

Auswertung Archivdokumente zu Kunstobjekten



Henry Rotzoll (DVZ MV GmbH)



Gliederung

- 1. Einführung und Fragestellung, Hintergrund
- 2. Datengrundlage und Umsetzung
- 3. Ergebnisse
- 4. Erkenntnisse und Empfehlungen
- 5. Fazit und Ausblick



Einführung und Fragestellung

Die Digitalisierung stellt Unternehmen und Organisationen vor große Herausforderungen

- sie verfügen über Archive mit großen Mengen an Dokumenten, die in eine gemeinsame Datenbank überführt werden sollen
- Dokumente wie Word- und PDF-Dateien enthalten unstrukturierte Daten, die vor der Datenintegration in ein strukturiertes Format überführt werden müssen
- Zusammenführung heterogener Daten ist ein zeitintensiver und fehleranfälliger Prozess, der mit hohem manuellen Aufwand verbunden ist

Kann der Prozess der Datenstrukturierung und -integration mithilfe von generativer KI automatisiert werden?



- Masterarbeit von Henrik Bongertmann (Universität Rostock)
- In Zusammenarbeit mit dem Sachgebiet GEO in der Zeit von April bis Oktober 2024
- Bestehendes Projekt mit der Nordkirche: Kunstguterfassung der Evangelisch-Lutherischen Kirche in Norddeutschland
 - Ziel: erfasste Kunstgüter in eine relationale Datenbank zu integrieren
 - Zweck: die in den Kirchengemeinden vorhandenen Kunstobjekte über eine interaktive Karte zugänglich zu machen
- 4 verschiedene Open Source Sprachmodelle verglichen
 - 30 Testdokumente verwendet
 - Gewinner: Llama-3.1-8B
- Laptop mit Intel Core i7, 32 GB Arbeitsspeicher und Nvidia GeForce GTX 1070 mit 8 GB VRAM für die Prozessierung verwendet



Datengrundlage

6509 Dokumente in den Dateiformaten PDF, DOCX und DOC

Dokumente enthalten Objektbeschreibungen aus der Kunstguterfassung

Propose Strateural

Alvenshagen

Kangslutur Mitte 16. Jb. (Reinser 16. 1) gestiffet vermallich von Hans Albrecht von Phischoe (Plinkow) d. J.



Datierung: um 1741/50

Material: Holz, fetting gefassil und tellvergolder,

Fassung: Fassung in 20. Jt.; 1960er oder Yor-Jahre?) errecert, Astritektur in versithederen hellen Graufforen int engoldelen Lieders, die Stalverschafte blau mit helte Mamorienium und vergodelen Kapitellen, Orsanentik feltrengistet, Williten auf Versitherung blau gebatert, Figuren monochrom welft, deutstet Reite seine Stenier Fassung?

Maller, H.B., TS nr. (R. Reimer, Tugeroffigur finite 125 pm, rechts 127 cm teachraftes: or inschrößbartusche am Scholdenket, Helige mistri in demer / Waterheir, Pusto Innic. "Du sollat Gott demen Herre laden von ganzem Herzen", Pullto rechts, "Du sollat deinen nächsten leben als dath selbet"

Wapper: Predefections size von Pilantee (Pilatine; das Wapper int dunts die Nieubauung der geben, die der reine der State (der Pilatine) der State (der S

Beschwärung/Bereichsager, breigeschseiger zuchleichsischer Auftau, als mit dem Gefreisager, Placialie- und Dilmer-Ormwert geschmigteit in Ausstalier im Meterote met von geweit zurs Stallen zu der Sollen gerahmt, welche von stalle profilierne werkrightes Gebüld troper, zweichen dem Stallen zu dem Sollen gerahmt, welche ein stalle profilierne werkrightes Gebüld troper, zweichen dem Stallen zur dem Faunde und der Stallen zur dem Stallen geschlich und der Vertrech. Vermutlich handet die sich von False und Spes (Stallet und Haufburger). Die Oberprecht und von einer Wöldungspielen mit dem Auge Gebähn mit dem Australizaben Gebüldung der Stallen und dem Gebüldungspielen und dem Gebüldungspielen der dem Auge der Vertrecht und dem Gebüldungen der der Stallet und dem Gebüldungen einem Außer mit dem Gebüldungen einem Auge der Vertrecht und dem Gebüldungen der dem Stallet und dem Gebüldungen einem Ausstallet vierande siehen dem Augen der dem Gebüldungen der dem Augen der dem Gebüldungen der der dem Gebüldungen der der dem Gebüldungen der dem Gebüldungen der der dem Gebüldungen der der der dem Gebüldungen der der der dem Gebüldungen der dem Gebüldungen der der der dem Gebüldungen der der dem Gebüldung

AALLING 2007, S. 567, delite die "harcole Fartemposch" des Klangelijtens als setten haraus. ³ Siete <u>Impachte wildendie organisch" PECTAD Cables. (Adeleganthoute.</u> (Abert 02.10.2017). Kircheskreis Diffmarichen, Hennie, St. Minies.

Beichemül/Pactoreumth

Standsen Kirche

Made: H cz. 229 cm, B cz. 250 cm; T cz. 137,5 cm

Somen der Wandräfelung, Rückwand: H cz. 98,2 czz. B ca. 45,8-46 czz.

Tegenden (Front): 1f ca 50 - 86,5 cm; B ca 46,7 - 48 cm

Tugenden (Tür): H ca. \$5.5 cm. B ca. 50.5 cm

Manerial: Holz, Essen, Farbfassing

Datierung: 1622/1699 (ungeben der Kinnstopographie), 1709 (Audemeite Zispang)

Inschriften

Überschriffen der Tugenden: Der Globe i die Liebe / die Hoffinng / die gedoß / die demath / die Voruchtiebeith

Subscriptio (Innouraum); GENES XXL1 / GENES XXXXIX / GENES XXXVII.

Baschriften: Der Glaube / die Liebe / die Hoffmag / die gedalt / die demath / die Verschrigbeith

Werkstatt Kinsder; unbekamt

Beschreibung Bemerkungen:

Das einzehlige Knobengestild mit vergittertem Sichtlichter wurde auf der Evangeliemerbe aus Chorbogen im Chorbereich aufgestellt. Ein archärktreisieber Auffrun mit kassettierten Feldem und Abschlausgesims bestimmt die optische Erscheimung. Das Gesims ist mit Zubmicharlt und Daminifrossen shytlamisiert. In die Britistungsfelder sind farbige schlichte Personsfikationen der Togenden auf das Tragermaterial gerault wurden. Die christlichen Tragenden wurden mit zwei der Kardinsfitugenden erweitert. Glaube, Liebe, Hoffmung, Geduld, Demuit und im Inneutzum die Vorsicht.

Der obere Abschaft des uszumholten Stulis ist zur Gemeissle, zu der Lingsveite wie zum Altur vergittert. Die versikal und horizontal ausgeführten Stülee und gefant. Die vergitterte Lingsveite ist den Ertisstangsfeldern entsprechend viergetreilt, wobei die beiden zurfleren über eine Scharbefführung zu öfflines und, ähnlich verhält es sich mit den beiden zur Gemeinde ausgesichteten Kompartamenten. Die zum Altur ausgenichtete Tilt ist über einem Hebel-Bügel-Verschäuse zu betätigen. Umpränglich wur ein Schäuss angebracht.

Die Rückvand des Innerannen beidet vor Felder der Wandtäfelung unt ein, und thenanisiert den Kaupf Jocobs nut dem Engel - Josef wird von seinen Bridern in den Brunnen gestessen - Josef soll

Kunst- und Kulturgut im Bereich des Kirchenkreises Mecklenburg

Ort / Örtliche Kirche

Goldberg

Ortskennzahl

407

Kirchengemeinde

Bezeichnung

ausgebautes Grabplattenfragment

Anzahi

Inventarnummer

060

Verweis auf zugehörige Objekte

059

Aufbewahrungsort / Standort

auf dem oberen Absatz der Treppe des Nordanbaus auf dem Fußboden lagernd

Abweichender Herkunftsort

vorher eingelassen auf der Empore in die Nordwand östlich des Eingangs zum

Treppenaufgang

Material/Technik

Stein

Maße

H: 67 cm

.

B: 102,5 cm

T: cm

L: cm

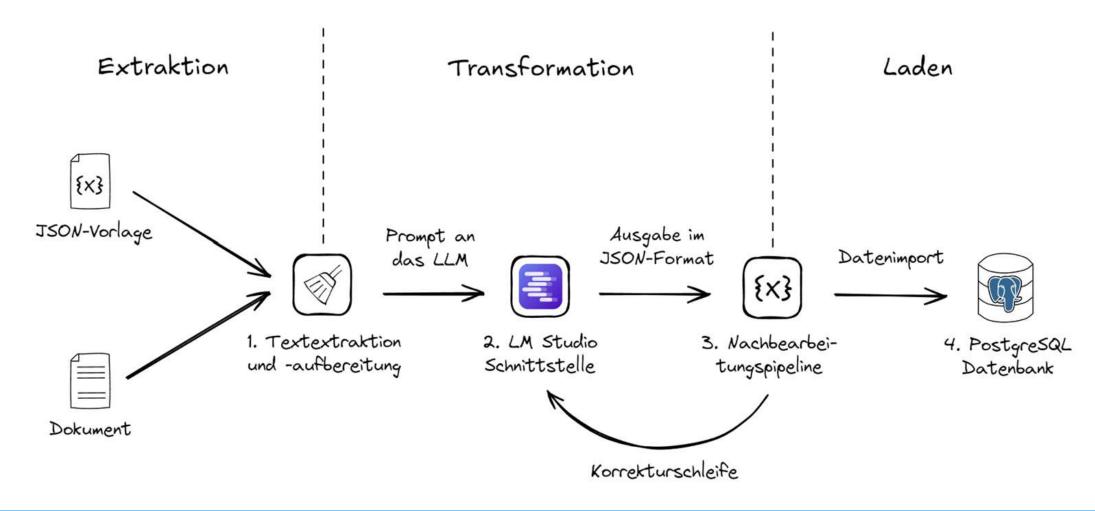
Gewicht: g

Künstler/Werkstatt



Umsetzung

Die Schritte im Implementierungsprozess





Prompt zur Strukturierung von Dokumenten in ein gültiges JSON-Format

In der folgenden Objektbeschreibung wurden Informationen über ein
bestimmtes Kunstgut erfasst:
{document_text}

Aufgabe: Übertrage alle Informationen aus der Objektbeschreibung in ein gültiges JSON-Objekt anhand der folgenden Vorlage: {json_template}

Hinweis: Die Informationen sollen ungekürzt und unverändert übertragen werden. Übertrage für jedes Merkmal den gesamten Text bis das nächste Merkmal beginnt. Es dürfen keine Informationen verloren gehen.

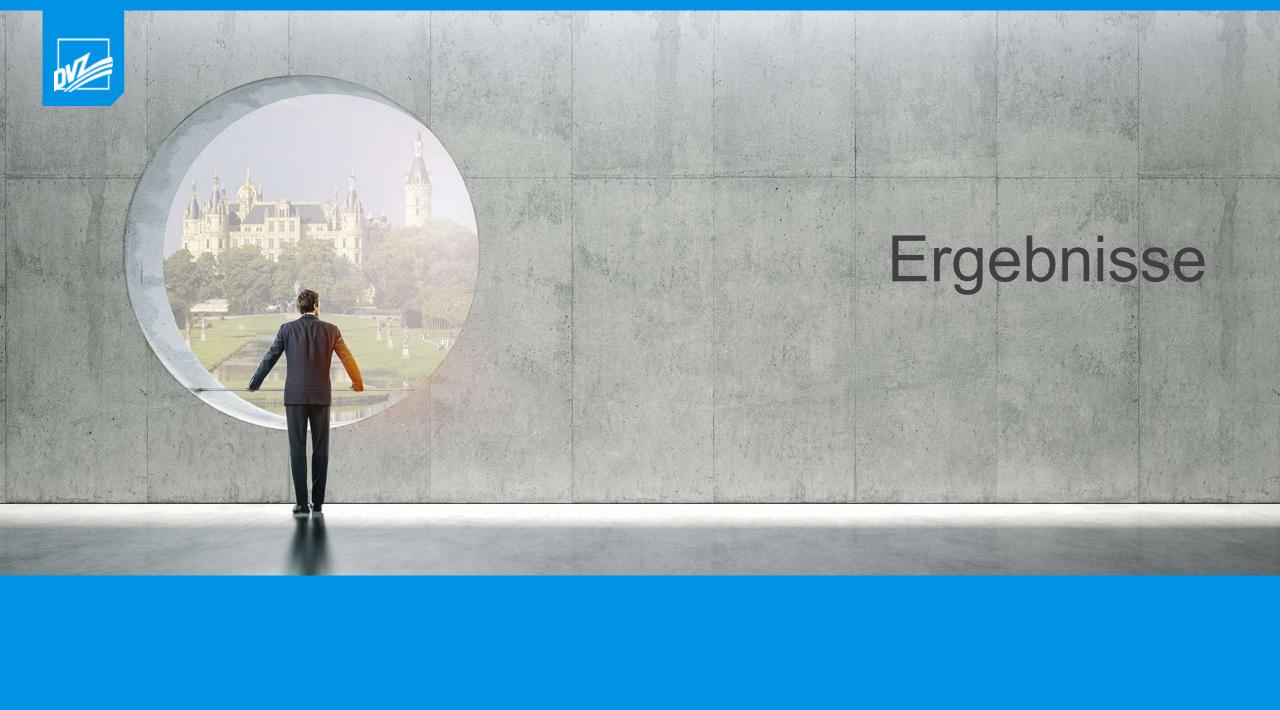
Verwende die vorgegebenen Schlüssel aus der Vorlage für das JSON-Objekt in der Antwort. Formatiere alle Werte im JSON-Objekt als Text-String. Setze die Schlüssel und Werte in doppelte Anführungszeichen, um sicherzustellen, dass die Antwort der korrekten JSON-Syntax entspricht. Die Antwort soll ausschließlich das JSON-Objekt und keinen zusätzlichen Text enthalten.



Datenmodell

Vorlage im JSON Format (29 Merkmale)

```
| json_template = [
      "Ortoname:";
      "Kirche": "",
      "Ortskenngahl": "",
      "Kirchengemeinde": "",
      "Bezeichnung"; "",
      "Anxabl": "",
      "Inventarnummer": "",
      "Verweis auf zugehörige Objekte": "",
      "Aufbewahrungsort/Standort": "",
      "Material/Technik": "",
      "Halle": "",
      "Gewicht": "",
      "Künstler/Werkstatt": "",
      "Entstehungsort/Herkunft": "",
      "Datierung": "",
      "Beschreibung/Bemerkungen": "",
      "Bildthena/Bildprogramm": "",
      "Inschriften": "",
      "Hetallmarken und Auflösung": "",
      "Zustand/Schaden": "",
      "Fassung": "",
      "Restaurierung": "",
      "Leibvertrag": "",
      "Literatur": "",
      "Bistorische Bildquellen/Archivalische Quellen/Biografische Verweise": "",
      "Stifter": "",
      "Foto Nummern": "",
      "Erfasser": "",
      "Erfansungsdatum/Letzte Anderung": ""
```





Finetuning Modell-Parameter

Temperatur-Parameter bestimmt die Modellkreativität bei der Antwortgenerierung

- Kreativität zwischen 0 und 2, wobei 0 möglichst deterministische Antworten erzeugt
- Testdurchlauf mit 30 Dokumenten ohne Nachbearbeitung der Ausgabe
- Niedrige Temperatur-Parameter (0,1 0,4) konnten höhere Erfolgsrate aufweisen

Parameter	Erfolgsrate	Antwortzeit	Antwortlänge	Token/sec
temp = 1.0	17 von 30 (57%)	72 s	1134 Token	16 t/s
temp = 0.7	18 von 30 (60%)	75 s	1187 Token	16 t/s
temp = 0.4	22 von 30 (73%)	75 s	1181 Token	16 t/s
temp = 0.1	21 von 30 (70%)	75 s	1175 Token	16 t/s

Tabelle 5: Vergleich der Erfolgsraten bei der Strukturierung von 30 Dokumenten in das Zielformat JSON mit unterschiedlichen Anpassungen der Modell-Kreativität



Nachbearbeitungspipeline

Vergleich der Erfolgsraten mit und ohne Nachbearbeitung der Modellantworten

- Bereinigung der generierten Ausgabe soll gültiges JSON-Format sicherstellen
- Korrekturschleife bei fehlerhaften Modell-Ausgaben, die nicht korrigiert werden können
- Effektive Nachbearbeitung verbessert Erfolgsrate signifikant

Performance/	Erfolgsrate	Korrektur-	Ungültige	Laufzeit
Strukturierung	Gesamt (in %)	läufe	Antworten	gesamt
ohne Nachbearbeitung	22 von 30 (73%)	-	8 von 30 (27%)	37 min
mit Korrekturschleife	23 von 30 (77%)	8 (27%)	15 von 38 (40%)	44 min
mit Nachbearbeitung	29 von 30 (97%)	-	1 von 30 (3%)	38 min
beides zusammen	30 von 30 (100%)	1 (3%)	1 von 31 (3%)	39 min

Tabelle 6: Vergleich der Erfolgsraten bei der Strukturierung von 30 Dokumenten in das Zielformat JSON mit Korrekturschleife und Nachbearbeitung der Ausgabe



Ergebnisse der Fallstudie

5443 Dokumente konnten mit einer Erfolgsrate von 96% importiert werden

Quell- Ordner	Anzahl Dokumente	Erfolgreich importiert	Korrektur- läufe	Antwort- länge	Antwort- zeit	Laufzeit gesamt
A	842	829 (98%)	55 (7%)	698 Token	38 s	9h 46min
В	1729	1642 (95%)	79 (5%)	651 Token	36 s	ca. 18h
\mathbf{C}	1199	1124 (94%)	69 (6%)	786 Token	48 s	ca. 17h
D	1900	1848 (97%)	85 (5%)	817 Token	49 s	ca. 27h
Gesamt	5670	5443 (96%)	288 (5%)	742 Token	43s	ca. 71h
Schnitt	100	96 (96%)	5 (5%)	742 Token	43s	1h 15min

Tabelle 8: Performance des LLM-basierten ETL-Prozesses für verschiedene Datenquellen



Fehlerursachen

515 Fehlversuche auf Dokumente (24%) und LLM-Antworten (76%) zurückzuführen

Quell- Ordner	Durch- läufe	Fehl- versuche	Leere Datei	Token -limit	JSON ungültig	Spalte ungültig
A	897	68 (7%)	0 (0%)	2 (3%)	22 (32%)	44 (65%)
В	1808	166 (9%)	60 (36%)	2 (1%)	59 (36%)	45 (27%)
C	1268	144 (11%)	47 (33%)	3 (2%)	59 (41%)	35 (24%)
D	1985	137 (7%)	6 (5%)	3 (2%)	36 (26%)	92 (67%)
Gesamt	5958	515 (9%)	113 (22%)	10 (2%)	176 (34%)	216 (42%)

Tabelle 9: Fehlerursachen bei der LLM-basierten Strukturierung von Dokumenten



Erkenntnisse





Die Antwortqualität von LLMs wird durch viele Faktoren beeinflusst

- Viele Stellschrauben, die optimiert werden können
 - Textaufbereitung, Datenmodell, Prompt-Engineering und Modell-Parameter
 - Auswirkungen sind nicht vorhersehbar, da die Funktionsweise eine Black-Box ist
 - Optimierungsprozess nach dem Trial-and-Error-Prinzip
- Nachbearbeitungspipeline konnte die Erfolgsrate signifikant erhöhen
 - LLMs sind keine Alleskönner, sondern nur ein Werkzeug
 - Eine effektive Nachbearbeitung der Modellausgabe ist daher unerlässlich
- Optimierungen müssen individuell auf den Anwendungsfall abgestimmt sein



Fazit und Ausblick

Ergebnisse zeigen Potenzial von LLMs bei der automatisierten Datenstrukturierung

- Erfolgreiche Fallstudie: Über 5000 Dokumente wurden in die Datenbank importiert
- Open-Source-LLMs k\u00f6nnen unstrukturierte Texte in strukturierte Daten \u00fcberf\u00fchren
- LLM-basierter ETL-Prozess kann manuellen Aufwand erheblich reduzieren
- Optimierungsmaßnahmen konnten Ergebnisqualität signifikant steigern
- Ansatz ist auf andere Anwendungsfälle übertragbar (mit kleineren Anpassungen)
- Es können Empfehlungen für ein allgemeines Vorgehen abgeleitet werden





Fragen?



Kontakt









www.dvz-mv.de/twitter

Henry Rotzoll

+ 49 385 4800 532

h.rotzoll@dvz-mv.de

IHR ANSPRECHPARTNER



www.dvz-mv.de







