

# Dokumentation der Code-Entwicklung Masterarbeit

Sven Burkhardt

11. April 2025

## 1 Übersicht

Dieses Dokument enthält eine Dokumentation der am 11. April 2025 erstellten oder aktualisierten Code-Komponenten für die Masterarbeit. Diese Module unterstützen den Datenverarbeitungsprozess für historische Dokumente und automatisieren die Extraktion und Anreicherung von Metadaten.

## 2 Module

### 2.1 Dokumenttyp-Erkennung (`type_matcher.py`)

Das Modul `type_matcher.py` ermöglicht die automatische Erkennung und Zuordnung von Dokumenttypen anhand von Dateinamen und/oder XML-Metadaten:

## Dokumenttyp-Erkennung

```
1 def get_document_type(filename: str, xml_path: Optional[str]
  = None, debug: bool = False) -> str:
2     match = re.match(r"(\d{7})_Akte_.*?p(?:age)?(\d+)",
  filename, re.IGNORECASE)
3     if not match:
4         if debug:
5             print(f"[DEBUG] Kein Match f r Dateiname: {
  filename}")
6         return ""
7
8     transkribus_id = match.group(1)
9     page_number = match.group(2).zfill(3)
10
11     row = type_df[
12         (type_df["Transkribus-ID"] == transkribus_id) &
13         (type_df["csv_page_number"] == page_number)
14     ]
15
16     if not row.empty:
17         doc_type = row.iloc[0].get("Dokumententyp", "").strip
  ()
18         if debug:
19             print(f"[DEBUG] Typ aus CSV: {doc_type} f r ID {
  transkribus_id}, Seite {page_number}")
20         return doc_type
21
22     # Fallback auf XML-Metadaten, falls keine
  bereinstimmung in der CSV gefunden wurde
23     if xml_path:
24         try:
25             tree = ET.parse(xml_path)
26             root = tree.getroot()
27             for elem in root.iter():
28                 if elem.tag.endswith("CSVData"):
29                     typ = elem.findtext("Dokumententyp",
  default="").strip()
30                     if typ:
31                         if debug:
32                             print(f"[DEBUG] Fallback-Typ aus
  XML: {typ}")
33                         return typ
34         except Exception as e:
35             if debug:
36                 print(f"[DEBUG] Fehler beim XML-Fallback: {e}
  ")
37
38     return ""
```

## 2.2 Ortserkennung (place\_matcher.py)

Das Modul `place_matcher.py` implementiert einen Fuzzy-Matching-Algorithmus zur Erkennung und Standardisierung von Ortsangaben in historischen Doku-

menten:

## Ortserkennung mit Fuzzy-Matching

```
1 class PlaceMatcher:
2     def __init__(self, csv_path, threshold=90):
3         self.threshold = threshold
4         try:
5             self.places_df = pd.read_csv(csv_path, sep=";")
6             self.known_name_map = self._build_known_place_map
7         ()
8         except Exception as e:
9             logging.error(f"Fehler beim Laden der Ortsdaten
10             aus {csv_path}: {e}")
11             self.places_df = pd.DataFrame()
12             self.known_name_map = {}
13
14     def match_place(self, input_place: str):
15         """Fuzzy-Matching gegen alle bekannten Namen &
16         Alternativnamen"""
17         if not input_place or not input_place.strip():
18             return None
19
20         if not self.known_name_map:
21             logging.warning("Keine bekannten Orte zum
22             Abgleich verf gbar.")
23             return None
24
25         try:
26             # Verschiedene Fuzzy-Matching-Methoden
27             kombinieren f r bessere Ergebnisse
28             match, score, _ = process.extractOne(
29                 input_place.strip(),
30                 list(self.known_name_map.keys()),
31                 scorer=fuzz.token_sort_ratio
32             )
33
34             # Verschiedene Vertrauensstufen zur ckgeben
35             confidence = "low"
36             if score >= self.threshold:
37                 confidence = "high"
38             elif score >= 75: # Mittlere Vertrauensstufe
39                 confidence = "medium"
40
41             if score >= 75: # Wir akzeptieren auch mittlere
42                 Matches
43
44             return {
45                 "matched_name": match,
46                 "score": score,
47                 "confidence": confidence,
48                 "data": self.known_name_map[match]
49             }
50         except Exception as e:
51             logging.error(f"Fehler beim Matching des Ortes '{
52             input_place}': {e}")
53
54         return None
```

### **2.3 Rollenerkennung (`Assigned_Roles_Module.py`)**

Das Modul `Assigned_Roles_Module.py` erkennt Rollen und zugehörige Organisationen für identifizierte Personen:

## Rollenerkennung für Personen

```
1 # Ground Truth Mapping laut CSV (nur deutsche
  Rollenbezeichnungen)
2 ROLE_MAPPINGS_DE = {
3     "ehrenpr sident": "Ehrenpr sident",
4     "ehrenmitglied": "Ehrenmitglied",
5     "vorstand": "Vorstand",
6     "schriftf hrer": "Schriftf hrer",
7     "kassierer": "Kassierer",
8     "sachwalter": "Sachwalter, Notenwart",
9     "notenwart": "Sachwalter, Notenwart",
10    "zweiter vorstand": "ZweiterVorstand",
11    "dirigent": "Dirigent",
12    "chorleiter": "Chorleiter",
13    "ehrenf hrer": "Ehrenf hrer",
14 }
15
16 # Regex zur Rollenerkennung (Textrollen)
17 POSSIBLE_ROLES = list(set(ROLE_MAPPINGS_DE.keys()) | {
18     "vereinsf hrer", "leiter", "obmann", "pr sident" #
    zus tzliche Rollen
19 })
20
21 ROLE_ORG_REGEX = re.compile(
22     r"(?P<n>[A-Z      ][a-z      ]+(?:\s+[A-Z      ][a-
23     z      ]+)?)\s*,?\s*(?P<role>" +
24     "|".join(POSSIBLE_ROLES) + r")\s*(des|der|vom)?\s*(?P<
25     organisation>[A-Z      ][\w\s\~]+)?",
26     re.IGNORECASE | re.UNICODE
27 )
28
29 def assign_roles_to_known_persons(persons: List[Dict[str, str
30     ]], full_text: str) -> List[Dict[str, str]]:
31     """
32     Reiche Rollen und Organisationen f r bekannte Personen
33     anhand des Kontexts im Transkripttext an.
34     """
35     for match in ROLE_ORG_REGEX.finditer(full_text):
36         name = match.group("name")
37         raw_role = match.group("role")
38         organisation = match.group("organisation") or ""
39
40         name_parts = name.strip().split(" ")
41         if len(name_parts) >= 2:
42             forename_candidate = " ".join(name_parts[:-1])
43             familyname_candidate = name_parts[-1]
44
45             for person in persons:
46                 if (person.get("familyname") ==
47                     familyname_candidate and
48                     forename_candidate in person.get("
49                     forename", "")):
50                     person["role"] = raw_role # Original aus
51                     Text
52                     person["role_schema"] =
53                     map_role_to_schema_entry(raw_role)
54                     person["associated_organisation"] =
55                     organisation.strip()
56
57     return persons
```

## 2.4 LLM-Anreicherung (`llm_enricher.py`)

Das Modul `llm_enricher.py` implementiert eine API-basierte Anreicherung historischer Dokumente mit Hilfe von Large Language Models:

## LLM-basierte Datenanreicherung

```
1 def enrich_document_with_llm(json_data: dict, client: openai.  
  OpenAI, model="gpt-4", temperature=0.0) -> Dict:  
2     prompt = f"""  
3     Temperatur: 0,4  
4     Du bekommst ein vollständiges JSON-Dokument aus einem  
      historischen Transkriptionsworkflow.  
5     Deine Aufgabe ist es, folgende Felder **zu ergänzen oder  
      zu korrigieren**:  
6  
7     - 'author'           Wer hat den Text verfasst? Suche nach  
      Gru formeln  
8     - 'recipient'       An wen ist der Text gerichtet?  
      Analysiere das Adressfeld  
9     - 'creation_date'    Nutze Datumsangaben im Text  
10    - 'creation_place'    Oft steht der Ort vor dem Datum  
11    - 'content_tags_in_german' Themen oder Gefühle im  
      Text  
12    - 'mentioned_persons', 'mentioned_organizations', 'mentioned_places' Dubletten entfernen  
13  
14    Besondere Regeln:  
15    - ** Laufenburg (Baden) Rhina** oder ähnliche  
      Kombinationen sind **in der Regel ein Ortsname**  
16    - ** Männerchor Murg** oder ähnliche Begriffe sind  
      **in der Regel eine Organisation**  
17  
18    Wenn ein Feld **nicht eindeutig bestimmbar ist**,  
      verwende "[...]"'.  
19    """  
20  
21    response = client.chat.completions.create(  
22        model=model,  
23        temperature=temperature,  
24        messages=[{"role": "user", "content": prompt}]  
25    )  
26  
27    output = response.choices[0].message.content  
28    input_tokens = response.usage.prompt_tokens  
29    output_tokens = response.usage.completion_tokens  
30  
31    try:  
32        enriched_data = json.loads(output)  
33    except Exception as e:  
34        print("Fehler beim Parsen der LLM-Antwort:", e)  
35        enriched_data = json_data  
36  
37    enriched_data["llm_metadata"] = {  
38        "input_tokens": input_tokens,  
39        "output_tokens": output_tokens,  
40        "cost_usd": round((input_tokens / 1000 * INPUT_COST_PER_1K) +  
41                          (output_tokens / 1000 *  
42                          OUTPUT_COST_PER_1K), 4),  
43        "model": model  
44    }  
45  
46    return enriched_data
```



### 3 Zusammenfassung

Die am 11. April 2025 implementierten oder aktualisierten Module bilden zusammen eine Pipeline zur Verarbeitung und Anreicherung historischer Dokumente aus dem Transkribus-System:

1. **Dokumenttyperkennung:** Automatische Erkennung des Dokumenttyps (Brief, Postkarte, etc.) auf Basis von Dateinamen und Metadaten
2. **Ortserkennung:** Fuzzy-Matching von Ortsnamen gegen eine Ground-Truth-Datenbasis zur Standardisierung
3. **Rollenerkennung:** Erkennung und Standardisierung von Personenrollen in Vereinen/Organisationen
4. **LLM-Anreicherung:** Nutzung von Large Language Models zur automatischen Anreicherung der JSON-Daten mit bisher unerkannten Metadaten

Die Module sind Teil eines größeren Systems zur Erfassung und semantischen Anreicherung historischer Dokumente und unterstützen die digitale Erschließung der Sammlung des Männerchors Murg.