

Automatischer Ingest aus Goobi Ein Praxiseinblick

Frankfurt, 01. Juni 2017



Was war Goobi nochmal?

- ▶ Workflowtool für Digitalisierungsprojekte
- ▶ Mittlerweile auch Präsentationstool
- ▶ Entwickelt in Java als Open Source
- ▶ Verbreitung zunächst vor allem in Deutschland
- ▶ Seit 2009 außerhalb Deutschlands in 56
Einrichtungen in 12 Ländern
- ▶ Jährliche User Meetings in Göttingen und London



Goobi - Workflowsteuerung

The screenshot shows the Goobi 3.0 web interface for managing production templates. The main heading is "Produktionsvorlagen". A table lists 15 templates with columns for ID, Titel, Status, Projekt, and Aktionen. A detailed view for ID 434 is shown below, listing 9 steps: 1. Preparation, 2. Scanning, 3. Image QA, 4. Automatic image conversion (JPEG2000), 5. Automatic JPEG 2000 validation, 6. Metadata indexing, 7. OCR, 8. NER, 9. Automatic export to viewer.

ID	Titel	Status	Projekt	Aktionen
541	> Beispielworkflow_OLR	<div style="width: 100%; height: 10px; background: linear-gradient(to right, green, yellow, red);"></div>	Sample_Project	[edit] [copy] [refresh]
343	> Example_Workflow	<div style="width: 100%; height: 10px; background: linear-gradient(to right, green, yellow, red);"></div>	Sample_Project	[edit] [copy] [refresh]
462	> Example_Workflow	<div style="width: 100%; height: 10px; background: linear-gradient(to right, green, yellow, red);"></div>	Wiener_Library	[edit] [copy] [refresh]
345	> Example_workflow_LayoutWizard	<div style="width: 100%; height: 10px; background: linear-gradient(to right, green, yellow, red);"></div>	Wiener_Library	[edit] [copy] [refresh]
410	> Example_workflow_LayoutWizard_Stanford	<div style="width: 100%; height: 10px; background: linear-gradient(to right, green, yellow, red);"></div>	Sample_Project	[edit] [copy] [refresh]
516	> Example_workflow_LayoutWizard_table_cropping	<div style="width: 100%; height: 10px; background: linear-gradient(to right, green, yellow, red);"></div>	Wiener_Library	[edit] [copy] [refresh]
434	> Example_workflow_NER_draft	<div style="width: 100%; height: 10px; background: linear-gradient(to right, green, yellow, red);"></div>	Sample_Project	[edit] [copy] [refresh]

Nr.	Titel	Status
1	Preparation	<div style="width: 100%; height: 10px; background-color: green;"></div>
2	Scanning	<div style="width: 100%; height: 10px; background-color: yellow;"></div>
3	Image QA	<div style="width: 100%; height: 10px; background-color: red;"></div>
4	Automatic image conversion (JPEG2000)	<div style="width: 100%; height: 10px; background-color: red;"></div>
5	Automatic JPEG 2000 validation	<div style="width: 100%; height: 10px; background-color: red;"></div>
6	Metadata indexing	<div style="width: 100%; height: 10px; background-color: red;"></div>
7	OCR	<div style="width: 100%; height: 10px; background-color: red;"></div>
8	NER	<div style="width: 100%; height: 10px; background-color: red;"></div>
9	Automatic export to viewer	<div style="width: 100%; height: 10px; background-color: red;"></div>

Workflows organisiert in Arbeitsschritten

Editor für Paginierung, Strukturdaten und Metadaten

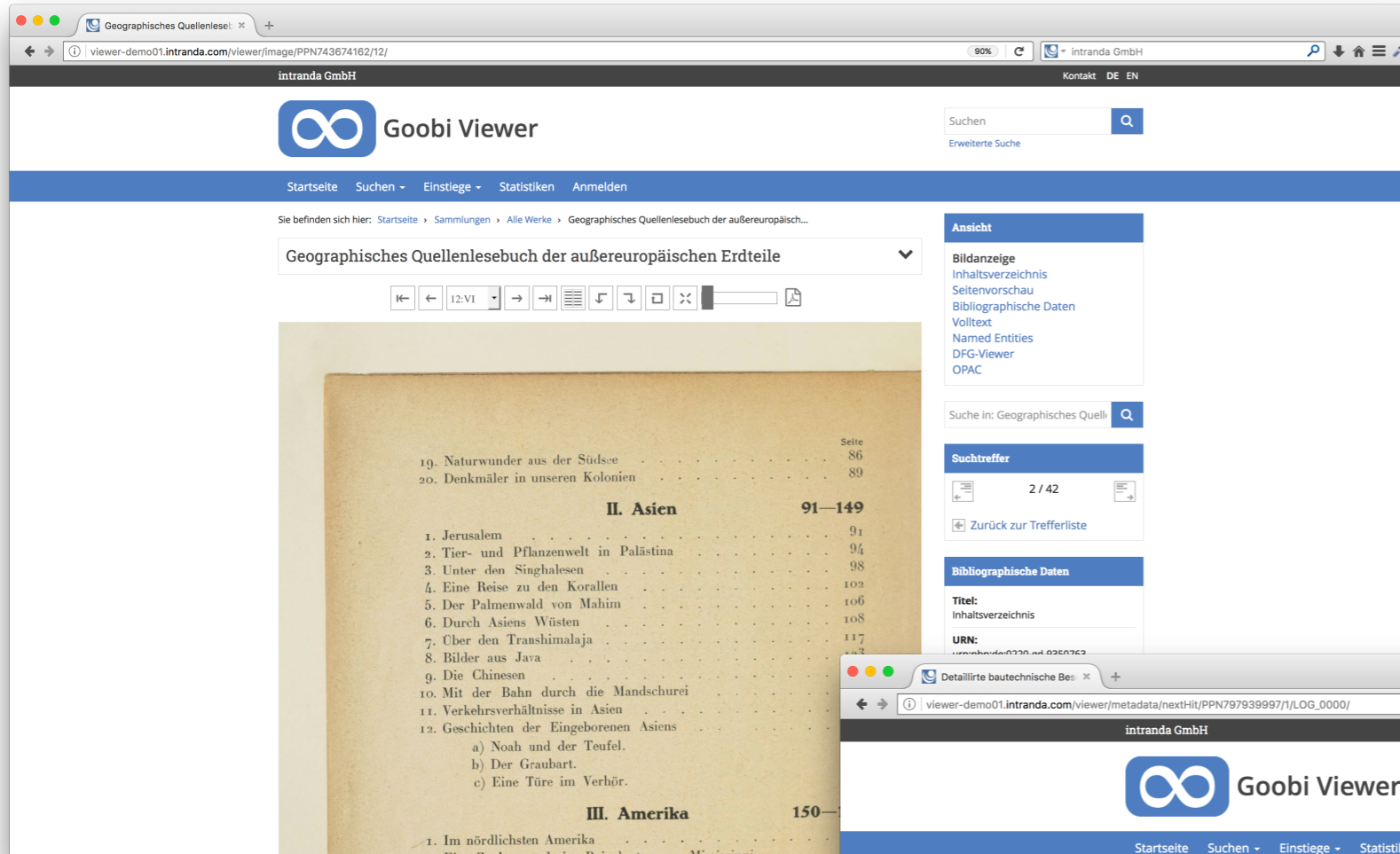
The screenshot shows the Goobi 2.2 metadata editor. It features a sidebar with a tree view containing "Monographie", "Einband", "Titelblatt", "Einleitung", and "Kapitel". The main area is divided into sections for "Personen" (with fields for Vorname, Nachname, Rolle, Normdaten) and "Metadaten" (with fields for Alle Autoren u. Beteiligte (orig.), Digitale Kollektion, Entstehungsjahr, Erscheinungsort, Haupttitel, Identifer (digital), Identifer (Vorlage), Schlagwort, and Sprache).

The screenshot shows a scanned manuscript page with a digital overlay. The page number is 18. The text includes a list of species: 9) Cordylus Monitor, Lacerta Monitor Gm., 10) — bimaculatus, Lacerta bimaculata Gmel., 11) — Dracaena, Lacerta Dracaena Gm., 12) — Caudiuerbera, Lacerta Caudiuerbera Gm. Below this is the heading "VII. CROCODYLVS." followed by a Latin description: "Corpus callosissimum, feguentis instructum. Linguae loco, valua elastica inter angulos maxillarum, fundo oris. Cauda anceps, in segmenta diuisa. a) Synonima. Crocodylus. Corpus quadrupes, feguentis instructum. Cauda anceps in segmenta diuisa. Pedum digiti aliquot vnguiculati, reliquis inermibus. Pedes posteriores quandoque natatorii. Lingua nulla. Gronouii Mus, ichth. T. II. P. 74. Crocodylus. Rostrum longissimum, hinc rictus omnium animalium maximus, lingua nulla, sed eius loco musculi maxillae inferioris admodum tu-".

01.06.2017

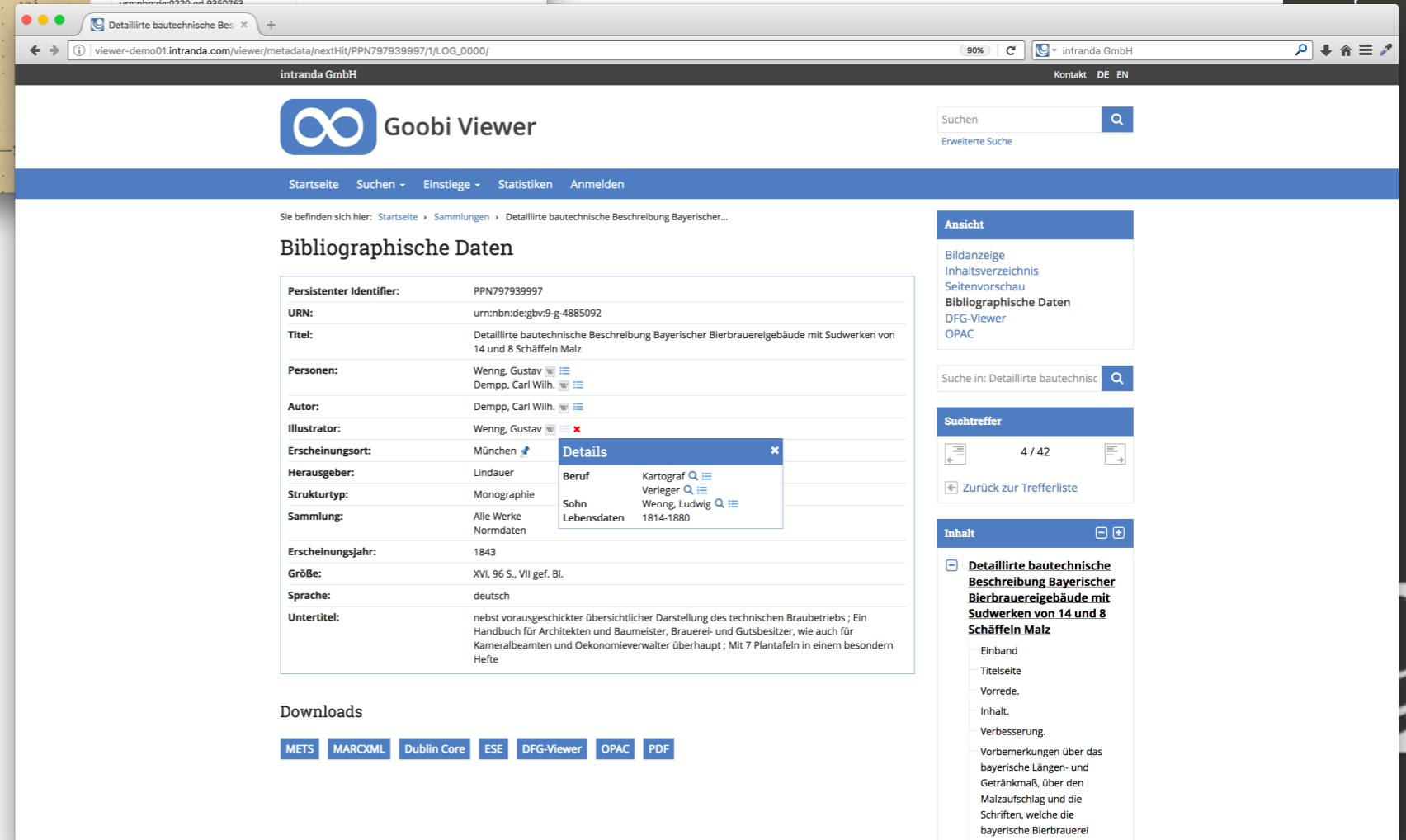
anda GmbH

Goobi viewer



IIF kompatible
Bildanzeigen

Metdatenanzeige inkl.
Integration von Live-
Daten aus der GND



01.06.2017

intra GmbH

Was war Goobi nochmal?

- ▶ Zumeist Retrodigitalisierungsprojekte in Bibliotheken
 - ▶ Arbeitsorganisation (Arbeitsschritte)
 - ▶ Metadatenerfassung
 - ▶ Automatisierung (Scripte etc.)
 - ▶ Gehostet oder selbst betrieben
 - ▶ Koordinierung der Datenmassen
 - ▶ Veröffentlichung der Ergebnisse
 - ▶ Datenübernahme aus anderen Systemen
 - ▶ Datenübergabe an andere Systeme



Was war Goobi nochmal?

Input / Import

Workflow step 1
Scannen

Workflow step 2
Qualitätskontrolle

Workflow step 3
Bildoptimierung

Workflow step 4
Metadatenerfassung

Workflow step 5
OCR

Workflow step 6
Katalogupdate

Output / Export

Daten kommen aus einer externen Quelle

Arbeitsschritte im Workflow:

- ▶ Schritte werden nacheinander ausgeführt
 - ▶ Schritte werden ausgeführt als Aufgaben
 - ▶ Aufgaben können manuell oder automatisch sein
 - ▶ Aufgaben können Tools, Scripts o.ä. aufrufen
 - ▶ Aufgaben können als Plugins über eigene Nutzeroberfläche verfügen
 - ▶ Benutzer haben eine Aufgabenliste
 - ▶ Abhängig von Workflow, Projekt und Nutzergruppe können Aufgaben ausgeführt werden
 - ▶ Validierungsplugins vermeiden Fehler
-

Daten werden zu einem anderen System geliefert



Zusammenspiel mit MyCoRe bisher

- ▶ Beispiel Thulb Jena
 - ▶ im Wesentlichen nur Workflow zur Koordinierung einiger Arbeitsschritte in Goobi
 - ▶ Metadatenerfassung innerhalb von MyCoRe
- ▶ Beispiel Braunschweig TU
 - ▶ Workflow vollständig via Goobi
 - ▶ Manueller Upload von PDF-Dateien nach MyCoRe



Zusammenspiel mit MyCore bisher

- ▶ Beispiel Rostock UB
 - ▶ Vollständige Workflowsteuerung durch Goobi
 - ▶ Metadatenerfassung in Goobi
 - ▶ Automatischer Datenabzug der internen METS-Datei sowie der Bilder nach MyCoRe über Eigenimplementierung „von außen“
 - ▶ Kommunikation mit Goobi via Web-API
 - ▶ Besonderheit außerdem: Bereitstellung der Digitalisate via IIF in MyCoRe und Anzeige in Landesportal auf Basis des Goobi viewers



Zusammenspiel mit MyCore bisher

01.06.2017

The screenshot shows a web browser window with the URL www.digitale-bibliothek-mv.de/viewer/image/83206579X/14/. The page header features the logo for Mecklenburg-Vorpommern with the slogan "MV tut gut." and navigation links for "KONTAKT", "Suchen", and "Anmelden".

The main content area displays the title "Zweyter Theil: Aus mehr als 750. Speisen bestehend, Nebst einer kleinen Hauß-Apotheck und andern Dem Frauenzimmer dienlichen raren Kunst-Stücken, Absonderlich Welchen an Säuberung Edelsteine, Perlen, und andern Sachen gelegen". Below the title is a viewer interface with navigation controls and a thumbnail of a manuscript page.

The manuscript page shows the following text:

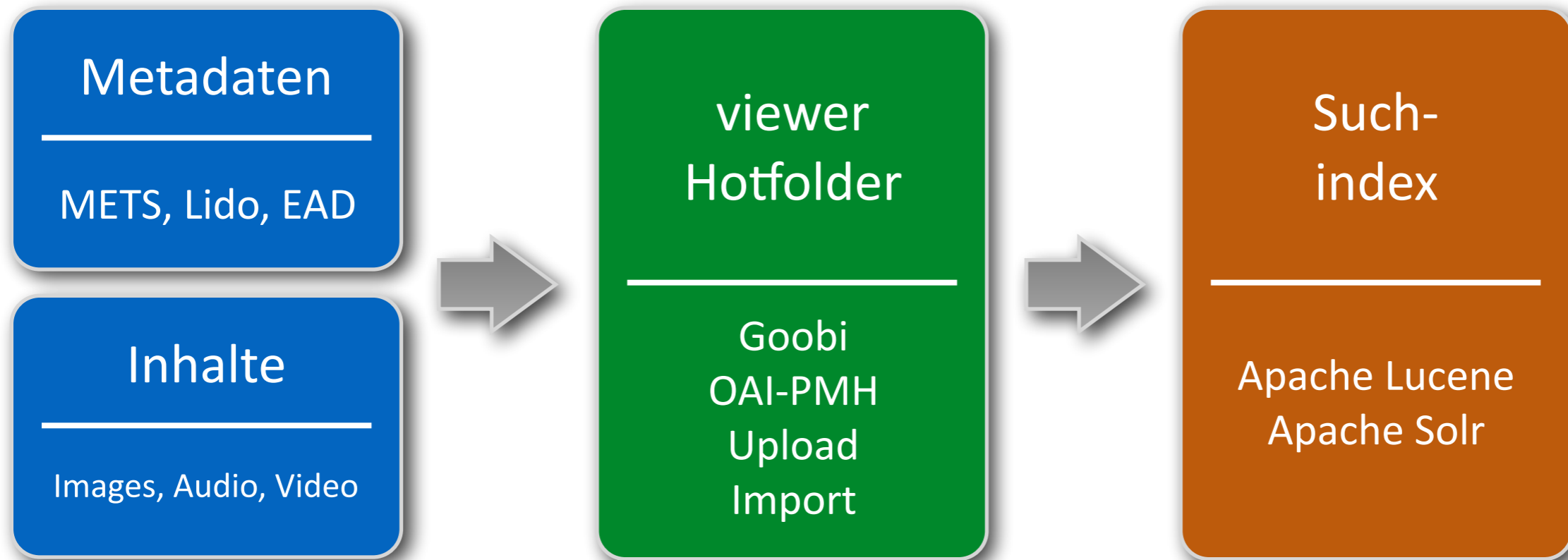
4 Gründliche Nachrichten
8. Reiß-Schnitten.
Roche im ersten Theil gezeigter massen einen Reiß-Brey, aber so dick, daß der Löffel darinnen stecken bleibt, giesse ihn in ein flaches Zinn, und laß ihn kalt werden; schneide aus dem Brey so grosse Stücke, als sonst die Weißbrodt-Schnitten seyn; mache von ein paar Aiern, Zucker, Zimmet und Rosentwasser ein kleines Taiglein, kehre die Breystücke darinnen um und backe sie im Schmalz.
9. Leber Schnitten.
Hacke die Leber mit Peterling, Zwibeln und eingeweichtem, aber wieder ausgedrucktem Weißbrodt, rühre es mit ein wenig Mutschel-

Verglichen dazu: Die Injests der anderen?

- ▶ Bisherige andere Injests in andere Systeme:
 - ▶ Goobi viewer
 - ▶ Preservica
 - ▶ Rosetta
 - ▶ Fedora



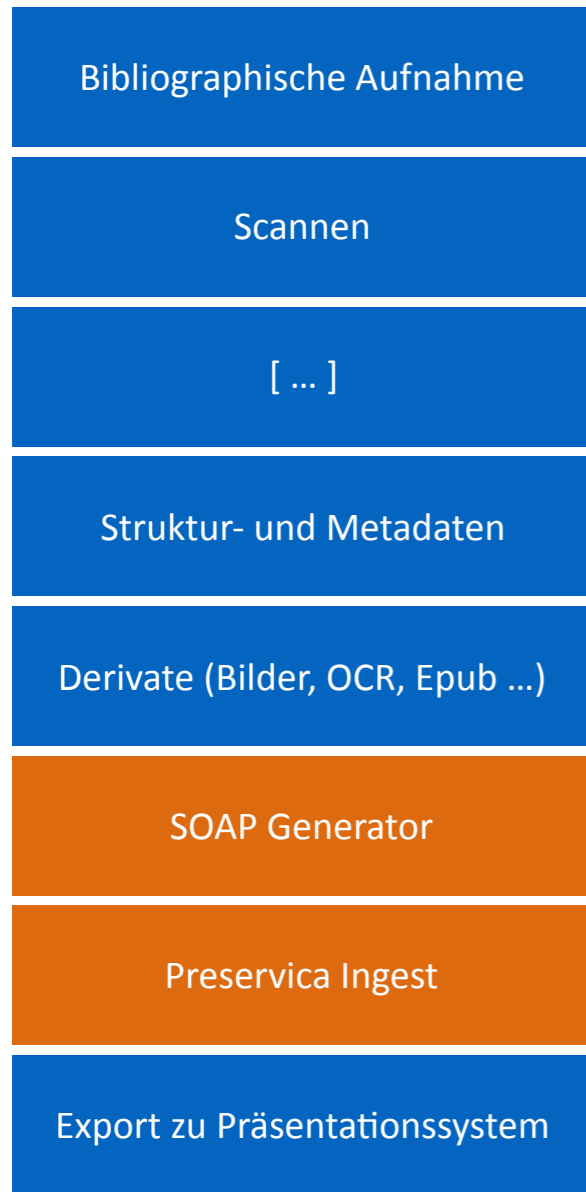
Ingest in Goobi viewer



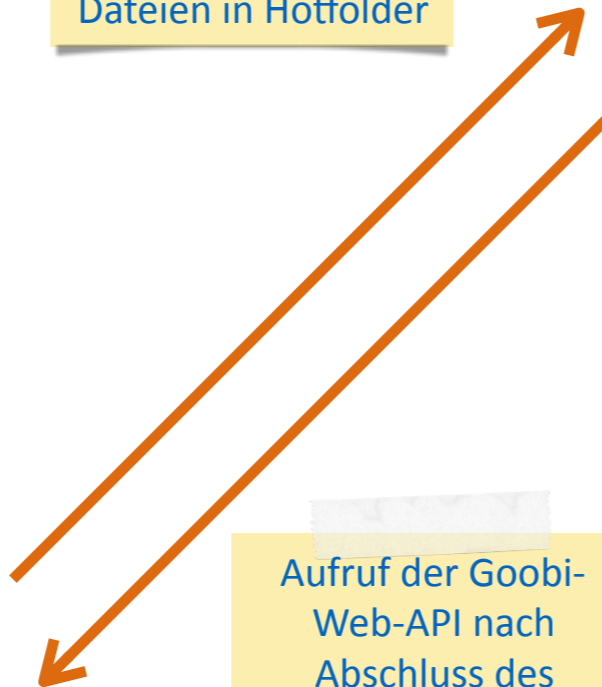
Arbeitsweise bisher über das Dateisystem



Goobi

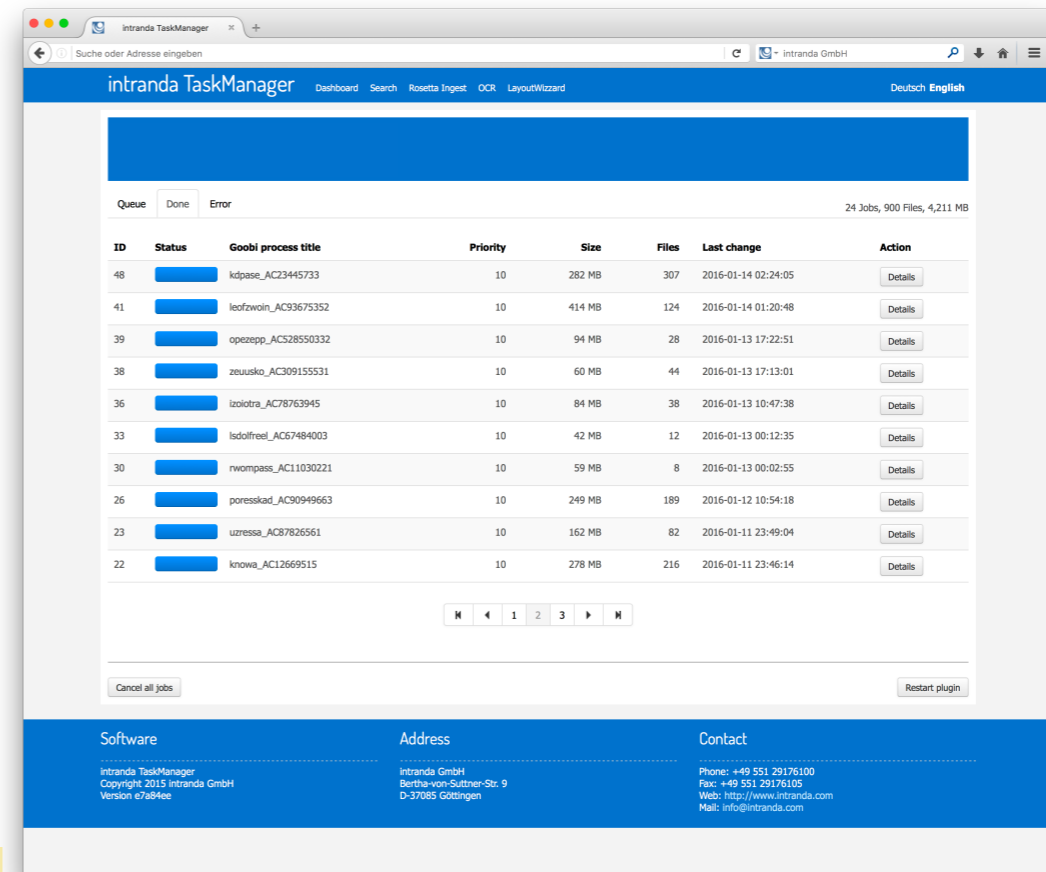


Separater Aufruf des SOAP-Generators und erst im Anschluß Bereitstellung der Dateien in Hotfolder



Aufruf der Goobi-Web-API nach Abschluss des erfolgreichen Ingests, Rückgabe einer xml-Datei mit Hash-Werten und URLs zur Ergänzung der METS-Datei

TaskManager

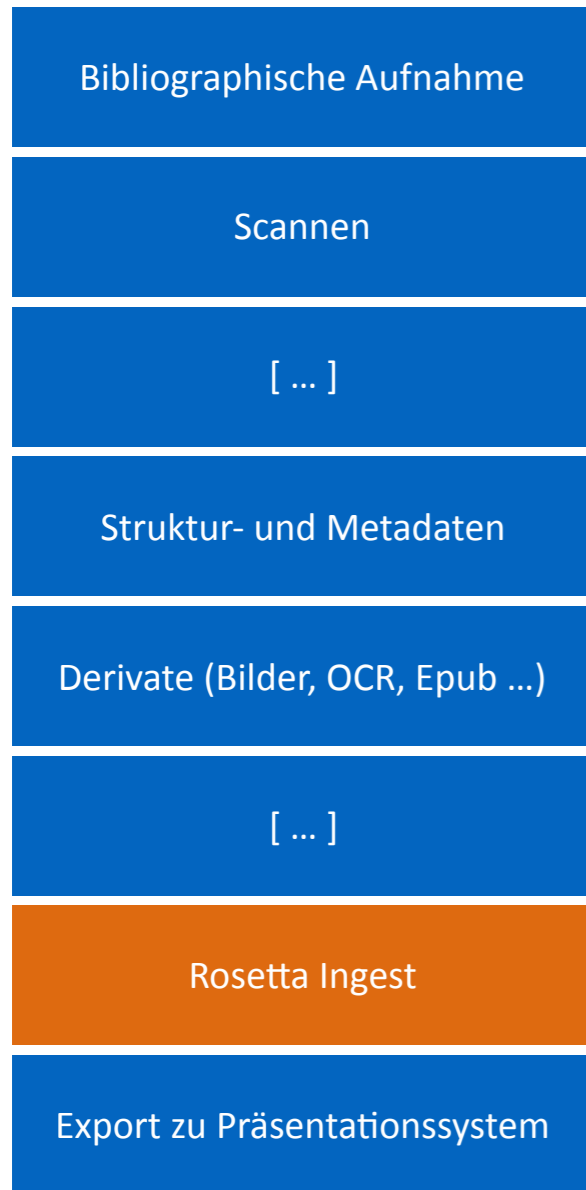


Bereitstellung der Inhalte ohne Metadaten; Ungewisse Wartezeit bei teilweise ungewissem Ergebnis

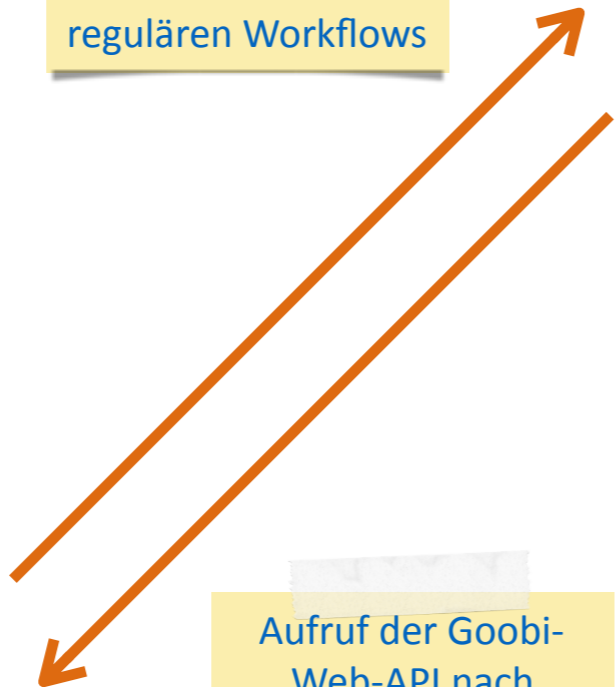
Preservica



Goobi

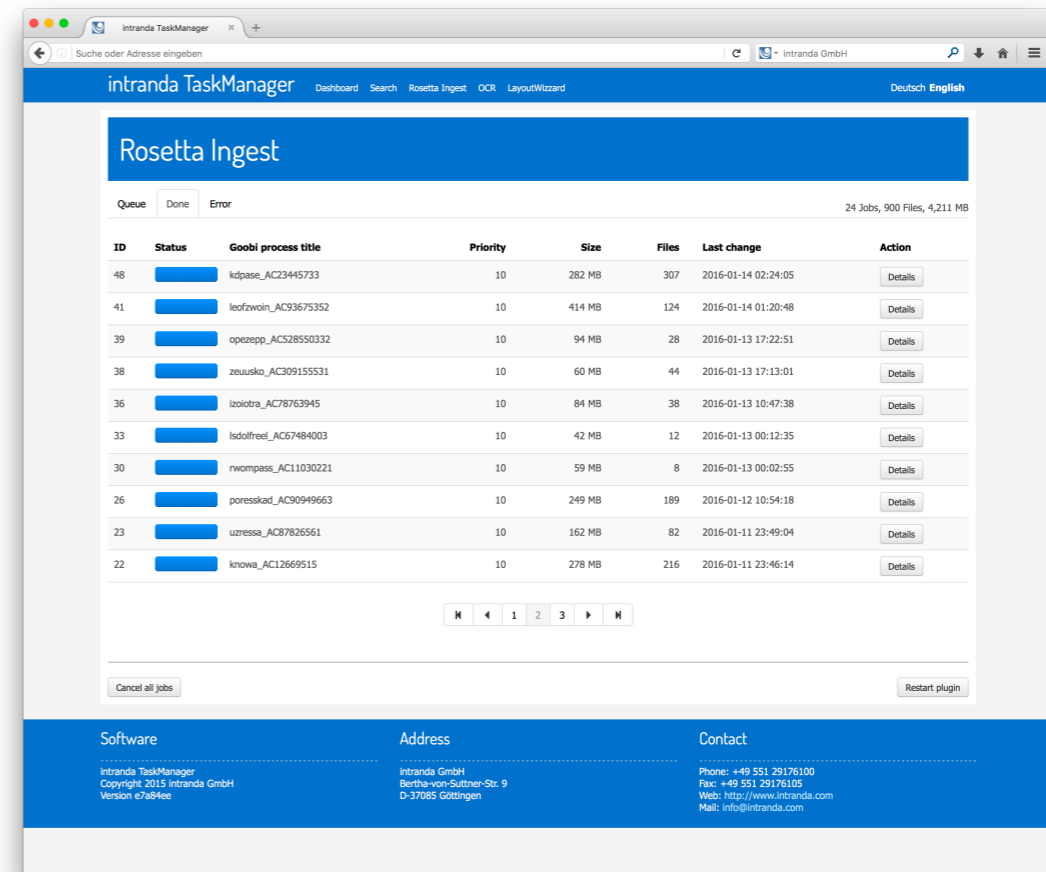


Aufruf des TaskManagers via TaskClient als Commandline Aufruf innerhalb des regulären Workflows



Aufruf der Goobi-Web-API nach Abschluss des erfolgreichen Ingests

TaskManager



Bereitstellung einer aufbereiteten METS-Datei und der zugehörigen Streams für den Ingest in Rosetta

Rosetta



Lessons learned: Der Status quo

- ▶ Große Datenmengen
 - ▶ Lange Transferzeit
 - ▶ Überwachung des Ingestfortschritts
 - ▶ Abschließende Statusüberprüfung
-
- ▶ Verlässliche Auskunft, ob die Daten vollständig und intakt angekommen sind



Goobi - 2.2

goobi01.fritz.box/goobi/ui/process_all.xhtml

intranدا GmbH

Steffen Hankiewicz

4 Aktive Benutzer

17:58:58
Sonntag, 11. September 2016

Vorgangsdetails

Dashboard > Vorgänge > Vorgangsdetails

Vorgang

Vorgangstitel:	demo_mycore_ocr
Projekt:	Tests
Erstellungsdatum:	16.08.2016
Regelsatz:	Standard
Laufzettel:	Standard
In Auswahlliste anzeigen:	<input type="checkbox"/>
Ist eine Vorlage:	<input type="checkbox"/>
ID:	38666
Batch:	

Vorgangsllog

[Nachricht hinzufügen](#)

Abfolge der Aufgaben

Nr.	Titel	Status	Aktionen
0	> Vorgang anlegen	■ < >	
1	> Einspielen der Images	■ < >	
2	> Tiff-Header erzeugen	■ < >	
3	> Erstellung der komprimierten Derivate	■ < >	
4	> Qualitätskontrolle	■ < >	
5	> Struktur-u. Metadaten	■ < >	
6	> OCR	■ < >	
7	> MyCore preparation	■ < >	
8	> MyCore ingest	■ < >	

[+ Aufgabe hinzufügen](#)

Physische Vorlagen

Autoren:	Knowles, Lees; Winkel, G. G.
Erscheinungsjahr:	1914

Aufgabendetails

4 Aktive Benutzer | 18:00:20 Sonntag, 11. September 2016

Dashboard > Vorgänge > Vorgangsdetails > Aufgabendetails

Details - demo_mycore_ocr

Titel *	MyCore preparation
Reihenfolge *	7
Priorität *	0
Metadaten	<input type="checkbox"/>
Images lesen	<input type="checkbox"/>
Images schreiben	<input type="checkbox"/>
Beim Abschließen verifizieren	<input type="checkbox"/>
Export DMS	<input checked="" type="checkbox"/>
Aufgabe überspringen	<input type="checkbox"/>
Automatische Aufgabe	<input checked="" type="checkbox"/>
Skript Schritt	<input type="checkbox"/>
Status *	Offen
Batch Schritt	<input type="checkbox"/>
Plugin für Arbeitsschritt	plugin_intranda_mycore_export
Validierungsplugin	Nicht ausgewählt
Plugin für Zeitverzögerung	<input type="checkbox"/>
Metadatenindex beim Abschließen aktualisieren	<input type="checkbox"/>

Löschen | **Abbrechen** | **Speichern**

Berechtigte Personen

- Administration

Start / / Stuart Lewis

Ein anderer Tag diesmal ohne Corps-Studen... 6 von 9 IIIF-Image Test

Strahlsundische Zeitung

Vorschau

Strukturübersicht

- Strahlsundische Zeitung
 - Ausgabe
 - Ausgabe
 - Ausgabe
 - Ausgabe

31% 0°

Dateien

Datei	Datum	Größe	Aktionen
00000001.tif	31.08.2016	2.41 MB	
00000002.tif	31.08.2016	2.66 MB	
00000003.tif	31.08.2016	2.78 MB	
00000004.tif	31.08.2016	2.44 MB	

auf die Merkliste | Aktionen

Zitieren

tweet | teilen | +1 | i

Zitierform:
Stuart Lewis.

Zitier-Link kopieren

Export

BibTeX, MODS, RIS, ISI, PICA, DC

Systeminformation

Publikationsstatus: eingereicht
Erstellt am: 19.05.2016 - 16:38:00
Erstellt von: administrator
Letzte Änderung: 31.08.2016 - 11:39:01
Zuletzt geändert von: administrator
MyCoRe ID: mir_mods_00000018
Version: 2
(Gesamte Versionsgeschichte)

STEFAN TRANKIEWICZ, INTRANDA GMBH



Lessons learned: Der Status quo

- ▶ Implementiert:
 - ▶ Erzeugung einer (einzigen) großen zip-Datei
 - ▶ Generierung von Dublin-Core Daten
 - ▶ Meldung der Dublin-Core Daten an MyCoRe
 - ▶ Hochladen der zip-Datei
- ▶ Derzeit noch problematisch für Werke größer als 10 GB:
 - ▶ (lange) warten
 - ▶ selbst manuell prüfen



Offene Fragen zur Vorgehensweise?

Fragen

