



Ein Open Source Package
Analyse von komplex strukturierten Dokumenten

01.03.2022 – Dr. Janis Meyer

Beleg

Extraktion aus Rechnungen

Entitätenerkennung, Gruppierung von Positionen

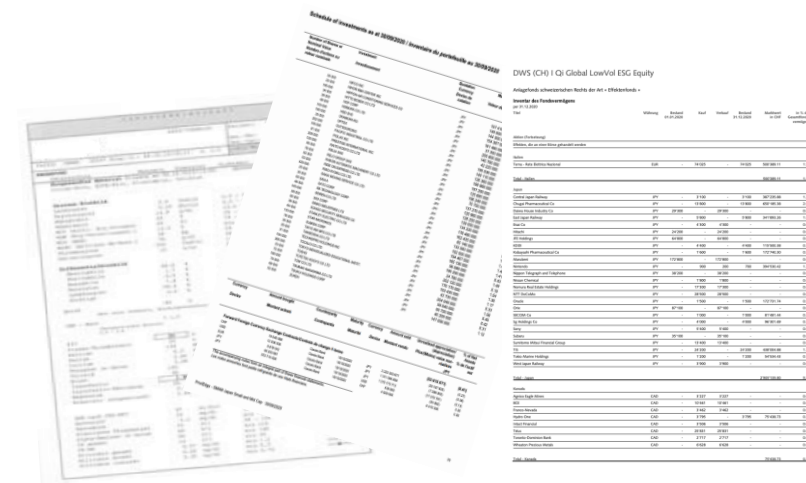
1 x	Nasi Campur Bali	75,000
1 x	Bbk Bengil Nasi	125,000
1 x	MilkShake Starwb	37,000
1 x	Ice Lemon Tea	24,000
1 x	Nasi Ayam Dewata	70,000
3 x	Free Ice Tea	0
1 x	Organic Green Sa	65,000
1 x	Ice Tea	18,000
1 x	Ice Orange	29,000
1 x	Ayam Suir Bali	85,000
2 x	Tahu Goreng	36,000
2 x	Tempe Goreng	36,000
1 x	Tahu Telor Asin	40,000
1 x	Nasi Goreng Samb	70,000
3 x	Bbk Panggang Sam	366,000
1 x	Ayam Sambal Hija	92,000
2 x	Hot Tea	44,000
1 x	Ice Kopi	32,000
1 x	Tahu Telor Asin	40,000
1 x	Free Ice Tea	0
1 x	Bebek Street	44,000
1 x	Ice Tea Tawar	18,000
Sub-Total		1,346,000
Service		100,950
PBI		144,695
Rounding		-45
Grand Total		1,591,600

1 x	Nasi Campur Bali	75,000
1 x	Bbk Bengil Nasi	125,000
1 x	MilkShake Starwb	37,000
1 x	Ice Lemon Tea	24,000
1 x	Nasi Ayam Dewata	70,000
3 x	Free Ice Tea	0
1 x	Organic Green Sa	65,000
1 x	Ice Tea	18,000
1 x	Ice Orange	29,000
1 x	Ayam Suir Bali	85,000
2 x	Tahu Goreng	36,000
2 x	Tempe Goreng	36,000
1 x	Tahu Telor Asin	40,000
1 x	Nasi Goreng Samb	70,000
3 x	Bbk Panggang Sam	366,000
1 x	Ayam Sambal Hija	92,000
2 x	Hot Tea	44,000
1 x	Ice Kopi	32,000
1 x	Tahu Telor Asin	40,000
1 x	Free Ice Tea	0
1 x	Bebek Street	44,000
1 x	Ice Tea Tawar	18,000
Sub-Total		1,346,000
Service		100,950
PBI		144,695
Rounding		-45
Grand Total		1,591,600

results: {
 line_items:
 [
 { quantity: 2,
 menu_item: "burger",
 price: 13.50
 }, ...
],
 sub_total_net: 150,
 service: 5,
 tax: 28,
 total: 183
}

Positionen aus Fondsberichten

Document to Database: Extraktion von Fondsdaten aus Jahresberichten



- Document Layout Analysis
 - Tabellenerkennung
 - Textextraktion
 - Dokument parsen



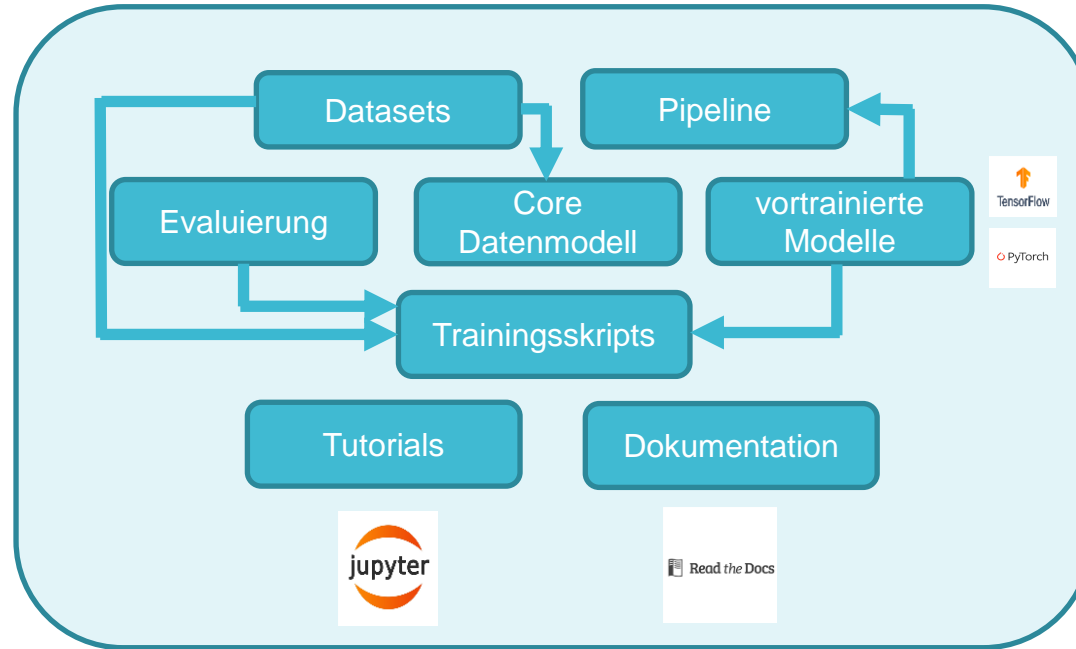
- Tabellenklassifikation
- Tabellennormalisierung

Datum	Fonds	Asset	Marktwert
31.12.21	HV Funds	VW AG	100

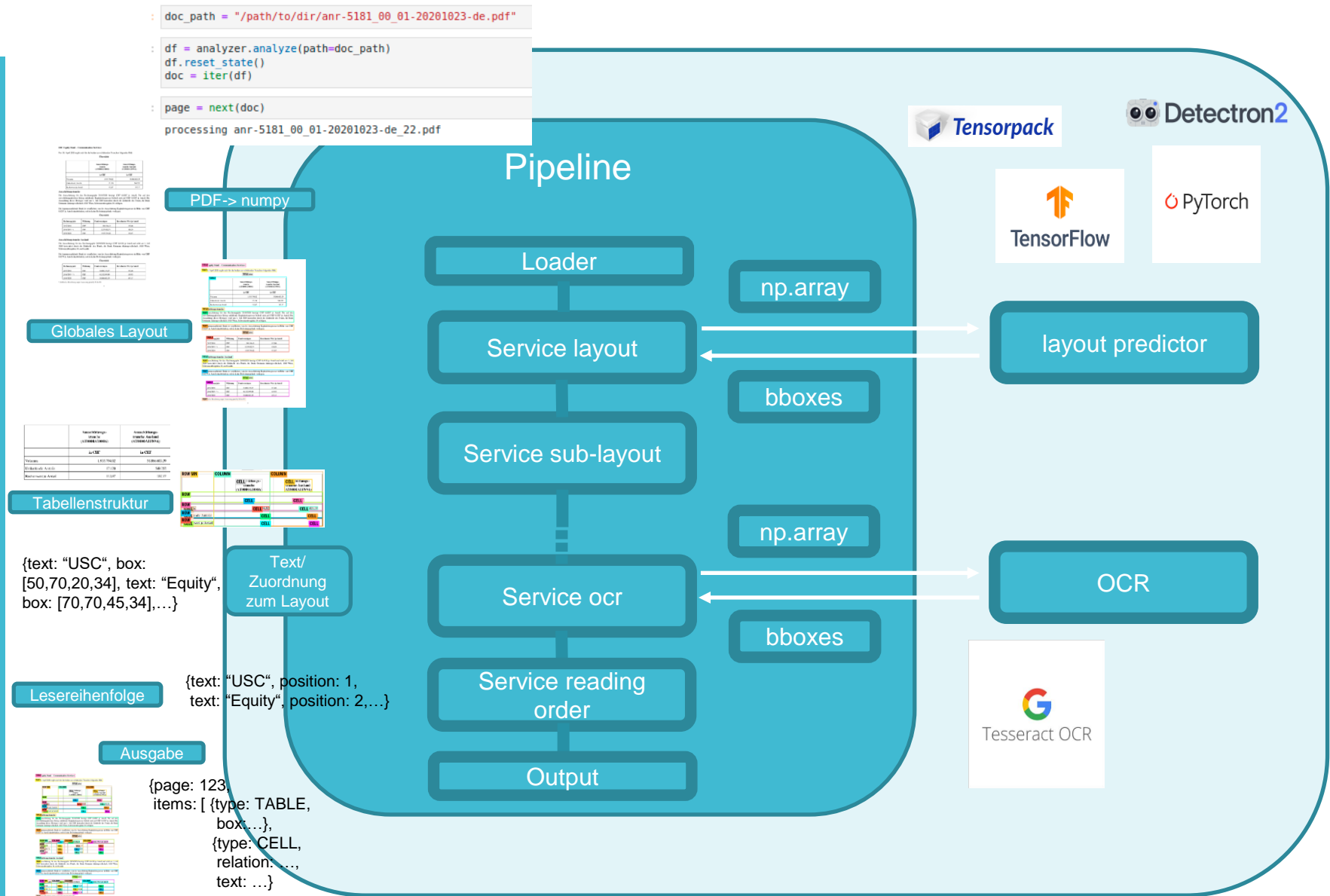
Open Source Tool-Box

Ein Python Package für Document Layout Analysis und Anwendungen

 deepdoctection



Pipelines



Fine-tuning

[illegible]

deepdoctection's pre-trained model

Nach Fine-tuning mit Domain-Dataset (ca. 400 Seiten)

Ausblick

Planung und Ideen

Schnittstellen zu NLP-Packages (z.B. SpaCy)

PDF- Miner

Weitere OCR Tools

Pre-Training/ Integration neuer Modelle

Nützliche Features, mehr Tutorials, ...

Anregungen und PRs willkommen !

<https://github.com/deepdoctection/deepdoctection>