Trainingsdatenvorverarbeitung | 27 |
| 4.2 | NER-Modelltraining | 31 |
| 4.3 | NER-Klassifikation | 33 |
| 4.4 | NER Evaluation | 37 |

Abbildungsverzeichnis

| | | |
|---|--------------------------------------------------------------------------------------------------------------------------------------------------------------|----|
| 1 | Die Trainingsdaten wurden aus Texten von einem bestimmten Absender erstellt, während die Evaluationsdaten von einem einzigen Absender stammen | 22 |
| 2 | Die Trainingsdaten wurden aus einer Mischung von Daten mehrerer Absender erstellt, während die Evaluationsdaten von einem einzigen Absender stammen. | 23 |
| 3 | Sowohl die Trainingsdaten als auch die Evaluationsdaten wurden aus einer Mischung von Daten verschiedener Absender zusammengestellt. | 23 |

Listings

| | | |
|---|----------------------------------------------------|----|
| 1 | Datenverarbeitungsskript für NER-Modelle | 27 |
|---|----------------------------------------------------|----|

| | | |
|---|--------------------------------------------------|----|
| 2 | NER-Modelltraining mit Flair | 31 |
| 3 | NER-Klassifikation mit Flair und Spacy | 33 |
| 4 | NER Evaluation Pipeline | 37 |

1 Grundlage von NER-Konzepte

1.1 Intro

Derzeit ist eine große Menge an digitalen Daten (wie E-Mails, soziale Anwendungen, Zeitungen und Instagram) in verschiedenen Sprachen verfügbar. Diese Informationen werden sowohl in strukturierter als auch unstrukturierter Form gesammelt, um die Daten zu verarbeiten und nützliche Informationen zu extrahieren. Es ist jedoch eine der größten Herausforderungen, bedeutungsvolle Kenntnisse aus solchen riesigen Datenmengen zu extrahieren. Das Hauptziel der natürlichen Sprachverarbeitung (Natural Language Processing, NLP) besteht darin, nützliche Informationen aus den menschlichen natürlichen und Informationssprachen zu gewinnen, damit Maschinen die Informationen der menschlichen Sprache besser verstehen und entsprechend handeln können. Viele Informationssysteme zur Extraktion von Informationen in NLP wurden entwickelt, die automatisch Fragen beantworten und Texte zusammenfassen [8]. Dazu gehört das sogenannte Information Extraction (IE), eine der wichtigen Aufgaben in der Textanalyse und der Verarbeitung natürlicher Sprache (Natural Language Processing, NLP), die das Extrahieren von bedeutungsvollen Wissensfragmenten aus unstrukturierten Informationsquellen umfasst, da unstrukturierte Daten rechnerisch undurchsichtig sind. Das Ziel von IE ist es, eine Wissensbasis zu erstellen, d.h. die Informationen so zu organisieren, dass sie für Menschen nützlich sind und in einer semantischen Weise angeordnet werden, sodass Algorithmen nützliche Schlussfolgerungen daraus ziehen können [5]. Eine der wesentlichen Unteraufgaben von IE ist das Named Entity Recognition (NER), bei dem es darum geht, Namen/Entitäten zu finden und zu klassifizieren. Sobald diese Named Entities (NE) extrahiert sind, können sie indexiert und durchsuchbar gemacht werden, Beziehungen abgeleitet, Fragen beantwortet und vieles mehr. Die NER-Techniken unterscheiden sich je nach Domain, aufgrund der Einzigartigkeit, die in jeder Domäne existiert, obwohl der Prozess auf einer Reihe grundlegender Schritte der natürlichen Sprachverarbeitung (NLP) basiert, wie z.B. Tokenisierung, Part-of-Speech-Tagging, Parsing und Modellbildung [5].

1.2 Motivation und Zielsetzung

Die zentrale Herausforderung dieser Arbeit liegt im "WossiDi"-Archiv, das die digitale Version des Zettelkastensystems von Richard Wossidlo (1859–1939) zur Verfügung stellt. Wossidlo, ein wegweisender Feldforscher, reiste durch ganz Mecklenburg und schuf – unterstützt von zahlreichen Sammelhelfern –

eine einzigartige ethnografische und regionalsprachliche Sammlung. Das Archiv umfasst eine Zettelsammlung, die nach Volksüberlieferungen, Arbeits- und Lebensbereichen systematisch verschlagwortet ist, sowie die dazugehörige Beiträgerkorrespondenz. Zudem enthält es wertvolle Hinweise zu überregionalen Traditionen und den alphabetischen Zettelkatalog Wossidlos, ergänzt durch den abgeleiteten Katalog von Hermann Teuchert, der als Grundlage für das Mecklenburgische Wörterbuch dient. Weitere Bestandsgruppen aus dem Umfeld des Wossidlo-Archivs bereichern das Material. WossiDiA stellt etwa zwei Millionen überwiegend handschriftliche Dokumente im freien Zugang online zur Verfügung. Diese sind durch Schlagwörter, Personen, Orte und Zeitangaben durchsuchbar. Transkriptionen sowie Erklärungen, etwa zu Kürzeln, erleichtern die Nutzung und ermöglichen eine benutzerfreundliche und effiziente Recherche [9]. Die Zielsetzung dieser Arbeit ist es, die Prinzipien der Named Entity Recognition (NER) zu erläutern und deren praktischen Nutzen anhand der Anwendung auf die transkribierten Texte der Korrespondenzen von Richard Wossidlo zu demonstrieren. Dabei wird untersucht, wie NER-Tools eingesetzt werden können, um historische Dokumente, insbesondere handschriftliche Transkripte, zu analysieren und strukturierte Informationen wie Namen von Personen, Orten und Zeitangaben automatisch zu extrahieren. Die Arbeit soll auch aufzeigen, wie NER zur Identifikation und Kategorisierung von Entitäten in den Korrespondenzen von Wossidlo beitragen kann, um deren Analyse und Interpretation zu erleichtern. Ziel ist es, die Möglichkeiten und Herausforderungen der Anwendung von NER auf historische Texte zu verdeutlichen und zu zeigen, wie diese Technologien die Forschung zu historischen und geisteswissenschaftlichen Datenbereichen verbessern können.

1.3 Definition der NER

Named Entity Recognition (NER) ist eine zentrale Aufgabe in der Verarbeitung natürlicher Sprache (Natural Language Processing, NLP) und Information Extraction (IE). Es handelt sich dabei um den Prozess, bei dem benannte Entitäten wie Personennamen, Organisationen, Orte, Zeitangaben und andere bedeutende Begriffe aus unstrukturierten Texten extrahiert und klassifiziert werden. Das Hauptziel von NER besteht darin, aus unstrukturierten Daten strukturierte Informationen zu gewinnen, die für die weitere Analyse und maschinelle Verarbeitung von Bedeutung sind. Diese Entitäten werden in vordefinierte Kategorien eingeteilt, um den Text in einer maschinenlesbaren Form darzustellen, die für Algorithmen und Wissensdatenbanken nutzbar ist. Die Identifikation von benannten Entitäten ermöglicht es, wichtige Informationen zu extrahieren, die für Anwendungen wie Fragebeant-

wortung, Wissensgraphen und Textzusammenfassungen verwendet werden können [8] [5] [2]. Ein Beispiel für die Anwendung von NER in der Praxis wäre der Satz: "Jhon kaufte 500 Aktien von Acme Corp. im Jahr 2016." In diesem Fall werden "Jhon" als Person, "Acme Corp." als Organisation und "2016" als Zeitangabe erkannt und entsprechend klassifiziert [8].

2 Einsatz von Tools

In diesem Kapitel wurden verschiedene Tools für die NER untersucht und auf die spezifischen Anforderungen der historischen Daten des Projekts angewendet. Ziel ist es, die Leistung und Funktionalität der Tools zu evaluieren und das beste Werkzeug für zukünftige Anwendungen auszuwählen.

2.1 Überblick

In der Abschnitt wird ein umfassender Einblick in die Zielsetzung und die zentralen Aspekte der NER für historische Texte gegeben. Zunächst werden die zu erkennenden Entitäten definiert und ihre Relevanz für die Analyse erläutert. Anschließend werden die verschiedenen Eingabe- und Ausgabeformate beschrieben, um eine strukturierte Verarbeitung der Daten zu ermöglichen. Dieser Überblick dient als Grundlage für die weitere Analyse und Bewertung der eingesetzten NER-Tools.

2.1.1 Zielsetzung

Das Ziel dieses Kapitels besteht darin, verschiedene Tools für die NER einzusetzen, um deren Funktionalität und Leistung auf den spezifischen historischen Daten unseres Projekts zu bewerten. Dabei liegt der Fokus auf:

- Der Evaluation der Genauigkeit und Effizienz der Tools bei der Erkennung und Klassifizierung benannter Entitäten in historischen Texten.
- Dem Vergleich mehrerer Tools, um deren Stärken und Schwächen im Hinblick auf die Verwendbarkeit für unsere Daten zu analysieren.
- Der Identifikation des besten Ansatzes für zukünftige Anwendungen auf ähnliche Datenquellen.

2.1.2 Zu erkennenden Entitäten

Im Rahmen des Projekts sind die folgenden Entitätstypen definiert, die aus den Texten extrahiert werden sollen. Jede Kategorie wird basierend auf ih-

rer Relevanz für die Analyse und den Kontext der historischen Dokumente beschrieben:

Orte (LOC)

- Beschreibung: Bezieht sich auf geografische Standorte wie Städte, Länder oder Regionen.
- Beispiele: `<LOC>Rostock</LOC>`, `<LOC>Mecklenburg</LOC>`.

Personen (PER)

- Beschreibung: Namen von Individuen, einschließlich historischer Persönlichkeiten oder Autoren.
- Beispiele: `<PER>E Geinitz</PER>`, `<PER>Herr Dr. Krause</PER>`.

Organisationen (ORG)

- Beschreibung: Institutionen oder Organisationen, die in den Texten erwähnt werden.
- Beispiele: `<ORG>Statistisches Amt</ORG>`.

Ereignisse (EVENT)

- Beschreibung: Zeitbezogene oder thematische Ereignisse, die in den Dokumenten erwähnt werden.
- Beispiele: `<EVENT>7.1.08</EVENT>`, `<EVENT>Kommissionssitzung</EVENT>`.

Tiere (ANIMAL)

- Beschreibung: Erwähnte Tiere oder Tierarten in den Texten.
- Beispiele: `<ANIMAL>Pferd</ANIMAL>`, `<ANIMAL>Schwan</ANIMAL>`.

Die Definition dieser Entitätenkategorien ermöglicht eine strukturierte Extraktion der wichtigsten Informationen aus den Texten und bildet die Grundlage für die nachfolgende Analyse.

2.1.3 Eingabedatenformat

Die Eingabedaten, die bei der Verarbeitung der historischen Texten bereitgestellt werden, liegen in zwei Formaten, wo jeder Text in den beiden Formaten verfügbar ist, vor:

- **Annotierte-Format:** Ein strukturiertes Textformat, bei dem benannte Entitäten bereits durch Tags markiert sind. Beispiel:

```
<LOC>Rostock</LOC>, <EVENT>7.1.08</EVENT>.  
Sehr geehrter <PER>Herr Dr.</PER>  
Möchten Sie die Güte haben, mir für den Druck der „Beispiel“-Seite  
Ihre festen Vorschläge aufschreiben?
```

- **Freitext-Format:** Texte ohne vorgegebene Struktur, die roh verarbeitet werden müssen. Beispiel:

```
4.11.14 Rostock, 6. 11. 08  
Die Kommissionssitzung ist nunmehr auf  
Sonnabend, d. 8. Nachm. nach 3 Uhr in Güstrow  
festgesetzt.
```

2.1.4 Ausgabedatenformat

Die Ergebnisse der NER-Tools werden in einem JSON-Format bereitgestellt. Dieses Format strukturiert die extrahierten Informationen und macht sie für die Weiterverarbeitung nutzbar. Jedes Dokument wird durch die folgenden Attribute beschrieben:

- **ID (id):** Eine eindeutige Identifikationsnummer des Dokuments, die es ermöglicht, die Ergebnisse den entsprechenden Eingabedaten zuzuordnen.
- **Absender (sender):** Der Absender des Dokuments basierend auf dem Dokumentnamen.
- **Text (text):** Der Originaltext des Dokuments, der für die NER verwendet wurde.
- **Entitäten (entities):** Eine Liste der erkannten benannten Entitäten, wobei jede Entität durch die folgenden Attribute beschrieben wird:

- **Textausschnitt:** Der exakte Text, der als benannte Entität erkannt wurde.
- **Kategorie:** Die Klassifikation der Entität, wie z. B. LOC (Ort), PER (Person) oder EVENT (Ereignis).
- **Startposition:** Der Index des ersten Zeichens der Entität im Originaltext.
- **Endposition:** Der Index des letzten Zeichens der Entität im Originaltext.

Beispielausgabe:

```
{
  "results": [
    {
      "id": "21776",
      "sender": "Geinitz",
      "text": "Rostock, 6. 11. 08...",
      "entities": [
        [
          "Rostock",
          "LOC",
          0,
          7
        ],
        [
          "Herrn Peltz",
          "PER",
          199,
          210
        ]
      ]
    }
  ]
}
```

Die JSON-Struktur erlaubt es, die Ergebnisse einfach zu speichern, zu durchsuchen oder in anderen Anwendungen weiterzuverarbeiten.

2.2 Auswahl der Tools

Für NER in diesem Projekt wurden verschiedene Tools analysiert, um diejenigen auszuwählen, die den Anforderungen am besten entsprechen.

2.2.1 Kriterien für die Auswahl

Die folgenden Kriterien wurden verwendet, um die am besten geeigneten Tools für die NER auszuwählen. Diese Kriterien gewährleisten, dass die Tools sowohl den technischen Anforderungen als auch den spezifischen Bedürfnissen des Projekts gerecht werden:

- **Genauigkeit:** Wie präzise und vollständig das Tool benannte Entitäten erkennt und klassifiziert.
- **Anpassungsfähigkeit:** Die Möglichkeit, das Tool an domänenspezifische oder historische Texte anzupassen.
- **Unterstützung der deutschen Sprache:** Das Tool muss vortrainierte Modelle oder eine einfache Anpassung für die Verarbeitung deutscher Texte bieten, da die Analyse historischer Dokumente in deutscher Sprache erfolgt.
- **Unterstützung neuer Entitäten:** Das Tool sollte es ermöglichen, zusätzliche oder benutzerdefinierte Entitätsklassen (z. B. **ANIMAL** oder **EVENT**) effizient zu integrieren und zu trainieren.
- **Benutzerfreundlichkeit:** Wie einfach die Installation, Konfiguration und Nutzung des Tools ist, insbesondere für Benutzer ohne tiefgehende technische Vorkenntnisse.
- **Performance:** Wie schnell und ressourcenschonend das Tool große Mengen an Textdaten verarbeiten kann, ohne die Genauigkeit der Analyse zu beeinträchtigen.

2.2.2 Ausgewählte einzusetzende Tools

Für die Erkennung von den Entitäten in den Texten werden gezielt folgende Tools basierend auf den oben vorgestellten Kriterien ausgewählt und eingesetzt:

- **Spacy:**

- **Beschreibung:** Spacy ist ein Open-Source-Softwarepaket für fortgeschrittene natürliche Sprachverarbeitung in Python. Es bietet schnelle und effiziente vortrainierte Modelle für NER und andere natürliche Sprachverarbeitungsaufgaben [3].
 - **Stärken:**
 - * Sehr schnelle Verarbeitungsgeschwindigkeit, auch für große Textmengen [3].
 - * Benutzerfreundlich und gut dokumentiert mit einer aktiven Community [3].
 - **Schwächen:**
 - * Begrenzte Flexibilität für domänenspezifische Anpassungen ohne zusätzliches Training [3].
 - * Modelle sind primär auf moderne und nicht auf historische Texte optimiert [3].
- **Flair:**
- **Beschreibung:** Flair ist ein einfach zu verwendendes Framework für modernste natürliche Sprachverarbeitung, entwickelt von der Humboldt-Universität zu Berlin. Es bietet vortrainierte Modelle für NER und ermöglicht benutzerdefinierte Anpassungen[1].
 - **Stärken:**
 - * Hohe Genauigkeit durch kontextuelle Einbettungen wie *Flair Embeddings* [1].
 - * Flexibilität für domänenspezifische Anwendungen durch Training eigener Modelle [1].
 - **Schwächen:**
 - * Vergleichsweise langsame Verarbeitung bei großen Datenmengen [1].
 - * Höhere Anforderungen an die Rechenleistung, insbesondere bei der Nutzung von GPU [1].

Beide Tools bieten spezifische Vorteile: Flair zeichnet sich durch hohe Genauigkeit aus und ist besonders für anspruchsvolle NER-Aufgaben geeignet, während Spacy für schnelle und effiziente Verarbeitung moderner Texte ideal ist. Der Vergleich und die Bewertung dieser Tools werden im weiteren Verlauf detailliert behandelt.

2.3 Einsatz des Spacy-Tools

In diesem Abschnitt wird untersucht, wie Spacy für die Extraktion benannter Entitäten in historischen Texten eingesetzt wurde. Dabei werden die verfügbaren Modelle vorgestellt und die konkrete Implementierung des gewählten Modells für die Analyse unserer Daten beschrieben. Ziel ist es, die Effizienz und Genauigkeit von Spacy in diesem speziellen Anwendungsfall zu bewerten.

2.3.1 Verfügbare Modelle vom Spacy

Spacy bietet vier vortrainierte Modelle für die deutsche Sprache, die auf Nachrichten- und geschriebene Texte optimiert sind und die Erkennung von vier Entitätsklassen ermöglichen: **PER** (Personen), **ORG** (Organisationen), **LOC** (geografische Orte) und **MISC** (sonstige Kategorien) [4].

- **de_core_news_sm**: Ein kleines Modell, das grundlegende Funktionen wie Syntax- und Entitätserkennung bietet.
- **de_core_news_md**: Ein mittleres Modell mit zusätzlichen Wortvektoren zur Verarbeitung komplexerer Daten.
- **de_core_news_lg**: Ein großes Modell mit erweitertem Vokabular und detaillierten Wortvektoren.
- **de_dep_news_trf**: Ein Transformermodell für Abhängigkeitsparsing und Syntaxanalyse, basierend auf Nachrichten- und geschriebener Sprache.

Die vortrainierten Modelle von Spacy ermöglichen eine Erkennung von der Entitäten, ohne dass zusätzliches Training erforderlich ist.

2.3.2 Einsatz des Spacy-Modells

In diesem Projekt kam das Modell *de_core_news_lg* zum Einsatz, um benannte Entitäten wie Personen (**PER**), Orte (**LOC**) und Organisationen (**ORG**) aus den Texten zu extrahieren. Der Einsatz des Spacy-Modells erfolgt durch die Methode *classify_entities_std_Spacy*, die als Python-Quellecode in Listing 3 dargestellt ist und wie folgt arbeitet:

- **Modell laden**: Das Modell wird mit der Funktion *Spacy.load("de_core_news_lg")* geladen. Dieses Modell ist speziell auf die Verarbeitung der deutschen Sprache optimiert und enthält vortrainierte Pipelines für NER.

- **Metadaten extrahieren:** Vor der Textanalyse werden aus dem Dateinamen Metadaten wie *sender* und *id* der Datei extrahiert. Dies wird durch die Funktion *extract_meta_data(file_path)* erreicht, die den Dateinamen analysiert und relevante Informationen extrahiert.
- **Textverarbeitung:** Der Text der Eingabedatei wird in ein Spacy-Dokument mit *doc = nlp(text)* umgewandelt, das dann durchlaufen wird, um alle erkannten Entitäten zu extrahieren.
- **Extraktion der Entitäten:** Für jede erkannte Entität im Text werden die folgenden Informationen gesammelt:
 - **Text der Entität** (*ent.text*): Der erkannte Begriff oder Name.
 - **Kategorie der Entität** (*ent.label_*): Die Klassifikation der Entität, z. B. LOC (Ort), PER (Person) oder ORG (Organisation).
 - **Startposition im Text** (*ent.start_char*): Die Position, an der die Entität im Text beginnt.
 - **Endposition im Text** (*ent.end_char*): Die Position, an der die Entität im Text endet.
- **Speicherung der Ergebnisse:** Die extrahierten Entitäten werden in einem strukturierten JSON-Format in der Datei *std_Spacy_model_output.json* gespeichert.

Das Spacy-Modell wird direkt auf die Eingabedaten angewendet, ohne dass zusätzliches Training erforderlich war. Es bietet eine einfache Möglichkeit, benannte Entitäten wie Personen, Orte und Organisationen zu identifizieren und die Ergebnisse in einem strukturierten Format bereitzustellen.

2.4 Einsatz des Flair-Tools

In diesem Abschnitt wird untersucht, wie das Standardmodell von Flair zur Erkennung benannter Entitäten in historischen Texten eingesetzt wurde. Ziel ist es, die Funktionalität, die Genauigkeit sowie die Anwendbarkeit des Tools auf unsere spezifischen Daten zu bewerten und mit anderen Modellen zu vergleichen.

2.4.1 Verfügbare Modelle von Flair

Flair bietet eine Vielzahl vortrainierter Modelle für die Erkennung benannter Entitäten in verschiedenen Sprachen an. Diese Modelle erkennen ebenfalls dieselben Entitätsklassen bei Spacy. Für die deutsche Sprache stehen unter anderem folgende Modelle zur Verfügung [7]:

- **de-ner:** Ein Standardmodell für die Erkennung deutscher Entitäten, das Personen (PER), Orte (LOC), Organisationen (ORG) und weitere Kategorien identifiziert. Dieses Modell bietet ein ausgewogenes Verhältnis zwischen Genauigkeit und Geschwindigkeit.
- **de-ner-large:** Ein erweitertes Modell mit höherer Genauigkeit, das jedoch mehr Rechenressourcen benötigt. Es ist besonders geeignet für Anwendungen, bei denen die Präzision im Vordergrund steht.

Die Modelle von Flair basieren auf kontextualisierten String-Embeddings und wurden auf umfangreichen deutschen Textkorpora trainiert, um eine hohe Erkennungsrate für verschiedene Entitätstypen zu gewährleisten. [7]

2.4.2 Einsatz des Flair-Modells

Das vortrainierte Modell *de-ner* von Flair wurde aufgrund seines ausgewogenen Verhältnisses von Genauigkeit und Geschwindigkeit ausgewählt und verwendet, um benannte Entitäten aus deutschen Texten zu extrahieren. Der Einsatz des Flair-Modells erfolgt im Code durch die Methode *classify_entities_std_flair*, die als Python-Quellecode in Listing 3 dargestellt ist und wie folgt arbeitet:

- **Modell laden:** Das Modell *de-ner* wird mit der Funktion *Classifier.load("de-ner")* geladen. Dieses Modell ist speziell für die Verarbeitung der deutschen Sprache optimiert und enthält vortrainierte Pipelines für NER.
- **Metadaten extrahieren:** Vor der Textanalyse werden aus dem Dateinamen Metadaten wie *sender* und *id* der Datei extrahiert. Dies wird durch die Funktion *extract_meta_data(file_path)* erreicht, die den Dateinamen analysiert und relevante Informationen extrahiert.
- **Textverarbeitung:** Der Text der Eingabedatei wird in ein Flair-Sentence-Objekt durch *sentence = Sentence(text)* umgewandelt, das dann von der Modellinstanz mit *tagger.predict(sentence)* vorhergesagt wird.
- **Extraktion der Entitäten:** Für jede erkannte Entität im Text werden die folgenden Informationen gesammelt:
 - **Text der Entität** (*entity.text*): Der erkannte Begriff oder Name.
 - **Kategorie der Entität** (*entity.tag*): Die Klassifikation der Entität, z. B. LOC (Ort), PER (Person) oder ORG (Organisation).
 - **Startposition im Text** (*entity.start_position*): Die Position, an der die Entität im Text beginnt.

- **Endposition im Text** (*entity.end_position*): Die Position, an der die Entität im Text endet.
- **Speicherung der Ergebnisse:** Die extrahierten Entitäten werden in einem strukturierten JSON-Format in der Datei *std_flair_model_output.json* gespeichert.

Das Flair-Modell wird auch wie Spacy-Modell direkt auf die Eingabedaten angewendet, ohne dass zusätzliches Training erforderlich war. Es bietet eine Möglichkeit, benannte Entitäten wie Personen, Orte und Organisationen zu identifizieren und die Ergebnisse in einem strukturierten Format bereitzustellen.

2.5 Einsatz des weiter zu trainierenden Flair-Modells

Aufgrund der Ergebnisse der Standardmodelle von Flair und Spacy sind deutlich, dass diese für die spezifischen Anforderungen historischer Texte nicht optimal geeignet sind, wie in 2.6 vorgestellt wird. Daher wird im Rahmen dieses Projekts auch ein eigenes Flair-Modell mit spezifischen Daten und benutzerdefinierten Entitäten trainiert. Ziel ist es, die Erkennungsleistung für domänenspezifische Textdaten zu verbessern. Der Trainingsprozess umfasste mehrere Schritte, darunter die Aufbereitung der Trainingsdaten, das Modelltraining sowie die Evaluation und Anwendung des neu trainierten Modells auf bisher unbekannte Textdaten [6].

2.5.1 Vorverarbeitung der Trainingsdaten

Die Rohdaten bestanden aus annotierten Texten, in denen Entitäten mit Tags wie <LOC>, <PER> oder <ORG> markiert waren. Diese Texte wurden in das BIO-Format (Beginning, Inside, Outside) konvertiert, um sie für das Training des Flair-Modells vorzubereiten. Die Umsetzung dieser Konvertierung erfolgt durch die Funktion *parse_annotated_text*, die als Python-Quellcode in Listing 1 dargestellt ist. Das BIO-Format enthält folgende Kennzeichnungen:

- **B-(Beginning)**: Beginn einer Entität (z. B. B-LOC für den Anfang einer Ortsangabe).
- **I-(Inside)**: Innerhalb einer Entität (z. B. I-LOC für die Fortsetzung der Ortsangabe).
- **O-(Outside)**: Kein Teil einer Entität.

Ein Beispiel für das BIO-Format:

```

31    B-EVENT
1     I-EVENT
1999 I-EVENT
ist   0
eine  0
Stadt 0
.     0

```

Die in das BIO-Format konvertierten Daten werden anschließend mit der Funktion `split_train_file`, dargestellt in Listing 1, in drei Dateien unterteilt, die für das Modelltraining erforderlich sind:

- **train.txt**: Enthält die Trainingsdaten zur Modellanpassung.
- **dev.txt**: Dient der Validierung und Überwachung der Modellleistung während des Trainings.
- **test.txt**: Wird zur abschließenden Bewertung des trainierten Modells verwendet.

Diese drei Dateien bilden die Grundlage für das Training und die Evaluierung des Modells. Während *train.txt* zur Modelloptimierung genutzt wird, dient *dev.txt* zur laufenden Leistungsüberprüfung, und *test.txt* wird zur finalen Bewertung der Modellgenauigkeit herangezogen.

2.5.2 Training des Modells

Das Flair-Modell wird dann mit den vorbereiteten Daten trainiert. Hierbei werden *de* Word-Embeddings sowie Flair-Embeddings (*de-forward* und *de-backward*) kombiniert, um eine verbesserte Kontexterkenkung zu ermöglichen. Der Trainingsprozess erstreckt sich über 50 Epochen und wird unter Verwendung der Klasse `ModelTrainer` durchgeführt. Der gesamte Trainingsablauf ist in Listing 2 dargestellt und umfasst die folgenden Schritte:

- **Dateninitialisierung**: Die Trainings-, Test- und Validierungsdaten werden aus dem definierten Verzeichnis eingelesen und im CoNLL-Format verarbeitet. Dies geschieht durch die Funktion:

```

1     corpus: Corpus = ColumnCorpus(
2         self.train_data_folder,
3         column_format={0: "text", 1: "ner"},
4         train_file="train.txt",
5         test_file="test.txt",
6         dev_file="dev.txt")

```


- **Erstellung eines Label-Wörterbuchs:** Mit der Funktion `make_label_dictionary()` wird ein Wörterbuch aller in den Trainingsdaten vorkommenden Entitätsklassen erstellt:

```

1     label_type = "ner"
2     label_dict = corpus.make_label_dictionary(label_type=
        label_type)

```

- **Definition der Embeddings:** Die Trainingsdaten werden mit drei Arten von vortrainierten Embeddings ausgestattet: Word-Embeddings für die deutsche Sprache sowie zwei Flair-Embeddings für die Vorwärts- und Rückwärtsrichtung. Diese werden mittels der Funktion `StackedEmbeddings` kombiniert:

```

1     embedding_types_de = [
2         WordEmbeddings("de"),
3         FlairEmbeddings("de-forward"),
4         FlairEmbeddings("de-backward")
5     ]
6     stacked_embeddings = StackedEmbeddings(
        embedding_types_de)

```

- **Erstellung des Sequence-Taggers:** Das NER-Modell wird unter Verwendung der `SequenceTagger`-Klasse mit einer versteckten Größe von 256, CRF-Schicht und den zuvor definierten Embeddings konfiguriert:

```

1     tagger = SequenceTagger(
2         hidden_size=256,
3         embeddings=stacked_embeddings,
4         tag_dictionary=label_dict,
5         tag_type=label_type,
6         use_crf=True
7     )

```

- **Training des Modells:** Der Trainingsprozess wird mit einer Lernrate von 0.1, einer Mini-Batch-Größe von 16 und einer GPU-gestützten Embedding-Speicherung über 50 Epochen durchgeführt:

```

1     trainer = ModelTrainer(tagger, corpus)
2     trainer.train(
3         self.models_folder,
4         learning_rate=0.1,
5         mini_batch_size=16,
6         max_epochs=self.epochs,
7         embeddings_storage_mode="gpu"
8     )

```

- **Speichern des Modells** Nach Abschluss des Trainings werden zwei Versionen des trainierten Modells automatisch erzeugt und gespeichert, um eine flexible Nutzung je nach Anwendungsfall zu ermöglichen. Die gespeicherten Modellversionen sind:
 - **Final-Model:** Das Modell der letzten Trainingsepoche wird unter dem Namen *final-model.pt* gespeichert. Dieses Modell repräsentiert den Endzustand des Trainings und ermöglicht es, die zuletzt trainierte Version für zukünftige Anwendungen erneut zu laden und einzusetzen.
 - **Best-Model:** Zusätzlich wird das Modell mit der besten Leistung während des Trainingsprozesses unter dem Namen *best-model.pt* gespeichert. Dieses Modell wird anhand der besten validierten Metriken bestimmt und ist besonders für Anwendungen geeignet, bei denen maximale Genauigkeit erforderlich ist.

Die Speicherung beider Modelle stellt sicher, dass sowohl eine kontinuierliche Weiterentwicklung als auch eine optimale Modellnutzung je nach Anforderungen gewährleistet sind.

2.5.3 Anwendung des Modells

Die gespeicherten trainierten Modelle, *final-model.pt* und *best-model.pt*, werden verwendet, um benannte Entitäten aus neuen Texten zu extrahieren. Der gesamte Prozess der Modellanwendung ist im Python-Quellcode in Listing 3 durch die Methode *classify_entities* implementiert und umfasst folgende Schritte:

- **Laden des Modells:** Das gewünschte Modell wird mit der Methode *SequenceTagger.load(model_path)* geladen. Hier kann entweder das *final-model.pt* oder das leistungsstärkere *best-model.pt* gewählt werden.

```
1 self.model = SequenceTagger.load(self.model_path)
```

- **Lesen der Eingabedaten:** Alle Textdateien aus dem angegebenen Verzeichnis werden gelesen und die Zeichenkodierung wird mit der Funktion *detect()* automatisch erkannt.

```
1 txt_files = glob.glob(os.path.join(self.  
    input_data_folder, "*.txt"))  
2  
3 with open(file_path, "rb") as file:
```

```

4         raw_data = file.read()
5         encoding_info = detect(raw_data)
6         detected_encoding = encoding_info['encoding']
7
8         content = raw_data.decode(detected_encoding)

```

- **Extraktion von Metadaten:** Der Absender und die Dokumenten-ID werden aus dem Dateinamen extrahiert, um die Ergebnisse besser zuordnen zu können.

```

1         sender, file_id = self.extract_meta_data(file_path)

```

- **Entitätenerkennung:** Die gelesenen Texte werden mit der Methode *Sentence()* in Flair verarbeitet, und die erkannten Entitäten werden extrahiert.

```

1         sentence = Sentence(text, use_tokenizer=True)
2         self.model.predict(sentence)
3         for entity in sentence.get_spans("ner"):
4             entities.append([entity.text, entity.tag, entity.start_position, entity.end_position])

```

- **Speicherung der Ergebnisse** Die Ergebnisse werden in einem JSON-Format in der Datei *best-model.pt_output.json* oder *final-model.pt_output.json* je nach der geladenen Modellversion gespeichert, das die Entitäten mit Positionen und Kategorien darstellt.

```

1         with open(self.output_file, "w", encoding="utf-8") as
            json_file:
2             json.dump(all_results, json_file, indent=2,
                ensure_ascii=False)

```

2.5.4 Trainingsverlauf und Ergebnisse

Der Trainingsverlauf des benutzerdefinierten Flair-Modells wird mithilfe von unterschiedlichen Metriken überwacht. Diese folgenden Metriken dienen zur Bewertung der Modellleistung in Bezug auf die Erkennung und Klassifikation benannter Entitäten:

- **Präzision (Precision):** Der Anteil der korrekt erkannten Entitäten im Vergleich zu allen vom Modell vorhergesagten Entitäten.
- **Recall:** Der Anteil der korrekt erkannten Entitäten im Vergleich zu allen tatsächlich vorhandenen Entitäten.

- **F1-Score:** Das harmonische Mittel von Präzision und Recall, das eine ausgewogene Bewertung der Modellleistung ermöglicht.
- **Accuracy:** Der Anteil korrekt klassifizierter Wörter (Entität oder keine Entität) in Relation zu allen Wörtern im Text.

Während des Trainings wird die Leistung des Modells in jeder Epoche gemessen. Der folgende Verlauf zeigt die Entwicklung des F1-Score-Wertes über 50 Epochen:

```
Epoche 10: F1 = 40.78\%
Epoche 20: F1 = 53.33\%
Epoche 30: F1 = 57.14\%
Epoche 40: F1 = 59.77\%
Epoche 50: F1 = 63.53\%
```

Nach einem Trainingsverlauf sind beispielweise die folgenden Ergebnisse von (Precision, Recall, F1-Score und Accuracy) abhängig von den 5 Entitäten erzielt:

Results:

```
- F-score (micro) 0.5217
- F-score (macro) 0.4201
- Accuracy 0.3576
```

By class:

| | precision | recall | f1-score |
|--------|-----------|--------|----------|
| EVENT | 0.5745 | 0.7105 | 0.6353 |
| LOC | 0.4545 | 0.5000 | 0.4762 |
| PER | 0.5000 | 0.4783 | 0.4889 |
| ORG | 0.0000 | 0.0000 | 0.0000 |
| ANIMAL | 1.0000 | 0.3333 | 0.5000 |

Das trainierte Modell zeigt eine moderate Leistung mit einem finalen F1-Score von **63.53%**. Ereignisse (EVENT) wurden mit einem F1-Score von **63.53%** am besten erkannt, während Organisationen (ORG) überhaupt nicht identifiziert wurden (F1-Score: **0.00%**). Ortsangaben (LOC) und Personennamen (PER) erreichten mäßige Werte, was auf unzureichende Trainingsdaten hinweist. Tiere (ANIMAL) wurden mit hoher Präzision (**100.00%**), aber geringem Recall (**33.33%**) erkannt. Insgesamt zeigt das Modell Schwächen bei seltenen Entitäten und unausgewogene Trainingsdaten.

2.6 Evaluation der eingesetzten Modelle

In diesem Abschnitt werden die Leistungen der eingesetzten Modelle zur NER evaluiert. Ziel der Evaluation ist es, die Genauigkeit und Effizienz der Modelle bei der Erkennung und Klassifizierung benannter Entitäten in historischen Texten zu bewerten. Dazu wurden die 3 verschiedene Modelle auf identischen Evaluationsdaten getestet und anhand standardisierter Metriken wie Precision, Recall, F1-Score und Accuracy verglichen. Die gewonnenen Erkenntnisse helfen dabei, die Stärken und Schwächen der einzelnen Modelle zu identifizieren und Optimierungspotenziale für zukünftige Anwendungen aufzuzeigen.

2.6.1 Evaluationsmethodik und Vergleich der Modelle

Zur Evaluation der eingesetzten Modelle werden drei verschiedene Modelle auf denselben, nicht annotierten Texten anhand vordefinierter Metriken angewendet und bewertet. Der gesamte Evaluationsprozess ist als Python-Quellcode in Listing 4 implementiert und umfasst folgende Schritte:

- **Datenverarbeitung und Extraktion der Grundwahrheit (Ground Truth):** Annotierte Textdateien werden eingelesen, Metadaten extrahiert und Entitäten im JSON-Format gespeichert. Dies erfolgt durch die Methoden *extract_ground_truth*. Die Modelle werden auf dieselben Texte angewendet, und die extrahierten Entitäten werden mit der Grundwahrheit verglichen. Dabei wird die Methode *match_flexible* genutzt, um wahre und erkannte Entitäten zu vergleichen.
- **Berechnung der Metriken:** Die Leistung der Modelle wird anhand von Präzision, Recall, F1-Score und Accuracy evaluiert. Diese Berechnungen werden mit der Methode *calculate_metrics* durchgeführt.
- **Visualisierung der Ergebnisse:** Schließlich werden die Ergebnisse in Form von Diagrammen mit *plot_metrics* und *plot_individual_metrics* visualisiert.

Die Ergebnisse dieser Evaluation werden im JSON-Format gespeichert und erlauben eine direkte Gegenüberstellung der Modelle, wobei folgende Modelle evaluiert wurden:

- **Weiter trainiertes Flair-Modell (Final Model und Best Model):** Das Modell, das auf domänenspezifischen Texten trainiert ist, um die Erkennung von Entitäten zu optimieren.

- **Standardmodell von Flair** (Standard Flair Model): Ein vortrainiertes Modell, das ohne zusätzliche Anpassungen verwendet wird.
- **Standardmodell von Spacy** (Standard Spacy Model): Ein vortrainiertes Modell aus der Spacy-Bibliothek, das ebenfalls direkt eingesetzt wird.

Durch die Evaluierung dieser Modelle können wichtige Erkenntnisse über ihre Stärken und Schwächen bei der Erkennung und Klassifikation benannter Entitäten gewonnen werden.

2.6.2 Auswertung der Evaluation

Die Ergebnisse der Evaluation der drei Modelle sind in Abbildung 1, 2, 3 dargestellt. Die nachfolgenden Abbildungen zeigen eine detaillierte Beschreibung der Ergebnisse der Modelle hinsichtlich der Metriken *Precision*, *Recall*, *F1-Score* und *Accuracy*. Dabei werden verschiedene Kombinationen von Trainings- und Evaluationsdaten berücksichtigt, wie etwa Modelle, die mit Daten von Nachrichten eines einzelnen Absenders oder mit einer Mischung aus mehreren Absendern trainiert und evaluiert werden. Diese geben einen Überblick über die Leistungsunterschiede der Modelle unter den jeweiligen Bedingungen.

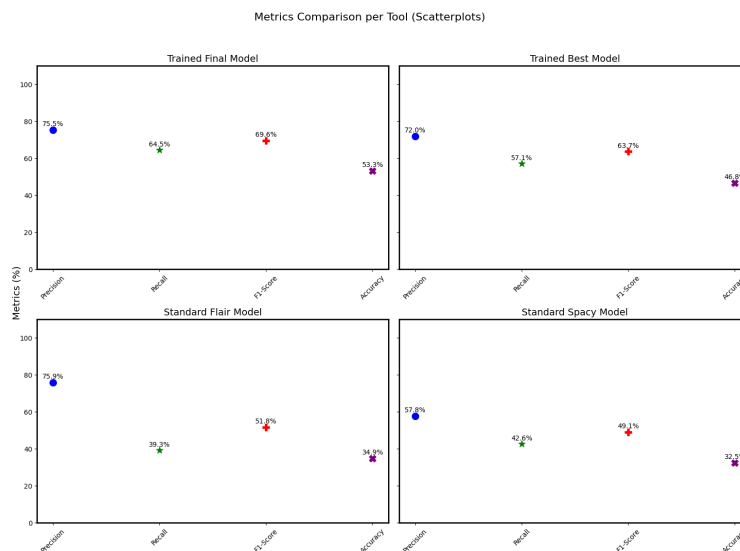


Abbildung 1: Die Trainingsdaten wurden aus Texten von einem bestimmten Absender erstellt, während die Evaluationsdaten von einem einzigen Absender stammen

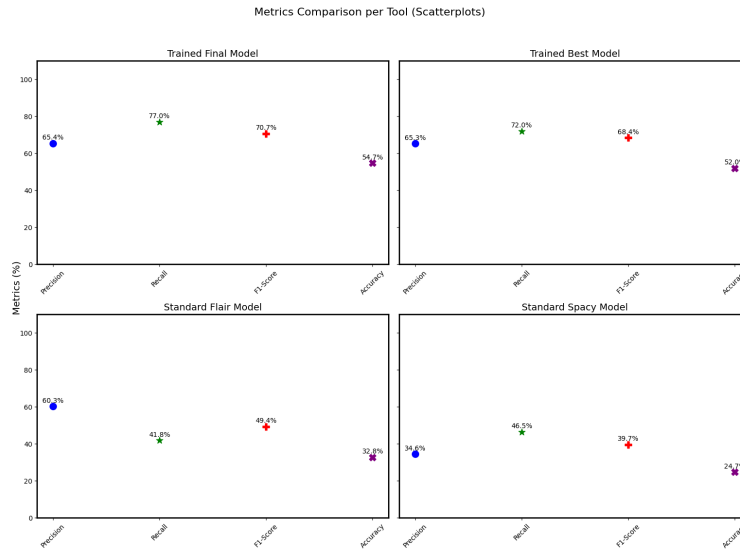


Abbildung 2: Die Trainingsdaten wurden aus einer Mischung von Daten mehrerer Absender erstellt, während die Evaluationsdaten von einem einzigen Absender stammen.

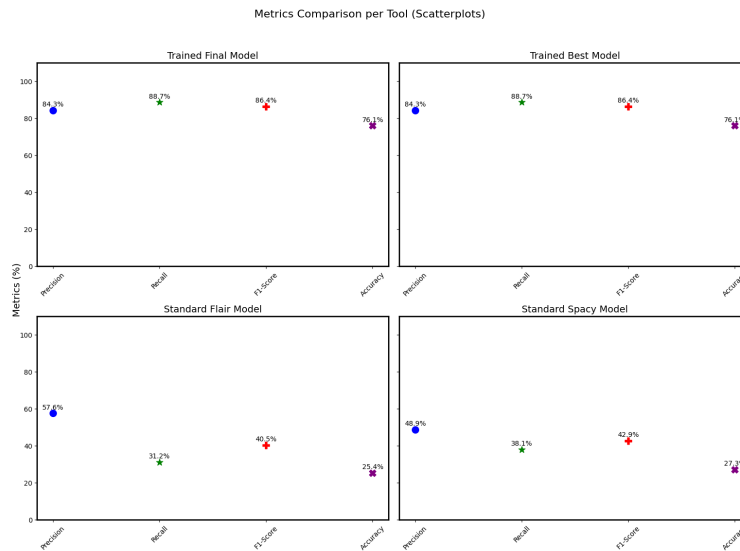


Abbildung 3: Sowohl die Trainingsdaten als auch die Evaluationsdaten wurden aus einer Mischung von Daten verschiedener Absender zusammengestellt.

2.6.3 Analyse der Ergebnisse und Schlussfolgerung

Die Analyse der Abbildungen zeigt, dass das weiter trainierte Flair-Modell (**Final Model** und **Best Model**) die beste Leistung unter den getesteten Modellen erzielt hat. Besonders der hohe Recall verdeutlicht, dass das Modell nahezu alle relevanten Entitäten erkennen konnte. Der F1-Score des trainierten Modells zeigt ein ausgewogenes Verhältnis zwischen Präzision und Recall, was auf die effektive Anpassung des Modells an die domänenspezifischen Trainingsdaten zurückzuführen ist. Die Abbildungen machen jedoch auch deutlich, dass die Art der Trainings- und Evaluationsdaten eine entscheidende Rolle für die Ergebnisse spielt. Flair-Modelle, die mit einer Mischung von Daten aus verschiedenen Absendern trainiert wurden, schnitten bei der Evaluation mit Daten eines einzelnen Absenders schlechter ab. Ebenso zeigen die Abbildungen, dass ein Modell, das mit Daten eines bestimmten Absenders trainiert und mit Daten eines anderen Absenders evaluiert wurde, in seiner Leistung limitiert ist. Diese Unterschiede verdeutlichen, dass eine Konsistenz zwischen Trainings- und Evaluationsdaten die Modellleistung stark beeinflusst. Im Vergleich dazu zeigen die Standardmodelle von Flair und Spacy in den Abbildungen geringere Leistungen. Das Standardmodell von Flair liefert zwar ähnliche Ergebnisse, erreicht jedoch weder die Präzision noch den Recall des weiter trainierten Modells. Das Standardmodell von Spacy zeigt eine insgesamt schwächere Leistung, was darauf hinweist, dass es weniger geeignet ist für die Anforderungen an domänenspezifische historische Texte. Zusammenfassend verdeutlichen die Abbildungen, dass die Anpassung durch eigenes Training auf spezifischen Daten essenziell ist, um eine hohe Erkennungsleistung zu erzielen. Gleichzeitig zeigt sich, dass die Art der Trainings- und Evaluationsdaten einen signifikanten Einfluss auf die Ergebnisse hat. Das weiter trainierte Flair-Modell bietet eine solide Grundlage für zukünftige Anwendungen und zeigt, dass benutzerdefinierte Trainingsansätze notwendig sind, um optimale Ergebnisse in spezialisierten Szenarien zu erreichen.

3 Zusammenfassung und Ausblick

In diesem Bericht wurden verschiedene Tools und Ansätze zur NER untersucht, evaluiert und miteinander verglichen. Ziel war es, die Erkennung benannter Entitäten in historischen, deutschsprachigen Texten zu optimieren. Die durchgeführten Analysen und Tests haben wertvolle Einblicke in die Leistungsfähigkeit der eingesetzten Modelle geliefert.

3.1 Zusammenfassung

Im Rahmen der Untersuchung wurden drei Ansätze eingesetzt und evaluiert: ein Standardmodell von Spacy, ein Standardmodell von Flair und ein weiter trainiertes Flair-Modell, das speziell an domänenspezifische Textdaten angepasst wurde. Die Ergebnisse der Evaluation zeigen, dass das weiter trainierte Flair-Modell die beste Gesamtleistung erzielte, insbesondere in Bezug auf Recall und F1-Score. Dies verdeutlicht die Relevanz einer domänenspezifischen Anpassung für die Verbesserung der NER-Leistung. Die wichtigsten Erkenntnisse lassen sich wie folgt zusammenfassen:

- Das weiter trainierte Flair-Modell erreichte mit einem Recall bis zum 88.0% und einem F1-Score bis zum 75.0% abhängig von den Trainings- und Evaluationsdatenart die beste Erkennungsleistung. Es konnte nahezu alle relevanten Entitäten erkennen und war somit für die spezifischen Anforderungen der Daten bestens geeignet.
- Die Standardmodelle von Spacy und Flair lieferten geringere Ergebnisse. Insbesondere das Spacy-Modell zeigte Schwächen bei der Präzision und Genauigkeit, da es ohne zusätzliche Anpassungen verwendet wurde.
- Die Ergebnisse unterstreichen, dass eine Anpassung des Modells durch Training mit spezifischen Daten entscheidend für die Leistungssteigerung ist.

3.2 Ausblick

Trotz der erzielten Ergebnisse bestehen weitere Möglichkeiten zur Optimierung der NER-Leistung:

- **Erweiterung der Trainingsdaten:** Die Verwendung eines größeren, ausgewogeneren Datensatzes könnte die Modellleistung insbesondere bei unterrepräsentierten Klassen wie **ORG** oder **ANIMAL** verbessern.
- **Anwendung weiterer Modelle:** Die Evaluation zusätzlicher NER-Ansätze, wie beispielsweise (Hugging Face), könnte weitere Erkenntnisse über die Leistungsfähigkeit moderner Ansätze liefern.
- **Integration in Workflows:** Die Ergebnisse könnten in einer praktischen Anwendung integriert werden, z. B. in der historischen Textanalyse oder bei der digitalen Archivierung.

Zusammenfassend lässt sich feststellen, dass die gezielte Anpassung von NER-Modellen an spezifische Anforderungen eine wesentliche Rolle für den Erfolg spielt. Die Erkenntnisse aus diesem Projekt bieten eine solide Grundlage für weiterführende Forschungen und die Entwicklung verbesserter NER-Methoden.

4 Python-Codes zum Einsatz der Tools

4.1 Trainingsdatenvorverarbeitung

```
1 import os
2 import random
3 from charset_normalizer import detect
4 from transformers import DataProcessor
5
6 class DataProcessor:
7
8     # Initialize file paths for training and base data.
9     def __init__(self):
10         self.output_file = "train_data/train.txt"
11         self.base_data_folder = "../base_data"
12         self.train_data_folder = "train_data"
13
14     # Read file content with automatic encoding detection.
15     def read(self, file_path):
16         try:
17             with open(file_path, "rb") as file:
18                 raw_data = file.read()
19                 encoding_info = detect(raw_data)
20                 detected_encoding = encoding_info['encoding']
21                 content = raw_data.decode(detected_encoding)
22                 print("File content loaded successfully.")
23                 return content
24         except FileNotFoundError:
25             print(f"Error: The file '{file_path}' was not found.")
26             return ""
27         except Exception as e:
28             print(f"Error: {e}")
29             return ""
30
31     # Convert annotated text to CoNLL format with BIO tagging
32     def parse_annotated_text(self, input_text):
33         sentences = []
34         current_sentence = []
35         current_tag = None
36         inside_tag = False
37         word_buffer = []
38         i = 0
39         while i < len(input_text):
40             char = input_text[i]
41             if char == "<":
42                 closing_index = input_text.find(">", i)
43                 tag_content = input_text[i + 1: closing_index]
```

```

]
44     if tag_content.startswith("/"):
45         if word_buffer:
46             word = "".join(word_buffer).strip()
47             if word:
48                 if inside_tag:
49                     current_sentence.append((word
50                                             , f"I-{{current_tag}}"))
51                 else:
52                     current_sentence.append((word
53                                             , "0"))
54                 word_buffer = []
55                 inside_tag = False
56                 current_tag = None
57             else:
58                 current_tag = tag_content
59                 inside_tag = True
60                 i = closing_index
61     elif char in {" ", "\n", ".", "!", "?"}:
62         if word_buffer:
63             word = "".join(word_buffer).strip()
64             if word:
65                 tag = f"B-{{current_tag}}" if not any(w
66                                                         [1].startswith("B-") for w in
67                                                         current_sentence) else f"I-{{
68                                                         current_tag}}"
69                 current_sentence.append((word, tag))
70                 word_buffer = []
71             if char in {".", "!", "?"}:
72                 current_sentence.append((char, "0"))
73                 sentences.append(current_sentence)
74                 current_sentence = []
75             else:
76                 word_buffer.append(char)
77                 i += 1
78     if word_buffer:
79         word = "".join(word_buffer).strip()
80         if inside_tag:
81             tag = f"B-{{current_tag}}" if not any(w[1].
82                                                         startswith("B-") for w in current_sentence
83                                                         ) else f"I-{{current_tag}}"
84             current_sentence.append((word, tag))
85         else:
86             current_sentence.append((word, "0"))
87     if current_sentence:
88         sentences.append(current_sentence)
89     return sentences
90
91 # Write the processed sentences in CoNLL format to an

```

```

    output_file.
85 def write_to_conll(self, sentences, output_file):
86     with open(output_file, "a", encoding="utf-8") as f:
87         for sentence in sentences:
88             for word, tag in sentence:
89                 f.write(f"{word} {tag}\n")
90             f.write("\n")
91
92 # Read, parse, and save text data in CoNLL format.
93 def process_text_file(self, file_path):
94     text = self.read(file_path)
95     if text:
96         sentences = self.parse_annotated_text(text)
97         self.write_to_conll(sentences, self.output_file)
98         print(f"CoNLL-formatted data saved to '{self.
          output_file}'.")
99
100 # Split the data into train, test, and dev sets.
101 def split_train_file(self, test_ratio=0.20, dev_ratio
    =0.1, seed=42):
102     with open(self.output_file, "r", encoding="utf-8") as
        f:
103         sentences = f.read().strip().split("\n\n")
104         random.seed(seed)
105         random.shuffle(sentences)
106         test_size = int(len(sentences) * test_ratio)
107         dev_size = int(len(sentences) * dev_ratio)
108         test_sentences = sentences[:test_size]
109         dev_sentences = sentences[test_size:test_size +
            dev_size]
110         train_sentences = sentences[test_size + dev_size:]
111         with open(f"{self.train_data_folder}/train.txt", "w",
            encoding="utf-8") as f:
112             f.write("\n\n".join(train_sentences))
113         with open(f"{self.train_data_folder}/test.txt", "w",
            encoding="utf-8") as f:
114             f.write("\n\n".join(test_sentences))
115         with open(f"{self.train_data_folder}/dev.txt", "w",
            encoding="utf-8") as f:
116             f.write("\n\n".join(dev_sentences))
117         print("Data file was split")
118
119 # Delete the existing training data file.
120 def delete_train_data(self):
121     with open(self.output_file, "w", encoding="utf-8") as
        file:
122         pass
123
124 # Process all text files from the base data folder.

```

```

125     def run(self):
126         for file_name in os.listdir(self.base_data_folder):
127             file_path = os.path.join(self.base_data_folder,
128                                     file_name)
129             if os.path.isfile(file_path):
130                 self.process_text_file(file_path)
131
132     # Initialize the DataProcessor class
133     data_processor = DataProcessor()
134
135     # Deletes the content of the training data file if it exists
136     data_processor.delete_train_data()
137
138     # Runs the entire data processing pipeline (reads and
139     # processes input files)
140     data_processor.run()
141
142     # Splits the processed data into training, test, and
143     # validation datasets
144     data_processor.split_train_file()
145
146     #-----

```

Listing 1: Datenverarbeitungsskript für NER-Modelle

4.2 NER-Modelltraining

```
1 from flair.data import Sentence, Corpus
2 from flair.datasets import ColumnCorpus
3 from flair.embeddings import WordEmbeddings, FlairEmbeddings,
  StackedEmbeddings
4 from flair.models import SequenceTagger
5 from flair.trainers import ModelTrainer
6
7 class Model:
8
9     # Initialize the model parameters such as data folders
10    # and number of epochs
11    def __init__(self, train_data_folder="train_data",
12                 test_data_folder="",
13                 models_folder="models/flair_model",
14                 model_path="models/flair_model/final-model.
15                 pt", epochs=50):
16        self.train_data_folder = train_data_folder
17        self.test_data_folder = test_data_folder
18        self.models_folder = models_folder
19        self.model_path = model_path
20        self.epochs = epochs
21
22    # Train the NER model using Flair embeddings and save it
23    def train(self):
24        columns = {0: "text", 1: "ner"}
25        corpus: Corpus = ColumnCorpus(
26            self.train_data_folder, column_format=columns,
27            train_file="train.txt", test_file="test.txt",
28            dev_file="dev.txt",
29        )
30        label_type = "ner"
31        label_dict = corpus.make_label_dictionary(label_type=
32            label_type)
33
34        embedding_types_de = [
35            WordEmbeddings("de"),
36            FlairEmbeddings("de-forward"),
37            FlairEmbeddings("de-backward"),
38        ]
39        stacked_embeddings = StackedEmbeddings(
40            embedding_types_de)
41
42        tagger = SequenceTagger(
43            hidden_size=256, embeddings=stacked_embeddings,
44            tag_dictionary=label_dict, tag_type=label_type,
45            use_crfs=True,
46        )
```

```

40
41     trainer = ModelTrainer(tagger, corpus)
42     trainer.train(
43         self.models_folder, learning_rate=0.1,
44         mini_batch_size=16, max_epochs=self.epochs,
45         embeddings_storage_mode="gpu",
46     )
47
48     # Evaluate the trained model using test data and display
49     # performance metrics
50     def evaluate(self):
51         columns = {0: "text", 1: "ner"}
52         corpus = ColumnCorpus(self.train_data_folder,
53                               column_format=columns, test_file="test.txt")
54         tagger = SequenceTagger.load(self.model_path)
55         result = tagger.evaluate(corpus.test, gold_label_type
56                                 ="ner")
57         print(result.detailed_results)
58
59     # Test the trained model with a given input text and
60     # print detected entities
61     def test(self, input_text):
62         tagger = SequenceTagger.load(self.model_path)
63         sentence = Sentence(input_text, use_tokenizer=True)
64         tagger.predict(sentence)
65         print(sentence.get_spans("ner"))
66         for entity in sentence.get_spans("ner"):
67             print(f"Text: {entity.text}, Type: {entity.tag},
68                   Confidence: {entity.score:.2f}")
69
70     # Model usage
71     model = Model()
72
73     # Train the model
74     model.train()
75
76     # -----

```

Listing 2: NER-Modelltraining mit Flair

4.3 NER-Klassifikation

```
1 import glob
2 import json
3 import os
4 import re
5 import spacy
6 from flair.data import Sentence
7 from flair.models import SequenceTagger
8 from charset_normalizer import detect
9 from flair.nn import Classifier
10
11 class NERClassifier:
12
13     # Initialize the classifier with model paths and
14     # configuration options
15     def __init__(self,
16                 std_spacy=False,
17                 std_flair=False,
18                 json_data_folder="json_data",
19                 input_data_folder="../input_data",
20                 model_path="models/flair_model/",
21                 model_type="final-model.pt"):
22
23         self.json_data_folder = json_data_folder
24         self.input_data_folder = input_data_folder
25         self.std_spacy = False
26         self.std_flair = False
27
28         if std_spacy:
29             self.output_file = os.path.join(self.
30                 json_data_folder, "std_spacy_model_output.json
31                 ")
32             self.std_spacy = True
33
34         elif std_flair:
35             self.output_file = os.path.join(self.
36                 json_data_folder, "std_flair_model_output.json
37                 ")
38             self.std_flair = True
39
40         else:
41             self.output_file = os.path.join(self.
42                 json_data_folder, model_type + "_output.json")
43             self.model_type = model_type
44             self.model_path = model_path + self.model_type
45             self.model = SequenceTagger.load(self.model_path)
46
47     # Extract metadata (sender and file ID) from file name
```

```

42     @staticmethod
43     def extract_meta_data(file_path):
44         file_name = os.path.basename(file_path)
45         match = re.match(r"([^\s]+\s+([^\s]+)", file_name)
46         sender = match.group(1) if match else "unknown"
47         file_id = match.group(2).split(".")[0] if match else
            "unknown"
48
49         return sender.replace("korrigiert", ""), file_id
50
51     # Read all text files from the input folder
52     def read_files(self):
53         txt_files = glob.glob(os.path.join(self.
            input_data_folder, "*.txt"))
54         files_content = {}
55
56         for file_path in txt_files:
57             with open(file_path, "rb") as file:
58                 raw_data = file.read()
59                 encoding_info = detect(raw_data)
60                 detected_encoding = encoding_info['encoding']
61                 try:
62                     content = raw_data.decode(
                        detected_encoding)
63                     files_content[file_path] = content
64                 except Exception as e:
65                     print(f"Error reading file {file_path}: {
                        e}")
66                     continue
67
68         return files_content
69
70     # Classify named entities using the trained Flair model
71     def classify_entities(self, file_path, text):
72         sender, file_id = self.extract_meta_data(file_path)
73
74         sentence = Sentence(text, use_tokenizer=True)
75         self.model.predict(sentence)
76
77         entities = []
78         for entity in sentence.get_spans("ner"):
79             print("Entity: ", entity)
80             entities.append([
81                 entity.text,
82                 entity.tag,
83                 entity.start_position,
84                 entity.end_position
85             ])
86

```

```

87         return {
88             "id": file_id,
89             "sender": sender,
90             "text": text,
91             "entities": entities
92         }
93
94     # Classify named entities using the standard Spacy model
95     def classify_entities_std_spacy(self, file_path, text):
96         sender, file_id = self.extract_meta_data(file_path)
97
98         nlp = spacy.load("de_core_news_lg")
99         doc = nlp(text)
100         entities = []
101         for ent in doc.ents:
102             label = "EVENT" if ent.label_ == "DATE" else ent.
103                 label_
104             entities.append([ent.text, label, ent.start_char,
105                             ent.end_char])
106
107         return {
108             "id": file_id,
109             "sender": sender.replace("korrigiert", ""),
110             "text": text,
111             "entities": entities
112         }
113
114     # Classify named entities using the standard Flair model
115     def classify_entities_std_flair(self, file_path, text):
116         sender, file_id = self.extract_meta_data(file_path)
117
118         model = "de-ner"
119         tagger = Classifier.load(model)
120         sentence = Sentence(text)
121         tagger.predict(sentence)
122
123         entities = []
124         for entity in sentence.get_spans('ner'):
125             label = "EVENT" if entity.tag == "DATE" else
126                 entity.tag
127             entities.append([entity.text, label, entity.
128                             start_position, entity.end_position])
129
130         return {
131             "id": file_id,
132             "sender": sender.replace("korrigiert", ""),
133             "text": text,
134             "entities": entities
135         }

```

```

132
133     # Process all text files and classify entities using the
        chosen model
134     def process_all_files(self):
135         all_results = {"results": []}
136
137         files_content = self.read_files()
138
139         for file_path, text in files_content.items():
140             if self.std_spacy:
141                 result = self.classify_entities_std_spacy(
                    file_path, text)
142             elif self.std_flair:
143                 result = self.classify_entities_std_flair(
                    file_path, text)
144             else:
145                 result = self.classify_entities(file_path,
                    text)
146             all_results["results"].append(result)
147
148             os.makedirs(self.json_data_folder, exist_ok=True)
149             with open(self.output_file, "w", encoding="utf-8") as
                json_file:
150                 json.dump(all_results, json_file, indent=2,
                    ensure_ascii=False)
151
152             print(f"All results were saved in {self.output_file}.
                ")
153
154
155     # Run classification with the trained Flair model
156     ner = NERClassifier()
157     ner.process_all_files()
158
159     # Run classification with the best-trained Flair model
160     ner = NERClassifier(model_type="best-model.pt")
161     ner.process_all_files()
162
163     # Run classification using the standard Spacy model
164     ner = NERClassifier(std_spacy=True)
165     ner.process_all_files()
166
167     # Run classification using the standard Flair model
168     ner = NERClassifier(std_flair=True)
169     ner.process_all_files()
170
171     #-----

```

Listing 3: NER-Klassifikation mit Flair und Spacy

4.4 NER Evaluation

```
1 import os
2 import json
3 import re
4 from datetime import datetime
5 from difflib import SequenceMatcher
6 from charset_normalizer import detect
7 from matplotlib import pyplot as plt
8
9
10 class Evaluator:
11
12     # Reads the content of a file with automatic encoding
13     # detection.
14     def read_file_with_encoding(self, file_path):
15         with open(file_path, "rb") as f:
16             raw_data = f.read()
17             encoding_info = detect(raw_data)
18             detected_encoding = encoding_info['encoding']
19
20             if detected_encoding:
21                 try:
22                     return raw_data.decode(detected_encoding)
23                 except Exception as e:
24                     print(f"Error decoding file {file_path}: {e}")
25                     return ""
26             else:
27                 print(f"Could not detect encoding for file {file_path}.")
28                 return ""
29
30     # Extracts sender name and document ID from the file name
31     @staticmethod
32     def extract_meta_data(file_path):
33         file_name = os.path.basename(file_path)
34         match = re.match(r"([^\s]+\s+)([^\s]+)", file_name)
35         sender = match.group(1) if match else "unknown"
36         file_id = match.group(2).split(".")[0] if match else "unknown"
37         return sender.replace("korrigiert", ""), file_id
38
39     # Extracts named entities from the given annotated text.
40     @staticmethod
41     def extract_entities_from_text(text):
42         entities = []
43         offset = 0
44         pattern = r"<(P<label>[^\<]+)>(P<text>[^\<]+)</\1>"
```

```

44         for match in re.finditer(pattern, text):
45             label = match.group("label")
46             entity_text = match.group("text")
47             start = match.start("text") - offset
48             end = match.end("text") - offset
49             entities.append([entity_text, label, start, end])
50         return entities
51
52     # Processes files to extract ground truth entities and
53     # save them in JSON format.
54     def extract_ground_truth(self):
55         output_file = "json_data/ground_truth.json"
56         results = []
57
58         file_paths = []
59         for root, dirs, files in os.walk("evaluate_data"):
60             for file in files:
61                 file_paths.append(os.path.join(root, file))
62
63         for file_path in file_paths:
64             if not os.path.exists(file_path):
65                 print(f"File not found: {file_path}")
66                 continue
67
68             sender, doc_id = self.extract_meta_data(file_path)
69             content = self.read_file_with_encoding(file_path)
70             entities = self.extract_entities_from_text(
71                 content)
72
73             results.append({
74                 "id": doc_id,
75                 "sender": sender,
76                 "entities": entities
77             })
78
79         with open(output_file, "w", encoding="utf-8") as
80             json_out:
81                 json.dump({"results": results}, json_out, indent
82                             =4, ensure_ascii=False)
83         print(f"Processed data saved to {output_file}")
84
85     # Loads JSON data from a given file.
86     @staticmethod
87     def load_json(file_path):
88         with open(file_path, "r", encoding="utf-8") as f:
89             return json.load(f)
90
91     # Extracts entity tuples (text, label) from JSON results.

```

```

88     @staticmethod
89     def extract_entities(results):
90         entities = []
91         for result in results:
92             for entity in result["entities"]:
93                 word = entity[0]
94                 label = entity[1]
95                 entities.append((word, label))
96         return set(entities)
97
98     # Checks if two words are similar based on a similarity
99     # threshold.
100    @staticmethod
101    def is_similar(word1, word2, threshold=0.8):
102        similarity = SequenceMatcher(None, word1, word2).
103            ratio()
104        return similarity >= threshold
105
106    # Matches predicted entities with ground truth using
107    # flexible matching criteria.
108    def match_flexible(self, true_entities,
109        predicted_entities, similarity_threshold=1.0):
110        true_positives = set()
111        false_positives = set(predicted_entities)
112        false_negatives = set()
113
114        for true_word, true_label in true_entities:
115            matched = False
116            for pred_word, pred_label in predicted_entities:
117                if true_label == pred_label and (
118                    true_word in pred_word or pred_word in
119                    true_word or
120                    self.is_similar(true_word, pred_word,
121                        similarity_threshold)
122                ):
123                    true_positives.add((pred_word, pred_label))
124                    matched = True
125                    false_positives.discard((pred_word,
126                        pred_label))
127                if not matched:
128                    false_negatives.add((true_word, true_label))
129
130        return {
131            "true_positives": true_positives,
132            "false_positives": false_positives,
133            "false_negatives": false_negatives,
134        }

```

```

129     # Calculates precision, recall, F1-score, and accuracy
130     for evaluation.
131     def calculate_metrics(self, true_entities,
132                           predicted_entities):
133         results = self.match_flexible(true_entities,
134                                       predicted_entities)
135         tp = len(results["true_positives"])
136         fp = len(results["false_positives"])
137         fn = len(results["false_negatives"])
138
139         precision = tp / (tp + fp) if (tp + fp) > 0 else 0.0
140         recall = tp / (tp + fn) if (tp + fn) > 0 else 0.0
141         f1 = 2 * (precision * recall) / (precision + recall)
142             if (precision + recall) > 0 else 0.0
143         total = tp + fp + fn
144         accuracy = tp / total if total > 0 else 0.0
145
146         return precision, recall, f1, accuracy
147
148     # Displays the evaluation results of different NER models
149     def display_results(self):
150         self.extract_ground_truth()
151
152         # Extract ground truth entities from annotated text
153         files
154         ground_truth_path = "json_data/ground_truth.json"
155         tool_final_path = "json_data/final-model.pt_output.
156             json"
157         tool_best_path = "json_data/best-model.pt_output.json
158             "
159         tool_std_path = "json_data/std_flair_model_output.
160             json"
161         tool_spacy_path = "json_data/std_spacy_model_output.
162             json"
163
164         # Load the ground truth and model outputs from JSON
165         files
166         ground_truth = self.load_json(ground_truth_path)["
167             results"]
168         tool_final_results = self.load_json(tool_final_path)[
169             "results"]
170         tool_best_results = self.load_json(tool_best_path)["
171             results"]
172         tool_std_results = self.load_json(tool_std_path)["
173             results"]
174         tool_spacy_results = self.load_json(tool_spacy_path)[
175             "results"]
176
177         # Extract entities from ground truth and model

```



```

163         results
164         true_entities = self.extract_entities(ground_truth)
165         tool_final_entities = self.extract_entities(
166             tool_final_results)
167         tool_best_entities = self.extract_entities(
168             tool_best_results)
169         tool_std_entities = self.extract_entities(
170             tool_std_results)
171         tool_spacy_entities = self.extract_entities(
172             tool_spacy_results)
173
174         # Calculate evaluation metrics (Precision, Recall, F1
175         # -Score, Accuracy) for each model
176         metrics = {
177             "Trained Final Model": self.calculate_metrics(
178                 true_entities, tool_final_entities),
179             "Trained Best Model": self.calculate_metrics(
180                 true_entities, tool_best_entities),
181             "Standard Flair Model": self.calculate_metrics(
182                 true_entities, tool_std_entities),
183             "Standard Spacy Model": self.calculate_metrics(
184                 true_entities, tool_spacy_entities),
185         }
186
187         # Print evaluation results for each model
188         for tool, (precision, recall, f1, accuracy) in
189             metrics.items():
190             print(f"{tool} - Precision: {precision:.2f},
191                   Recall: {recall:.2f}, F1-score: {f1:.2f},
192                   Accuracy: {accuracy:.2f}")
193
194         # plots evaluation results
195         self.plot_metrics(metrics)
196         self.plot_individual_metrics(metrics)
197
198         # Plots evaluation metrics for different models.
199         @staticmethod
200         def plot_metrics(metrics):
201             tools = list(metrics.keys())
202             precision = [metrics[tool][0] * 100 for tool in tools]
203             recall = [metrics[tool][1] * 100 for tool in tools]
204             f1_score = [metrics[tool][2] * 100 for tool in tools]
205             accuracy = [metrics[tool][3] * 100 for tool in tools]
206             x = range(len(tools))
207
208             fig, ax = plt.subplots(figsize=(10, 6))
209
210             ax.scatter(x, precision, label='Precision', marker='o')

```

```

    ', color='blue', linewidth=2)
198 ax.scatter(x, recall, label='Recall', marker='*',
    color='green', linewidth=2)
199 ax.scatter(x, f1_score, label='F1-Score', marker='P',
    color='red', linewidth=2)
200 ax.scatter(x, accuracy, label='Accuracy', marker='X',
    color='purple', linewidth=2)
201
202 for i, (p, r, f, a) in enumerate(zip(precision,
    recall, f1_score, accuracy)):
203     ax.text(i+0.1, p, f"{p:.1f}%", fontsize=10, ha='
    center', va='bottom', color='blue')
204     ax.text(i+0.1, r, f"{r:.1f}%", fontsize=10, ha='
    center', va='bottom', color='green')
205     ax.text(i+0.1, f, f"{f:.1f}%", fontsize=10, ha='
    center', va='bottom', color='red')
206     ax.text(i+0.1, a, f"{a:.1f}%", fontsize=10, ha='
    center', va='bottom', color='purple')
207
208 ax.set_xticklabels(labels=tools, rotation=45,
    fontsize=12)
209 ax.set_xticks(x)
210
211
212 ax.set_xlabel("NER Tools", fontsize=14)
213 ax.set_ylabel("Metrics (%)", fontsize=14)
214 ax.set_title("Comparison of NER Tools - Precision,
    Recall, F1-Score, and Accuracy", fontsize=16)
215
216 ax.legend(fontsize=12)
217
218 plt.tight_layout()
219 plt.show()
220
221 # Plots evaluation metrics for different models
    individual
222 @staticmethod
223 def plot_individual_metrics(metrics):
224     tools = list(metrics.keys())
225     num_tools = len(tools)
226
227     fig, axes = plt.subplots(2, 2, figsize=(16, 12),
        sharey=True)
228
229     colors = ['blue', 'green', 'red', 'purple']
230     markers = ['o', '*', 'P', 'X']
231     labels = ['Precision', 'Recall', 'F1-Score', '
        Accuracy']
232

```

```

233     for idx, tool in enumerate(tools):
234         row, col = divmod(idx, 2) # Bestimme Zeile und
           Spalte (2x2-Layout)
235         ax = axes[row, col] # Hole den aktuellen
           Achsenbereich
236
237         precision = metrics[tool][0] * 100
238         recall = metrics[tool][1] * 100
239         f1_score = metrics[tool][2] * 100
240         accuracy = metrics[tool][3] * 100
241
242         values = [precision, recall, f1_score, accuracy]
243         x = range(len(values)) # Positionen auf der X-
           Achse
244
245         for j, (value, color, marker) in enumerate(zip(
           values, colors, markers)):
246             ax.scatter(j, value, color=color, marker=
               marker, label=labels[j], s=100)
247
248             # Werte oberhalb der Punkte anzeigen
249             ax.text(j, value + 2, f"{value:.1f}%", ha='
               center', fontsize=10)
250
251         ax.set_title(f"{tool}", fontsize=14)
252         ax.set_xticklabels(labels, rotation=45, fontsize
           =10)
253         ax.set_xticks(range(len(labels)))
254         ax.set_ylim(0, 110)
255
256         ax.spines['top'].set_linewidth(2)
257         ax.spines['right'].set_linewidth(2)
258         ax.spines['bottom'].set_linewidth(2)
259         ax.spines['left'].set_linewidth(2)
260
261     for idx in range(len(tools), 4): # Wenn weniger als
           4 Tools vorhanden sind
262         row, col = divmod(idx, 2)
263         axes[row, col].axis('off') # Nicht genutzte
           Achsen deaktivieren
264
265     fig.text(0.04, 0.5, "Metrics (%)", va='center',
           rotation='vertical', fontsize=14)
266
267     fig.suptitle("Metrics Comparison per Tool (
           Scatterplots)", fontsize=16)
268
269     plt.tight_layout(rect=[0.04, 0.04, 1, 0.95])
270     plt.savefig(f"plots/evaluate_result_{datetime.now()}.

```

```

271         strftime("%Y_%m_%d%H_%M_%S").png")
272     plt.show()
273
274
275
276
277 # Run evaluation
278 evaluator = Evaluator()
279 evaluator.display_results()
280
281 # -----

```

Listing 4: NER Evaluation Pipeline

Literatur

- [1] Alan Akbik, Duncan Blythe und Roland Vollgraf. „Flair: An Easy-to-Use Framework for State-of-the-Art NLP“. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*. 2019, S. 54–59. URL: <https://aclanthology.org/N19-4010/>.
- [2] Papers with Code. *Named Entity Recognition (NER) - Papers with Code*. Zugriff: 07.11.2024. 2024. URL: <https://paperswithcode.com/task/named-entity-recognition-ner>.
- [3] Matthew Honnibal u. a. *spaCy: Industrial-strength Natural Language Processing in Python*. 2020. DOI: 10.5281/zenodo.1212303. URL: <https://github.com/explosion/spaCy>.
- [4] Matthew Honnibal u. a. *spaCy: Industrial-strength Natural Language Processing in Python*. 2020. URL: <https://spacy.io/models/de>.
- [5] Srinivasa Rao Kundeti u. a. „Clinical named entity recognition: Challenges and opportunities“. In: *2016 IEEE International Conference on Big Data (Big Data)*. 2016, S. 1937–1945. DOI: 10.1109/BigData.2016.7840814.
- [6] Flair NLP. *How Model Training Works*. Accessed: 2023-12-10. 2023. URL: <https://flairnlp.github.io/flair/master/tutorial/tutorial-training/how-model-training-works.html>.
- [7] Flair NLP. *Tagging entities*. 2023. URL: <https://flairnlp.github.io/docs/tutorial-basics/tagging-entities>.

- [8] Mudasar Ghafoor und Salman Naseer. *Named Entity Recognition (NER) in NLP: Techniques, Tools, Accuracy and Performance*. Zugriff: 16.11.2024. 2022.
- [9] *WossiDia*. Zugriff: 16.11.2024. URL: <https://www.wossidia.de/>.