

ENTWICKLUNG PYTHON SCRIPT ZUM PRÜFEN UND KORRIGIEREN VON RDF FILES

LEISTUNGSNACHWEIS 2

VON

ANDRÉS BAUMELER

TELEFON: 076 443 04 71, E-MAIL: ANDRES@BAUMELER.DEV

ALTE RIEDIKERSTRASSE 5C, 8610 USTER

BETREUER

ISMAIL PRADA



Universität
Zürich ^{UZH}



UNIVERSITÄT ZÜRICH, PHILOSOPHISCHE FAKULTÄT /
ZENTRALBIBLIOTHEK ZÜRICH

CAS DATENMANAGEMENT UND INFORMATIONSTECHNOLOGIEN

Inhaltsverzeichnis

1	Einleitung	2
2	Hauptteil	3
2.1	Hintergrund	3
2.2	Problemstellung	4
2.3	Lösungsansatz	5
2.4	Ergebnis	5
2.5	Dokumentation	5
3	Schluss	6
3.1	Offene Fragen	6
3.2	Ausblick	6
3.3	Fazit und Reflexion	6
	Abkürzungen	7
	Literaturverzeichnis	8

Kapitel 1

Einleitung

Banken sind gesetzlich verpflichtet gewisse Dokumente aus ihren Geschäftstätigkeiten aufzubewahren. Um die aufzubewahrenden Dokumente zentral zu verwalten werden digitale Archivsysteme eingesetzt. Ein solches System ermöglicht es alle in einer Bank produzierten Dokumente im Überblick zu behalten und den Lebenszyklus der Dokument zu verwalten. Das Archivsystem ist in der Lage Dokumente und Metadaten aus einer Vielzahl an unterschiedlichen Quellen entgegenzunehmen. Während die Akzeptierten Formate zwar klar definiert sind, ergeben sich beim Betrieb eines solchen Archivsystems immer wieder Herausforderungen mit Dokumenten welche die Anforderungen an den Archivierungsprozess nicht erfüllen. In dieser Arbeit wurde ein Python Script entwickelt welches den Betreiber einer Archivlösung dabei unterstützen kann solche Dokumente zu prüfen.

Konkret behandelt das in dieser Arbeit entwickelte Python Script Files welche Metadaten im Ressource Description Format (RDF) enthalten. Die Files sind in der eXtensible Markup Language (XML) strukturiert. In dieser Arbeit wird nur das Archivsystem Hyper Suite 5 (HS5) der Firma IMTF¹ betrachtet.

¹<https://imtf.com/>

Kapitel 2

Hauptteil

2.1 Hintergrund

In einer Bank gibt es in der Regel eine Vielzahl an Systemen welche Aufbewahrungspflichtige Dokumente produzieren. Die Aufbewahrung der Dokumente wird aber zentral in einem Archivsystem durchgeführt. Dadurch ergibt sich der Bedarf für Schnittstellen von produzierenden Systemen zum Archivsystem. HS5 bietet dazu die Möglichkeit Dokumente über eine File-basierte Schnittstelle als RDF und PDF Paar entgegen zunehmen. Das PDF enthält das eigentliche Dokument während das RDF File die Metadaten zum Dokument enthält. Der Import dieser File Paare erfolgt aus einem überwachten Ordner auf dem Filesystem des Archivservers. Sobald dort ein File Paar abgelegt wird, prüft das Archivsystem die Dokumente. Erfüllt das angelieferte File Paar die Anforderungen für die Archivierung wird das PDF archiviert. Bei der Archivierung werden die benötigten Metadaten aus dem RDF gelesen und in der Archiv Datenbank gespeichert. Das PDF wird auf eine Write Once Read Many (WORM) Storage geschrieben um sicherzustellen, dass das Dokument selbst nicht mehr verändert werden kann. Das originale RDF File wird einige Tage nach erfolgreicher Archivierung des PDFs gelöscht.

Da die Dokumente auf eine WORM Storage Lösung geschrieben werden, ist ein Löschen bis zum Ablauf der Aufbewahrungsfrist nicht möglich. Aus diesem Grund müssen die gelieferten Metadaten nicht nur syntaktisch

sondern auch inhaltlich korrekt sein. Das Archivsystem führt aus diesem Grund eine Datenbank mit Stammdaten der Bank. Darin enthalten sind etwa Kunden- und Kontonummern. Vor der Archivierung werden die im RDF angelieferten Dokumente gegen diese Datenbank geprüft. Wird beispielsweise eine Kontonummer in den Metadaten angegeben welche dem Archiv nicht bekannt ist, wird das Dokument nicht archiviert. Für Dokumente welche die Anforderung an die Archivierung nicht erfüllen, wird das RDF und PDF Paar zusammen mit einem kurzen Fehlerbeschrieb in ein Verzeichnis geschrieben. Dieses sogenannte failed Verzeichnis wird überwacht und die dort abgelegten Dokumente regelmässig geprüft.

2.2 Problemstellung

Die Prüfung von Dokumenten im failed Verzeichnis erfolgt von Hand und nimmt pro Dokument einiges an Zeit in Anspruch. Gewisse prüfungen können aber automatisiert werden. Wenn ein Dokument im failed Verzeicunis liegt kann drei Hauptgründe:

- Metadaten sind nicht vollständig (z.B. keine Kontonummer angegeben)
- Metadaten sind im falschen Format (z.B. Kontonr. hat zu wenig Stellen)
- Metadaten sind dem Archiv nicht bekannt (z.B. Kontonummer ist nicht in Archiv Datenbank)

Nicht vollständige Metadaten bzw. Metadaten im falschen Format können nicht automatisch geprüft werden und müssen weiterhin manuell abgeklärt werden. Grund hierfür ist, dass die Verantwortlichkeit über die Metadaten beim anliefernden System liegt.

Der dritte Punkt kann automatisiert werden. Aufgrund der asynchronen verarbeitung der Dokumente, kommt es immer wieder vor, dass Dokumente geliefert werden, bevor die Stammdaten dafür im Archiv sind. Wurde verifiziert, dass die Stammdaten im Archiv sind, kann der Archivierungsprozess für das Dokument nochmals gestartet werden. Die Prüfung ob Metadaten mittlerweile in der Archiv Datenbank bekannt sind, war der Hauptpunkt welcher durch das Script gelöst werden sollte.

2.3 Lösungsansatz

Als Lösung für das beschriebene Problem sollte ein Script entwickelt werden, welches die im RDF File gelieferten Metadaten gegen die Archivdatenbank prüft. Sind alle Daten dem Archiv bekannt, soll das Script die File Paare auch gleich in den korrekten Input Ordner verschieben können. Als Unterstützung sollte das Script auch anzeigen, wenn die Metadatenfiles im falschen Foramt sind (kein gültiges XML). Zur Umsetzung wurde Python gewählt. Grund für die Wahl von Python ist, dass eine entsprechende Laufzeitumgebung auf dem Zielsystem bereits vorhanden ist und die Sprache für den Umgang mit Dokumenten im XML Format sowie bei der Darstellung gegenüber Alternativen wie Shell Script einige Annehmlichkeiten bietet. Ziel war es ein Script zu entwickeln welches die folgenden Anforderungen erfüllt:

1. Automatisches Prüfen von Metadaten gegen Archiv Datenbank
2. Easy of Use: Als terminal UI Applikation oder via Commandline
3. Erweiterbarkeit und Wartbarkeit

2.4 Ergebnis

Das Ergebnis besteht aus einem Python Script

2.5 Dokumentation

Die aktuellste technische Dokumentation sowie eine Anleitung zum Aufsetzen einer Referenzumgebung ist in den Files *start.md* sowie *README.md* enthalten.

Der Source Code für das Script ist auf GitHub verfügbar.

Kapitel 3

Schluss

3.1 Offene Fragen

3.2 Ausblick

3.3 Fazit und Reflexion

Abkürzungen

HS5	Hyper Suite 5
RDF	Ressource Description Format
WORM	Write Once Read Many
XML	eXtensible Markup Language

Literaturverzeichnis