

Projektrealisierung

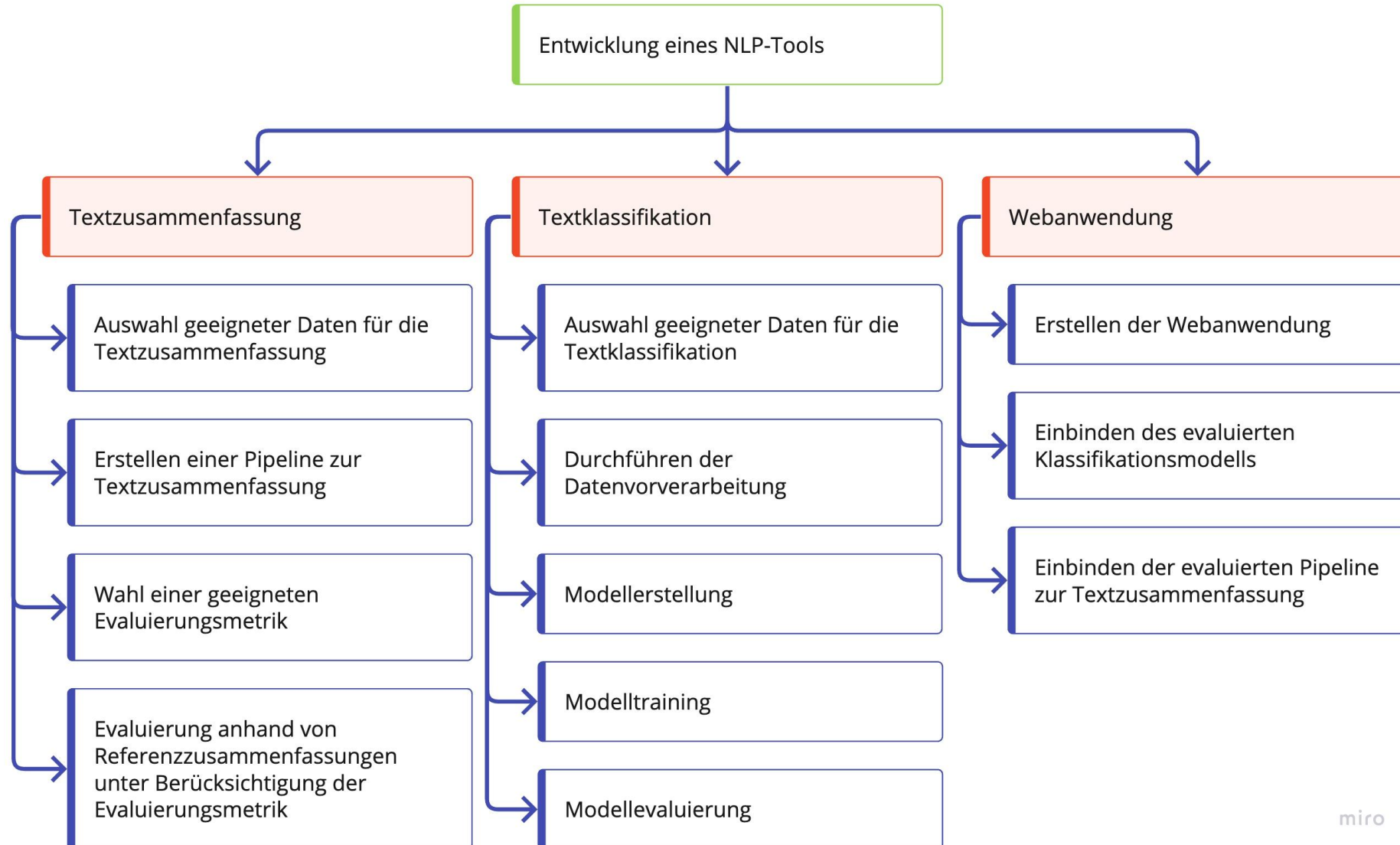
Zwischenstand – 20.06.23

Von Lukas Bonn, Amina Uicker-Darwish,
Jan Rüdert, Aymane Bouguern

Agenda

- ① Übersicht: Projektstrukturplan und Arbeitspakete
- ② Entwicklungsstand der Anforderungen
 - Klassifikation
 - Zusammenfassung
 - Frontend
- ③ Herausforderungen
- ④ Ausblick

Recap - Projektstrukturplan



Aktualisierung - Projektstrukturplan



Entwicklungsstand - Zusammenfassung

Aspekte aus dem PSP

Erstellung des Datensatzes	<ul style="list-style-type: none">• Texte aus 5 verschiedenen Kategorien, Texlänge je 500-800 Wörtern	✓
Referenzzusammenfassungen erstellen	<ul style="list-style-type: none">• Anhand von LSA, Textrank, Lexrank (Sumy Library), Unterschiedliche Kompressionsraten (random)	✓
Datenvorverarbeitung	<ul style="list-style-type: none">• Bereinigung der Datensätze• NLP-Pipeline zur Vorverarbeitung	✓ ✓
Funktionen zur Textzusammenfassung	<ul style="list-style-type: none">• 2 Extractive Ansätze: Textrank und LSA	✓
Geeignete Evaluierungsmetrik	<ul style="list-style-type: none">• ROUGE-Score 1	✓

Erfüllung der Anforderungen

Erfüllt

Kompressionsrate 20%-80%



Englische Sprache



Zusammenfassung unterschiedlicher
Textarten



Eigenständige Auswahl der Toolchain (Spacy,
NLTK, Sumy)



Noch offen

Extrahierung wichtigster Informationen



Entwicklungsstand - Klassifikation

Aspekte aus dem PSP

Erstellung der Datengrundlage

- 5 Kategorien: News Artikel, Literatur, Politische Rede, Blogartikel, Juristische Texte



Datenvorverarbeitung

- Vorverarbeitung der Datensätze für die Klassifikation
- NLP-Pipeline (klassische Algorithmen) -> Tokenizer, Lemmatizer, ..
- Datenvorverarbeitung (BERT-Modell) -> Tokenizer



Modellerstellung

- BERT-Model (Transfer Learning)
- Algorithmen, kein Deep Learning -> SVM, KNN, Naive Bayes (Modelltraining)



Modellevaluierung

- Richtig klassifizierte Texte im Vergleich zu allen klassifizierten Texten (Accuracy)



Erstellen der Klassifikationsfunktion

- Erstellen einer Funktion zur Nutzung des Modells in der Webanwendung



Erfüllung der Anforderungen

Erfüllt

Präzise Klassifikation in
Oberkategorien



Auswahl geeigneter
Modelle und Algorithmen



Eigenständige Auswahl der Toolchain
(Spacy, Scit-kit Learn)



Klassifikation englischer Texte



Noch offen



Entwicklungsstand - Webanwendung

Aspekte aus dem PSP

Erstellen der Webanwendung	<ul style="list-style-type: none">• Flask Anwendung (Benutzerfreundliches Frontend, verschiedene Funktionalitäten -> Demo)	✓
Einbinden des evaluierten Klassifikationsmodells	<ul style="list-style-type: none">• Einbindung des vorläufigen Modells (inkl. der Vorerarbeitungsfunktion, BERT)	✓
Einbinden der evaluierten Textzusammenfassungsfunktion	<ul style="list-style-type: none">• Einbindung der vorläufigen Zusammenfassungsfunktion (LSA)	✗
Evaluierung des Frontends nach Benutzerfreundlichkeit	<ul style="list-style-type: none">• Tools zur Evaluierung• Erprobung durch Testuser	✗

Erfüllung der Anforderungen

Erfüllt

Texte können in unterschiedlichen Formaten entgegengenommen werden (Textfeld, ODT, DOC(X), TXT)



Auswahl einer Kompressionsrate



Anzeigen der Resultate



Benutzerfreundlichkeit (Super Zoom, Kontraste)



Noch offen

Finale Textzusammenfassungsfunktion einbinden

Evaluierung der Benutzerfreundlichkeit und Barrierefreiheit

Kurzdemo...



Herausforderungen

- Datenerstellung:
 - Neue Anforderungen (5 Oberkategorien)
 - 5 geeignete Datensätze finden (+ Referenzzusammenfassungen erstellen)
 - Größe der Datensätze
 - Datenvorverarbeitung
- Barrierefreies Frontend

Ausblick



Die nächsten Schritte:

1. Evaluierung der Zusammenfassung
2. Einbinden der Zusammenfassungsfunktion
3. Evaluierung der Gesamtanwendung
4. Abschlussdokumentation/präsentation, Repository