

Integration von OCR-D in Kitodo: die ersten Implementierungsschritte und praktische Anwendung

Katja Rykhlinskaya, Sven Marcus, Michael Kotzyba, Robert Strötgen

Allgemeine Projektinformationen

- Projektförderung:
 - DFG, Dauer 2 Jahre
- Projektbeteiligte Institutionen:
 - UB Braunschweig
 - SLUB Dresden
 - UB Mannheim
- Projektaustausch:
 - BS-intern wöchentlich
 - Projektintern monatlich
 - OCR-D-Community







Was macht OCR-D?

Binarisierung (aus Graustufenbild wird S/W-Bild generiert) - Zuschneiden (Ränder werden entfernt) - Entrauschen (Flecken werden entfernt) - Entzerrung (Schräglage der Seite wird korrigiert) - Entschärfung (Textzeilen werden gestreckt) - Segmentierung (Regionen werden erkannt) - Erkennung (In jedem Region wird der Text erkannt)





Jehund folget

Der

Der Dir Cheil/

Daraus die Karte inventiret/und drauf gerichter worden:
da sich allezett ein Blat auf eine His
sort ist referiret: und nut emander.
XXXVI. benderseits anzu.
tressen seine.

(nierde/ daß ich in der Kars.

(Merce / daß ich in der Kars ten Kingel: Reimen allezeit alludiret babe auf den Spruch Rom. 12. v. 14. Inchet nicht/3c.) 40. 17ünliche Spiele Jekund folget Der

Mnder Theil/

Daraus die Karte inventiret/und drauf gerichtet worden: da sich allezeit ein Blat auf eine Hiflorie/ so mit gleicher Zahl verzeichnet ist/ referiret zund nut einander. XXVI. beyderseits anzutresten sennd.

(Merce / daß ich in der Karsten Kingel: Neimen allezeit alludiret babe auf den Spruch Rom. 12. v. 14. Juchet nicht / 3c.)

40. Müsliche Spiels Jekund folget Det

Ander Theil/

Daraus die Karte inventiret/und drauf gerichter worden:
da sich allezeit ein Blat auf eine Dia
florie / so mit gleicher Zahl verzeich
net ist referiret: und nur einander.
XXX VI. beyderseite anzutressen sond.

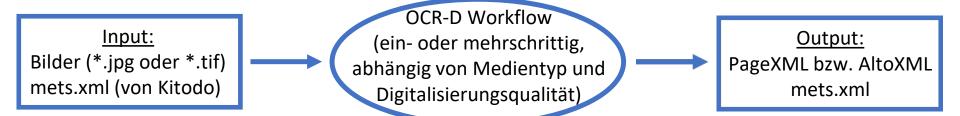
(Merce / daß ich in der Kar, ten Kingel : Reimen allezeit alludiret habe auf den Spruch Rom. 12. v. 14. Sluchet nicht / 3c.)

Ergebnis (Ausschnitt aus einer Alto-Datei):



Aktuelle OCR-D Anwendung

Die bereits digitalisierten Medien können auf dem OCR-D Server prozessiert werden!



Einschrittiger Workflow:

ocrd-tesserocr-recognize mit dem Model "Fraktur_GT4HistOCR" (kann gleichzeitig Binarisierung, Regionalsegmentierung, Tabellenerkennung, Liniensegmentierung und Texterkennung durchführen).

Merhschrittiger Workflow:

ocrd-sbb-binarize mit dem Model "default-2021-03-09"
ocrd-eynollah-segment mit dem Model "default"
ocrd-cis-ocropy-dewarp
ocrd-tesserocr-recognize mit dem Model "frak2021-0.905"

Page-to-Alto Umformung:

ocrd-fileformat-transform



Erkennungsbeispiel 1: Pharmazeutische Zeitung J. 1872

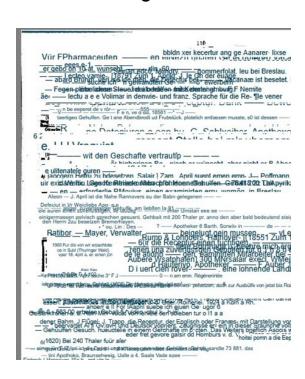
Originalbild



Mehrschrittiger Workflow



Einschrittiger Workflow





Erkennungsbeispiel 2: Kinderbuch J.1819

Originalbild

Rurze Satze zur Erweckung der Aufmerksams feit und des Nachdenkens. 3ch gehöre zu den Kindern. Kinder wissen noch nicht viel, und darum mussen sie unterrichtet werden und lernen. Dadurch werden sie verständig. 3ch werde in der Schule von Lehrern unterrichtet. 3ch bin meinem Lehrer Dankbarkeit und Gehorsam schuldig. So lange ich unterrichtet werde, bin ich ein Schüler. Ein guter Schuler ist aufmerksam; er hört nur auf das, was der kehrer sagt, und denkt nur an das, was er thun, oder begreifen und behalten soll.

Mehrschrittiger Workflow

Kurze Såtze zur Erweckung der Aufmerkſam» keit und des Nachdenkens.

C&

I gehöre zu den Kindern. Kinder willen noch nicht viel, und darum müllen lie unterrichtet werden und lernen. Dadurch werden lie verftändia.

Ich werde in der Schule von Lehrern unterrichtet. Ich bin meinem Lehrer Dankbarkeit und Gehorfam schuldig. So lange ich unterrichtet werde, bin ich ein Schüler. Ein guter Schüler ist aufmerksam; er hört nur auf das, was der Lehrer sagt, und denkt nur an das, was er thun, oder begreifen und behalten soll.

Einschrittiger Workflow

.Kurze Såtze zur Erweckung der Aufmerkſamkeit und des Nachdenkens.

l∍oa gehőre zu den Kindern. Kinder willen noch nicht

vi el, und darum mullen lie unterrichtet werben ind ler nen. Dadurch werden lie verständig.

Ich werde in der Schule von Lehrern unerrichtet. Ich bin meinem Lehrer Dankbarkeit und Gehorfam ſchuldig. So lange ich unterrichtet werde, bin ich ein Saler. Ein guter Schuler ift a ufmerkſam; er hort u auf das, was der Lehrer ſagt, und denkt nur an daß, baß thun/ oder begreifen und behalten ſoll.



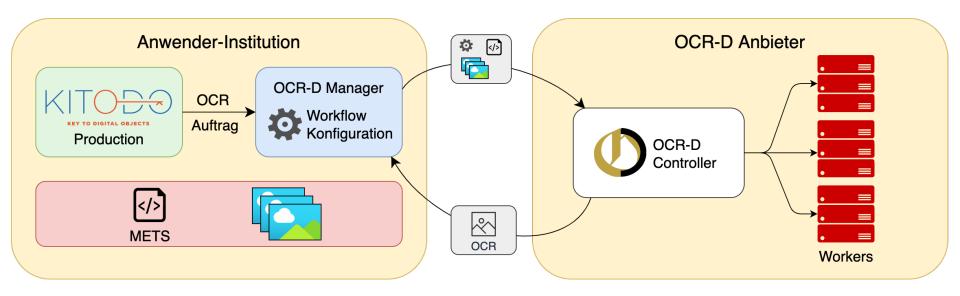
IV. Von der Erde und ihren Bewohnern

Daß die Erde fehr groß, aber doch nur ein kleiner Then

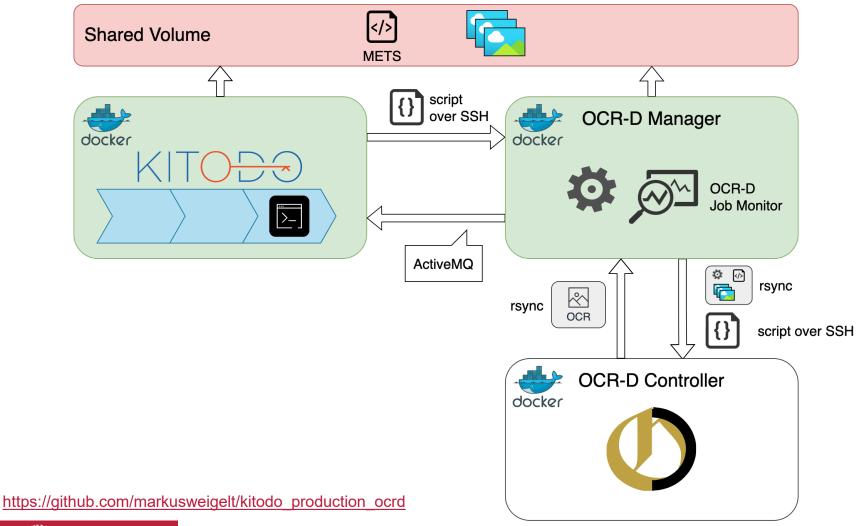
der Welt fei, haben wiv schon gehört. Was für eine Gesftalt die Erde habe, ift schwer auszumachen, weil man nur einen. Sehr kleinen Theil der Erde auf ein Mal übersehen kann, und weil sie uns, zu nahe ist. — Aus dem Schatten eines Kötbers kann man mit ziemticher Gewißheit erkennen, ob er rund, breit, oder ekkig und spitzig sei; And wenn der Schatten eines Körpers von allen Seiten alle Mal, so oft er sich zeigt, rund erscheint, sa ist nicht zu zweiseln, daß auch der Körper rund sei. Dies ist nun der Fall bei unserer Erde. Jyr habt wohl schon von Mondsinsternissen gehört? Bei diesen erblickt man in der Mondscheibe älle Mal einen runden Schatzen, und es ist ausgemacht, daß dieser Schatten von unserer lErde in den Mond geworsen wird, so vit sie bei ihrem Ams



Geplante Implementierung



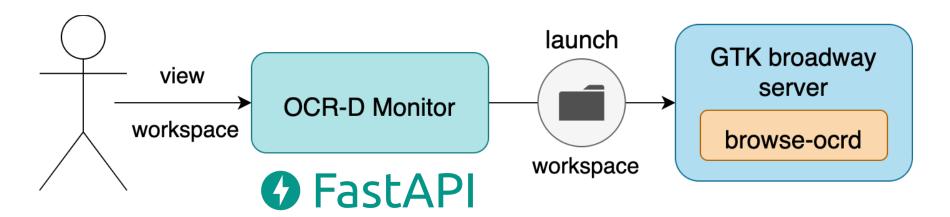
Prototyp



OCR-D Monitor

Teil des OCR-D Managers

- Status von OCR-D Jobs überwachen
- Mit Ergebnissen interagieren



Weitere Projektschritte

- Fertigstellung der automatisierten Kitodo-OCR-D Verbindung
- Optimierung der Workflowauswahl (Geschw.-Qualität Verhältnis!)
- MyCore-DFG Viewer Verknüpfung (seitens OCR-D)



Danke für Ihre Aufmerksamkeit!

