

TESSERACT: TURBO UND TRAINING

Daniel Brenn
Universitäts- und Landesbibliothek Sachsen-Anhalt
daniel.brenn@bibliothek.uni-halle.de



UNIVERSITÄTS- UND
LANDESBIBLIOTHEK
SACHSEN - ANHALT

Gliederung

- Warum OpenSource?
- Tesseract mit neuronalen Netzen
- Turbo in der Digitalisierung
- Training für besondere Schriften



Ausgangslage

- Digitalisierung läuft aktuell noch über Visual Library der Firma semantics
- OCR mit AbbyFine
- Volumenlizenzen zur Digitalisierung
- festgelegt auf bestimmte Sprachen und Schriften
- OCR nicht für alle Digitalisate möglich, kostenintensiv und eher unflexibel

- Entscheidung für Nutzung von Kitodo und DSpace
- Aufbau eines neuen, modularen Digitalisierungsworkflows mit OpenSource-Komponenten
- Plan zur Nutzung von AbbyCloud in Kombination mit Tesseract OCR



Gliederung

Warum OpenSource?

- Tesseract mit neuronalen Netzen
- Turbo in der Digitalisierung
- Training für besondere Schriften



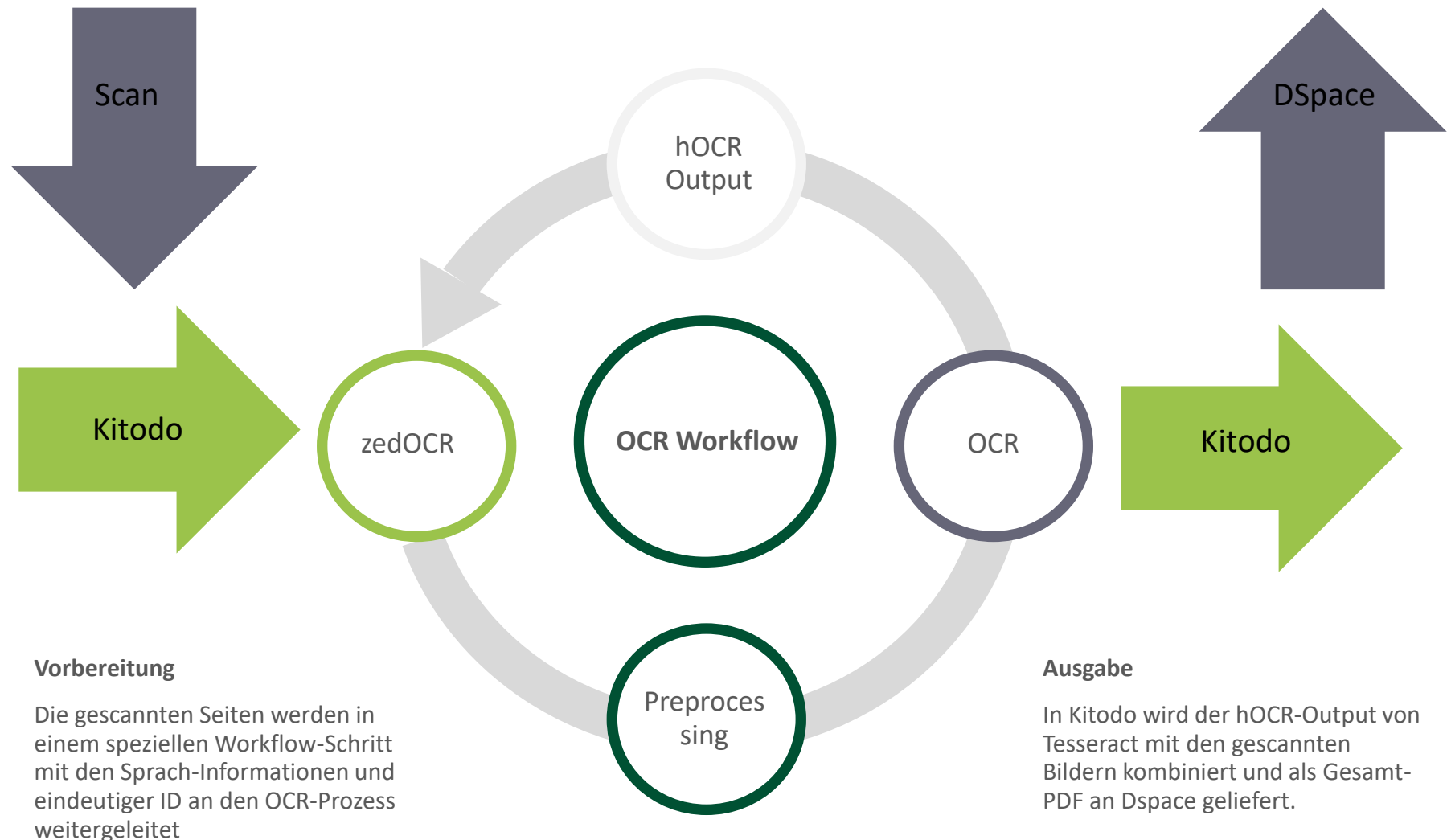
Warum OpenSource?

- Entscheidung zum Einsatz von OpenSource-Software für größere Flexibilität
- Insbesondere Tesseract bietet durch lange Entwicklung Zugriff auf sehr breites Feld von Trainingsdaten
- Nutzung für besondere Projekte durch eigene Weiterentwicklung
- Gute Möglichkeiten zur Kollaboration mit anderen Bibliotheken/ Projekten

- aber: OpenSource benötigt intensive Auseinandersetzung mit der Software und häufig mehr Arbeitszeit!
- zentral wichtig: transparente Dokumentation, um Wissen nicht nur bei einzelnen Kollegen zu bündeln



Warum Open Source? - Digitalisierungsworkflow

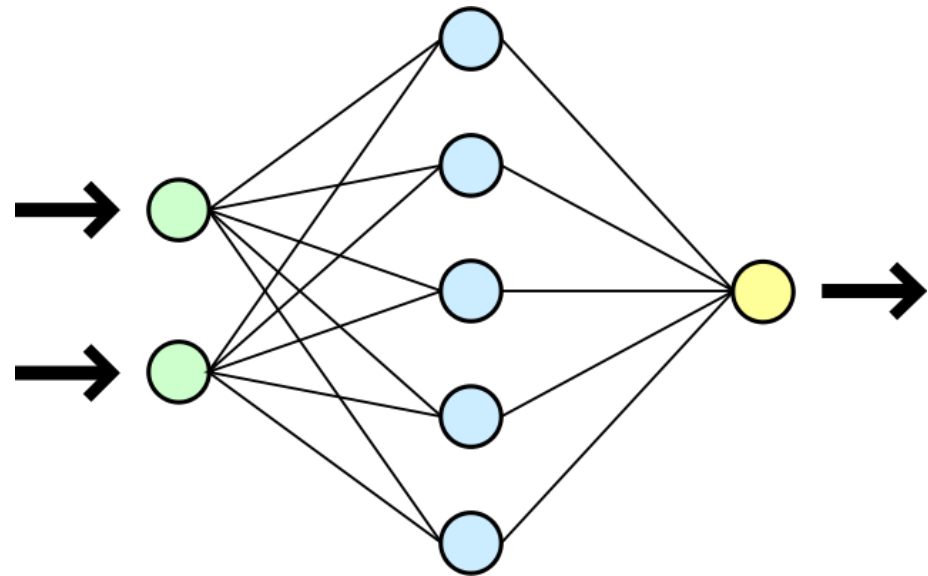


Gliederung

- Warum OpenSource?

Tesseract mit neuronalen Netzen

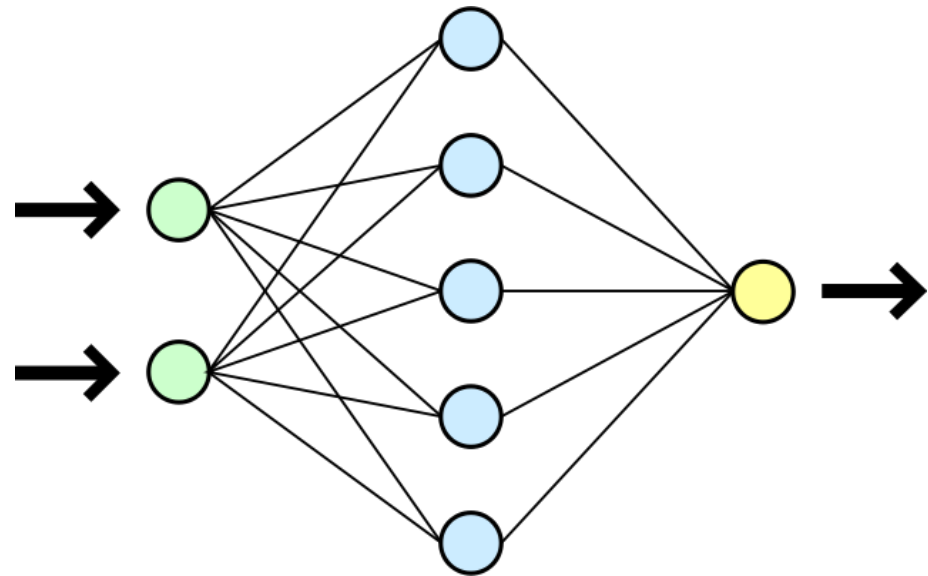
- Turbo in der Digitalisierung
- Training für besondere Schriften



Quelle: Dake, Mysid –
CC BY 1.0, <https://commons.wikimedia.org/w/index.php?curid=1412126>

Tesseract mit neuronalen Netzen

- Tesseract nutzt seit Version 4 neuronale Netze
- konkret: LSTM (Long Short-Term Memory)
- dadurch erhebliche Verbesserung der Genauigkeit bei schwierigerem Layout und schlechter lesbaren Texten erkennbar
- ermöglicht OCR auch bei Schriften, die vorher nur schwer oder gar nicht verwendbar waren
- erhöhter Leistungsbedarf im Vergleich zu früheren Versionen



Quelle: Dake, Mysid –
CC BY 1.0, <https://commons.wikimedia.org/w/index.php?curid=1412126>

Gliederung

- Warum OpenSource?
- Tesseract mit neuronalen Netzen
- Turbo in der Digitalisierung
- Training für besondere Schriften



UNIVERSITÄTS- UND
LANDESBIBLIOTHEK
SACHSEN - ANHALT

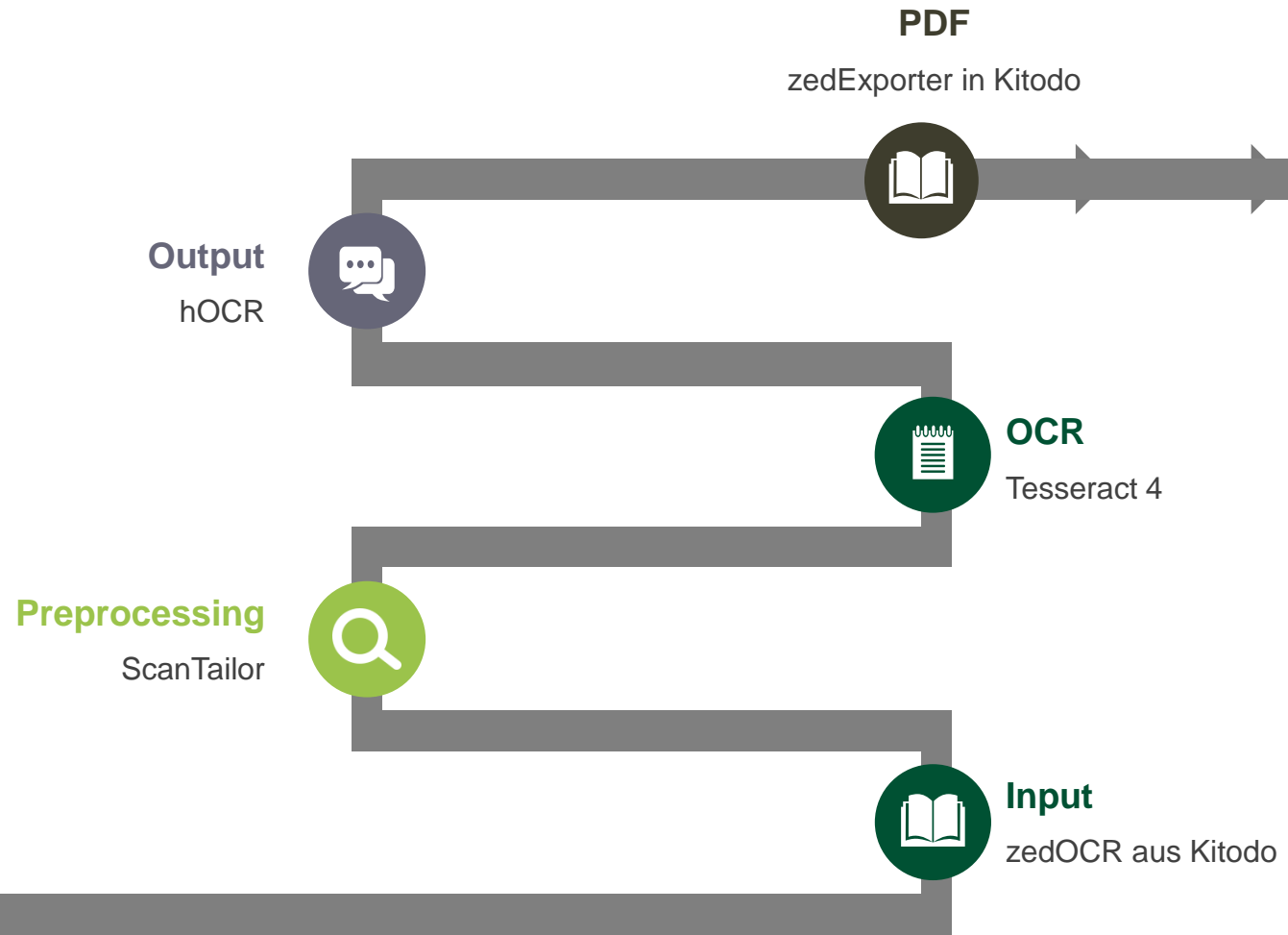
Turbo in der Digitalisierung

- für regulären Digitalisierungsworkflow ist relativ hohe OCR-Geschwindigkeit erforderlich
- nur Teilschritt der Digitalisierung
- Tesseract kann Prozesse auf einzelne Prozessorkerne auslagern und so parallelisieren
- „ab Werk“ bis zu 4 parallele Prozesse
- Zur Verbesserung der Geschwindigkeit ist weitere Parallelisierung mit pyesseract möglich
- durch Installation auf VM ist ein nachträgliches Erhöhen der Kapazität möglich, die Vorplanung sollte die Notwendigkeit dafür aber minimieren
- Qualitätssicherung erfolgt vor Freischaltung im Repository



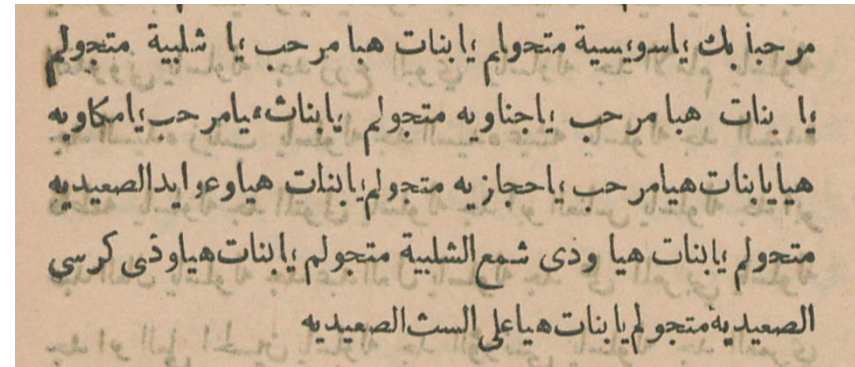
Automatisierter Digitalisierungsworkflow

Automatisierter Workflow,
eingebunden in die
Digitalisierung mit Kitodo



Gliederung

- Warum OpenSource?
 - Tesseract mit neuronalen Netzen
 - Turbo in der Digitalisierung
- Training für besondere Schriften

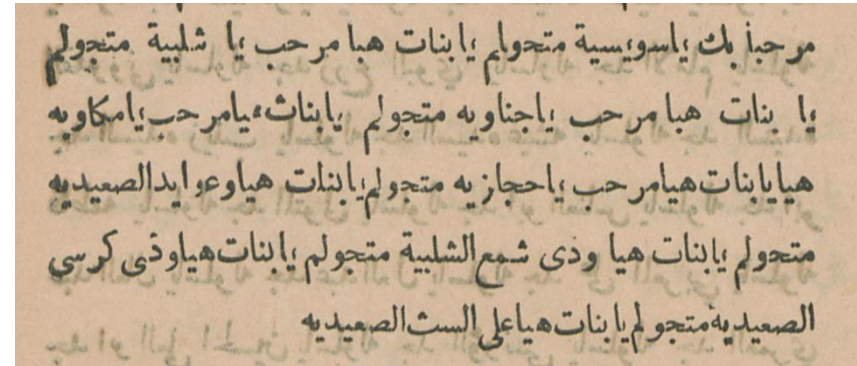


مرحباً بك: أسورية متحوام يابنات هبا مرحب وا شلبية متجولم
يابنات هيا مرحب يا جناويه متجولم يابنات هيا مرحب يا مكاويه
هيا يابنات هيا مرحب يا حجازيه متجولم يابنات هيا وعوايد الصعيديه
متجولم يابنات هيا ودي شمع الشلبية متجولم يابنات هيا ودي كرسى
الصعيديه متجولم يابنات هيا على السب الصعيديه

Quelle: Darwiš, Kāmal Afandī: Kitāb bida‘ al-fağār fi ḥaflat az-zār,
Kairo 1911, urn:nbn:de:gbv:3:5-86614, CC-BY-SA 4.0.

Training für besondere Schriften

- in ersten Tests erreichte Tesseract auch beeindruckende Ergebnisse für als schwierig einzustufende Schriften, insb. Arabisch und Persisch
 - Titel ab ca. 2. Hälfte des 20. Jhd. erreichten eine Genauigkeit von ca. 95%
 - Titel aus der ersten Hälfte des 20. Jhd. erreichten mit den Trainingsdaten von Tesseract 85-90%
- der neue Digitalisierungsworkflow kommt damit direkt unserem FID Nahost zugute
- im Rahmen der aktuellen FID-Förderphase ist zudem Training für schlechter zu lesende Schriftarten geplant
- Die hieraus gewonnen Erkenntnisse werden natürlich auch für weniger exotische Schriften nutzbar sein



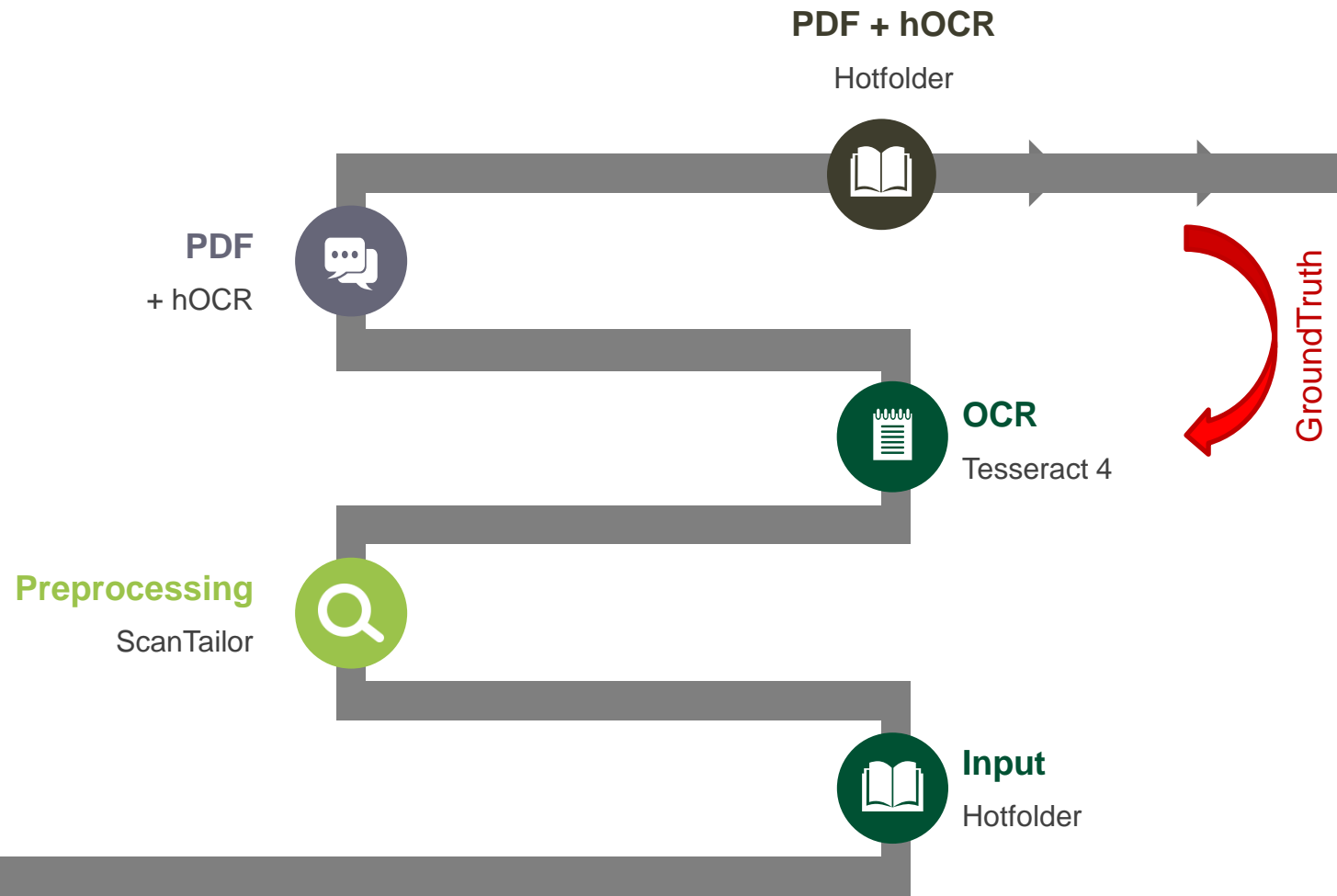
مرحبا بك : أسورسية متحوام يابنات هبا مرحب وا شلبية متجولم
يابنات هيا مرحب يا جناويه متجولم يابنات هيا مرحب يا مكاويه
هيا يابنات هيا مرحب يا حجازيه متجولم يابنات هيا وعوايد الصعديه
متجولم يابنات هيا ودي شممالشلبية متجولم يابنات هيا ودي كرسى
الصعديه متجولم يابنات هيا على الس الصعديه

Quelle: Darwiš, Kāmal Afandī: Kitāb bida‘ al-fağār fi ḥaflat az-zār,
Kairo 1911, urn:nbn:de:gbv:3:5-86614, CC-BY-SA 4.0.

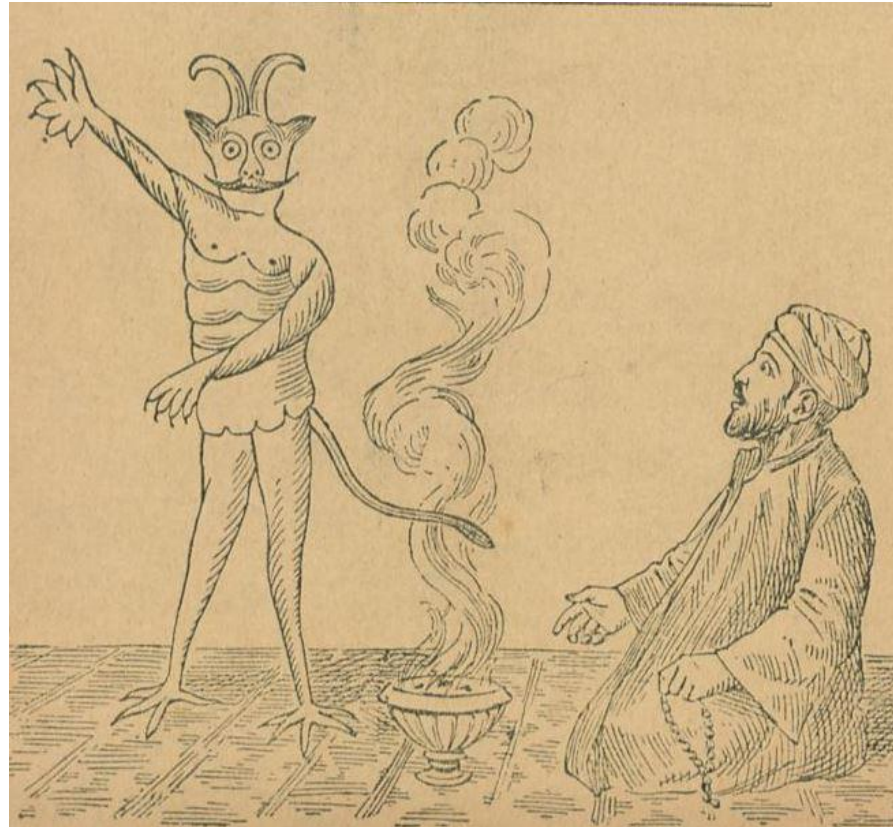
– ي ب ن ث ت ب – خ چ خ ح ج
ظ ط ض ص – ز ذ ر ز – ش ش ش س

Manuelle OCR mit Training

Manueller Workflow mit Erzeugung von GroundTruth-Daten



Vielen Dank, haben Sie noch Fragen?



Quelle: Darwiš, Kāmal Afandī: Kitāb bida' al-fağār fī ḥaflat az-zār,
Kairo 1911, urn:nbn:de:gbv:3:5-86614, CC-BY-SA 4.0.