



BERD
@NFDI

Focused Tutorial on Capturing, Enriching, Disseminating Research Data Objects

Extracting research data from historical documents with eScriptorium and Python

Jan Kamlah, Thomas Schmidt and Renat Shigapov
Universitätsbibliothek Mannheim

24.11.2022

1. Einführung & Digitalisierung
2. Layoutsegmentierung und OCR via eScriptorium
3. Datenextraktion und -strukturierung via Python
4. Fazit

1. Einführung & Digitalisierung

1. Einführung & Digitalisierung



- Anfrage zur Extraktion von Forschungsdaten aus "**Die Maschinen-Industrie im Deutschen Reich von 1937**" durch Prof. Jochen Streb (Professur für Wirtschaftsgeschichte @ Uni Mannheim)
- Kooperationsprojekt von BERD@NFDI und OCR-D Modulprojekt "Werkspezifisches Training" @ UB Mannheim
- Digitalisierung (654 Seiten) durch Digitalisierungswerkstatt der UB Mannheim

Phönix — Pincuss

Phönix-Werk G. m. b. H., Spezialfabrik moderner Trocken- Apparate, Meerane (Sa.).

Fernruf: 2424. **Drahtanschrift:** phönixwerk
Gründung: 1907.
Fabrikationsprogramm: Trockenapparate; Holzbearbeitungs-Maschinen.
Kapital: RM 17.500.—
Anteilhaber: G. E. Nestmann, Meerane (100%).
Geschäftsführer: Obering. A. Wackermann.
Prokurist: J. Frenzel.
Bankverbindungen: Meeraner Bank A.-G., Reichsbank, Meerane.
Postscheck-Konto: 115 485 Leipzig.
Geschäftsjahr: 1./1.—31./12.

Grundbesitz: 3700 qm, davon 1403 qm bebaut.
Anlagen: Fabrikationsräume mit Montagehalle, elektr. Schweißanlage; kaufm. u. techn. Büro.
Eigene Vertretung in Berlin: Willy Böckel, W 30, Martin-Luther-Str. 12.
Besondere Angaben: Gegründet als Spezialfirma moderner Trocknungsanlagen. Hergestellt wurden Schlangensammel-Trockenapparate System Otto und Getreidetrockner nach dem Zellen-system neben anderen allgemeinen Trocknungsanlagen. 1921 erfolgte neben dem Trocknerbau die Aufnahme der Fabrikation von Holzbearbeitungsmaschinen. 1932 übernahm das Werk zusätzlich die Herstellung von Bewoilt-Apparaturen für die Papierindustrie nach den Patenten Dr. Bruno Wiegler, Berlin.

Piccolo-Automaten G. m. b. H., Berlin W 35, Kurfürstenstraße 146.

Fernruf: 21 20 95. **Drahtanschrift:** piccoloautomat
Gründung: 1932.
Fabrikationsprogramm: Schokoladen-Verkaufsapparate (Tischautomaten, Kugelschapparate).
Kapital: RM 20.000.—
Anteilhaber: Dr. Richard Schönthal, Wien.
Geschäftsführer: Erwin Hantsch, Bln.-Steglitz.
Bankverbindung: Deutsche Bank u. Disconto-Ges., Berlin.
Postscheck-Konto: 8012 Berlin.
Geschäftsjahr: 1./1.—31./12.
Anlagen: Mechanische Werkstätten.
Gefolgschaft: 13 Arbeiter u. 6 Angestellte.

F. Piechatzek, Kran- u. Aufzug-Werke, Berlin N 65, Seestraße 51-56.

Fernruf: 46 43 11. **Drahtanschrift:** lüderszug
Gründung: 1885.
Fabrikationsprogramm: Krane u. Aufzüge, Hebezeuge u. Hebe-maschinen (Flaschenzüge, Laufkatzen, Winden, Elektro-Flaschenzüge).
Geschäftsleiter: Richard, Martin u. Paul Piechatzek.
Prokuristen: Paul Gräning, Alfred Knop, Otto Kuhwald.
Bankverbindung: Reichskredit-Gesellschaft A.-G., Berlin.
Postscheck-Konto: 4847 Berlin.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 9500 qm, davon 4500 qm bebaut.
Anlagen: Maschinenbau-Werkstätten (Dreherei, Schleiferei, Fräseerei, Presserei u. Schlosserei).
Eigene Vertretungen: Im Ausland.
Besondere Angaben: Der Export erstreckt sich auf alle Weltteile.
Gefolgschaft: 350 Mitglieder.

Otto Pieron
siehe Maschinenfabrik

Paul Pietzschmann Wasserwerksbau, Berlin-Spandau, Schönwalder Str. 34.

Fernruf: 37 68 71.
Gründung: 1903.
Fabrikationsprogramm: Entwurf und Bau von Trink- und Gebrauchswasserwerken jeder Größe, Enteisungsanlagen einschl. der erforderlichen Antriebsanlagen (Dampf — Diesel — Wasserkraft, Hoch- und Niederspannung).
Inhaber: Ing. Paul Pietzschmann.
Bankverbindungen: Dresdner Bank, Spandauer Bank, Spandau.

Besondere Angaben: Der Firmeninhaber beschäftigt sich hauptsächlich mit der Erstellung halb- u. vollautomatisch arbeitender Wasserwerke u. besitzt hierüber große u. langjährige Erfahrungen.

Anton Piller Maschinenfabrik, Osterode (Harz), Abgunst 24.

Fernruf: 211. **Drahtanschrift:** apo
Gründung: 1909.
Fabrikationsprogramm: Ventilatoren für Heizungs-, Lüftungs-, Absauge u. sonstige Zwecke.
Bankverbindung: Reichsbank, Städt. Sparkasse, Osterode a. H.
Postscheck-Konto: 40 278 Hannover. ×

Pilot, G. m. b. H.
siehe Maschinenfabrik

Friedrich Piltz & Sohn K.-G., Heidenheim a. d. Brenz, Friedrichstr. 9.

Fernruf: 637. **Drahtanschrift:** piltz
Gründung: 1863.
Fabrikationsprogramm: Genauigkeitswerkzeuge für Herstellung u. Kontrolle von Gewinden; Gewindeschleifeinrichtungen; Drehbank-Schleifapparate.
Gesellschafter: Otto u. Walther Piltz.
Geschäftsführer: Die Gesellschafter.
Bankverbindungen: Reichsbank, Deutsche Bank u. Disconto-Ges., Heidenheim.
Postscheck-Konto: 2178 Stuttgart.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 2600 qm, davon 2000 qm bebaut; gepachtet sind 1200 qm mit 700 qm bebauter Fläche.
Anlagen: Verwaltungs- u. Fabrikationsräume in Heidenheim u. München.
Besondere Angaben: In München befindet sich ein Zweigwerk der Gesellschaft.

Eduard Pincuss Armaturenfabrik, Sanitäre Einrichtungen, Berlin O 17, Gr. Frankfurter Str. 13.

Fernruf: 59 13 18. **Drahtanschrift:** epal
Gründung: 1859.
Fabrikationsprogramm: Wasserleitungs-Armaturen.
Inhaber: Arthur Landsberger, Bln.-Charlottenburg; Ernst Reichenbach, Bln.-Grünwald.
Prokuristen: Paul Derpsch, Frieda Thierschmann.

462

Forschungsfrage:

- Identifizierung von Rechtsformwechseln der gelisteten Unternehmen
- Hierfür Extraktion des **Firmennamens** sowie der **Rechtsformen** notwendig

**Piccolo-Automaten G. m. b. H.,
Berlin W 35, Kurfürstenstraße 146.**
Fernruf: 21 20 95. **Drahtanschrift:** piccoloautomat
Gründung: 1932.
Fabrikationsprogramm: Schokoladen-Verkaufsapparate
(Tischautomaten, Kugelstechapparate).
Kapital: RM 20 000.—.
Anteileigner: Dr. Richard Schönthal, Wien.
Geschäftsführer: Erwin Hantsch, Bln.-Steglitz.
Bankverbindung: Deutsche Bank u. Disconto - Ges.,
Berlin.
Postscheck-Konto: 8012 Berlin.
Geschäftsjahr: 1./1.—31./12.
Anlagen: Mechanische Werkstätten.
Gefolgschaft: 13 Arbeiter u. 6 Angestellte.

Forschungsfrage:

- Identifizierung von Rechtsformwechseln der gelisteten Unternehmen
- Hierfür Extraktion des **Firmennamens** sowie der **Rechtsformen** notwendig

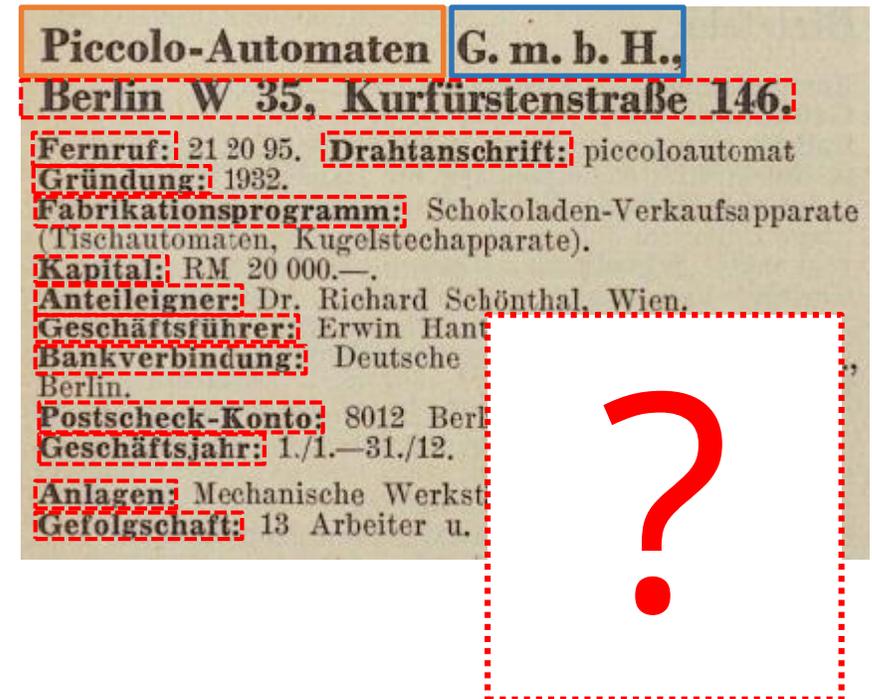


Piccolo-Automaten G. m. b. H.
Berlin W 35, Kurfürstenstraße 146.
Fernruf: 21 20 95. Drahtanschrift: piccoloautomat
Gründung: 1932.
Fabrikationsprogramm: Schokoladen-Verkaufsapparate
(Tischautomaten, Kugelstechapparate).
Kapital: RM 20 000.—.
Anteileigner: Dr. Richard Schönthal, Wien.
Geschäftsführer: Erwin Hantsch, Bln.-Steglitz.
Bankverbindung: Deutsche Bank u. Disconto - Ges.,
Berlin.
Postscheck-Konto: 8012 Berlin.
Geschäftsjahr: 1./1.—31./12.
Anlagen: Mechanische Werkstätten.
Gefolgschaft: 13 Arbeiter u. 6 Angestellte.

The image shows a historical document snippet with two boxes highlighting specific information: an orange box around 'Piccolo-Automaten' and a blue box around 'G. m. b. H.'. Dotted lines connect these boxes to the text in the list above. A blue dotted line also points from the 'Rechtsformen' box to the 'G. m. b. H.' box in the document snippet.

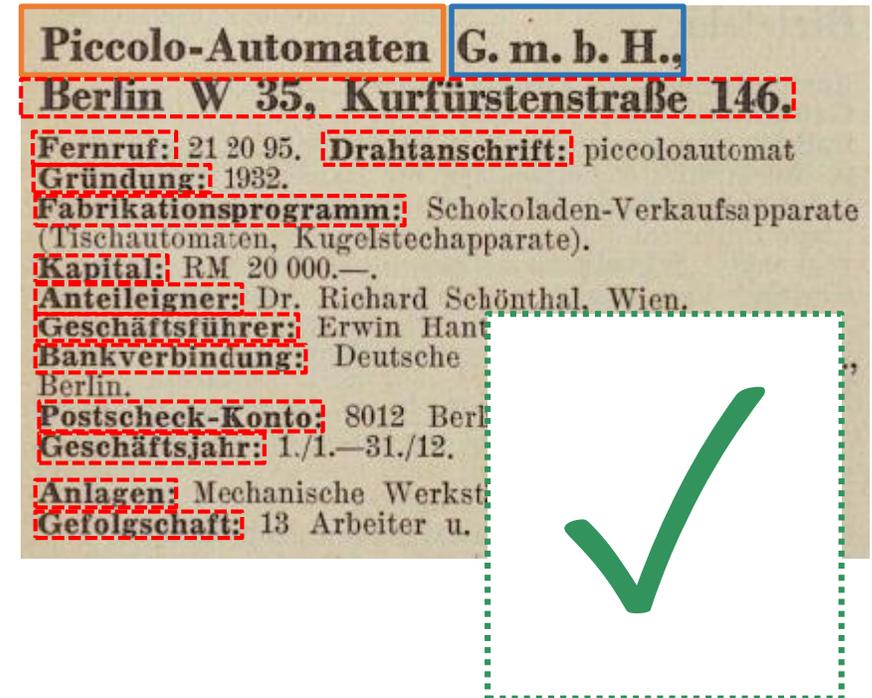
Forschungsfrage:

- Identifizierung von Rechtsformwechseln der gelisteten Unternehmen
- Hierfür Extraktion des **Firmennamens** sowie der **Rechtsformen** notwendig



Projektansatz und -ziel:

- Daten über die Forschungsfrage hinaus extrahieren und strukturieren (realistischer Mehraufwand)
- Wiederverwendbarkeit der Daten für andere Forschungszwecke sicherstellen
- Testworkflow, um Erfahrungen für vergleichbare Projekte zu sammeln



Projektansatz und -ziel:

JPG

**Piccolo-Automaten G. m. b. H.,
Berlin W 35, Kurfürstenstraße 146.**
Fernruf: 21 20 95. **Drahtanschrift:** piccoloautomat
Gründung: 1932.
Fabrikationsprogramm: Schokoladen-Verkaufsapparate
(Tischautomaten, Kugelstechapparate).
Kapital: RM 20 000.—.
Anteileigner: Dr. Richard Schönthal, Wien.
Geschäftsführer: Erwin Hantsch, Bln.-Steglitz.
Bankverbindung: Deutsche Bank u. Disconto - Ges.,
Berlin.
Postscheck-Konto: 8012 Berlin.
Geschäftsjahr: 1./1.—31./12.
Anlagen: Mechanische Werkstätten.
Gefolgschaft: 13 Arbeiter u. 6 Angestellte.

Transformation

Strukturierte Forschungsdaten

Firmenname	Piccolo-Automaten
Rechtsform	G.m.b.H.
Sitz	Berlin W 35 Kurfürstenstraße 146
Fernruf	212095
Drahtanschrift	piccoloautomat
Gründung	1932
Produkt	Schokoladen- Verkaufsapparate (Tischautomaten, Kugelstechapparate)

...

1. Einführung & Digitalisierung



Herausforderungen (OCR):

- Hohe Anforderungen an Qualität für Forschungsdaten
- Aufwand für OCR Nachtraining schwer abzuschätzen
- Layout (zweispaltig → Reading Order)

Phönix — Pincuss

Phönix-Werk G. m. b. H., Spezialfabrik moderner Trocken- Apparate, Meerane (Sa.).

Fernruf: 2424. **Drahtanschrift:** phönixwerk
Gründung: 1907.
Fabrikationsprogramm: Trockenapparate; Holzbearbeitungs-Maschinen.
Kapital: RM 17 500.—.
Anteilseigner: C. R. Nestmann, Meerane (100%).
Geschäftsführer: Oberg. A. Wackermann.
Prokurist: J. Frenzel.
Bankverbindungen: Meeraner Bank A.-G., Reichsbank, Meerane.
Postcheck-Konto: 115 485 Leipzig.
Geschäftsjahr: 1./1.—31./12.

Grundbesitz: 3700 qm, davon 1403 qm bebaut.
Anlagen: Fabrikationsräume mit Montagehalle, elektr. Schweißanlage; kaufm. u. techn. Büro.
Eigene Vertretung in Berlin: Willy Böckel, W 30, Martin-Luther-Str. 12.
Besondere Angaben: Gegründet als Spezialfirma moderner Trocknungsanlagen. Hergestellt wurden Schlangensammel-Trockenapparate System Otto und Getreidetrockner nach dem Zellsystem neben anderen allgemeinen Trocknungsanlagen. 1921 erfolgte neben dem Trocknerbau die Aufnahme der Fabrikation von Holzbearbeitungsmaschinen. 1932 übernahm das Werk zusätzlich die Herstellung von Bewoilt-Apparaturen für die Papierindustrie nach den Patenten Dr. Bruno Wieger, Berlin.

Piccolo-Automaten G. m. b. H., Berlin W 35, Kurfürstenstraße 146.

Fernruf: 21 20 95. **Drahtanschrift:** piccoloautomat
Gründung: 1932.
Fabrikationsprogramm: Schokoladen-Verkaufsapparate (Tischautomaten, Kugelschapparate).
Kapital: RM 20 000.—.
Anteilseigner: Dr. Richard Schönthal, Wien.
Geschäftsführer: Erwin Hantsch, Bln.-Steglitz.
Bankverbindung: Deutsche Bank u. Disconto-Ges., Berlin.
Postcheck-Konto: 8012 Berlin.
Geschäftsjahr: 1./1.—31./12.
Anlagen: Mechanische Werkstätten.
Gefolgschaft: 13 Arbeiter u. 6 Angestellte.

F. Piechatzek, Kran- u. Aufzug-Werke, Berlin N 65, Seestraße 51-56.

Fernruf: 46 43 11. **Drahtanschrift:** lüderszug
Gründung: 1885.
Fabrikationsprogramm: Krane u. Aufzüge, Hebezeuge u. Hebeemaschinen (Flaschenzüge, Laufkatzen, Winden, Elektro-Flaschenzüge).
Geschäftsleiter: Richard, Martin u. Paul Piechatzek.
Prokuristen: Paul Gräning, Alfred Knop, Otto Kuhwald.
Bankverbindung: Reichskredit-Gesellschaft A.-G., Berlin.
Postcheck-Konto: 4847 Berlin.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 9500 qm, davon 4500 qm bebaut.
Anlagen: Maschinenbau-Werkstätten (Dreherei, Schleiferei, Eiserei, Presserei u. Schlosserei).
Eigene Vertretungen: Im Ausland.
Besondere Angaben: Der Export erstreckt sich auf alle Weltteile.
Gefolgschaft: 350 Mitglieder.

Otto Pieron
siehe Maschinenfabrik

Paul Pietzschmann Wasserwerksbau, Berlin-Spandau, Schönwalder Str. 34.

Fernruf: 37 68 71.
Gründung: 1903.
Fabrikationsprogramm: Entwurf und Bau von Trink- und Gebrauchswasserwerken jeder Größe, Enteisungsanlagen einschl. der erforderlichen Antriebsanlagen (Dampf — Diesel — Wasserkraft, Hoch- und Niederspannung).
Inhaber: Ing. Paul Pietzschmann.
Bankverbindungen: Dresdner Bank, Spandauer Bank, Spandau.

Besondere Angaben: Der Firmeninhaber beschäftigt sich hauptsächlich mit der Erstellung halb- u. vollautomatisch arbeitender Wasserwerke u. besitzt hierüber große u. langjährige Erfahrungen.

Anton Piller Maschinenfabrik, Osterode (Harz), Abgunst 24.

Fernruf: 211. **Drahtanschrift:** apo
Gründung: 1909.
Fabrikationsprogramm: Ventilatoren für Heizungs-, Lüftungs-, Absauge u. sonstige Zwecke.
Bankverbindung: Reichsbank, Städt. Sparkasse, Osterode a. H.
Postcheck-Konto: 40 278 Hannover. ×

Pilot, G. m. b. H.
siehe Maschinenfabrik

Friedrich Piltz & Sohn K.-G., Heidenheim a. d. Brenz, Friedrichstr. 9.

Fernruf: 637. **Drahtanschrift:** piltz
Gründung: 1863.
Fabrikationsprogramm: Genauigkeitswerkzeuge für Herstellung u. Kontrolle von Gewinden; Gewindeschleifeinrichtungen; Drehbank-Schleifapparate.
Gesellschafter: Otto u. Walther Piltz.
Geschäftsführer: Die Gesellschafter.
Bankverbindungen: Reichsbank, Deutsche Bank u. Disconto-Ges., Heidenheim.
Postcheck-Konto: 2178 Stuttgart.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 2600 qm, davon 2000 qm bebaut; gepachtet sind 1200 qm mit 700 qm bebauter Fläche.
Anlagen: Verwaltungs- u. Fabrikationsräume in Heidenheim u. München.
Besondere Angaben: In München befindet sich ein Zweigwerk der Gesellschaft.

Eduard Pincuss Armaturenfabrik, Sanitäre Einrichtungen, Berlin O 17, Gr. Frankfurter Str. 13.

Fernruf: 59 13 18. **Drahtanschrift:** epal
Gründung: 1859.
Fabrikationsprogramm: Wasserleitungs-Armaturen.
Inhaber: Arthur Landsberger, Bln.-Charlottenburg; Ernst Reichenbach, Bln.-Grünwald.
Prokuristen: Paul Derpsch, Frieda Thierschmann.



Herausforderungen (Datenstrukturierung):

- Keine "all-in-one"-Lösung vorhanden: Software muss neu geschrieben werden
- Hohe Anforderungen an Qualität von Forschungsdaten
- Profiltrennung
- Inkonsistente Attributbezeichnungen:
 - *Postscheck-Konto, Postschekkonto, Postcheck-Konto*
 - *Geschäftsführer, Geschäftsleiter, Direktor, Betriebsführer*
 - ...

Phönix — Pincuss

Phönix-Werk G. m. b. H., Spezialfabrik moderner Trocken- Apparate, Meerane (Sa.).

Fernruf: 2424. **Drahtanschrift:** phönixwerk
Gründung: 1907.
Fabrikationsprogramm: Trockenapparate; Holzbearbeitungs-Maschinen.
Kapital: RM 17 500.—.
Anteilhaber: G. E. Nestmann, Meerane (100%).
Geschäftsführer: Ohering, A. Wackermann.
Prokurist: J. Frenzel.
Bankverbindungen: Meeraner Bank A.-G., Reichsbank, Meerane.
Postscheck-Konto: 115 485 Leipzig.
Geschäftsjahr: 1./1.—31./12.

Grundbesitz: 3700 qm, davon 1403 qm bebaut.
Anlagen: Fabrikationsräume mit Montagehalle, elektr. Schweißanlage; kaufm. u. techn. Büro.
Eigene Vertretung in Berlin: Willy Böckel, W 30, Martin-Luther-Str. 12.
Besondere Angaben: Gegründet als Spezialfirma moderner Trocknungsanlagen. Hergestellt wurden Schlangenbündel-Trockenapparate System Otto und Getreidetrockner nach dem Zellensystem neben anderen allgemeinen Trocknungsanlagen. 1921 erfolgte neben dem Trocknerbau die Aufnahme der Fabrikation von Holzbearbeitungsmaschinen. 1932 übernahm das Werk zusätzlich die Herstellung von Bewoilt-Apparaturen für die Papierindustrie nach den Patenten Dr. Bruno Wiegler, Berlin.

Piccolo-Automaten G. m. b. H., Berlin W 35, Kurfürstenstraße 146.

Fernruf: 21 20 95. **Drahtanschrift:** piccoloautomat
Gründung: 1932.
Fabrikationsprogramm: Schokoladen-Verkaufsapparate (Tischautomaten, Kugelschapparate).
Kapital: RM 20 000.—.
Anteilhaber: Dr. Richard Schönthal, Wien.
Geschäftsführer: Erwin Hantsch, Bln.-Steglitz.
Bankverbindung: Deutsche Bank u. Disconto-Ges., Berlin.
Postscheck-Konto: 8012 Berlin.
Geschäftsjahr: 1./1.—31./12.
Anlagen: Mechanische Werkstätten.
Gefolgschaft: 13 Arbeiter u. 6 Angestellte.

F. Piechatzek, Kran- u. Aufzug-Werke, Berlin N 65, Seestraße 51-56.

Fernruf: 46 43 11. **Drahtanschrift:** lüderszug
Gründung: 1885.
Fabrikationsprogramm: Krane u. Aufzüge, Hebezeuge u. Hebeemaschinen (Flaschenzüge, Laufkatzen, Winden, Elektro-Flaschenzüge).
Geschäftsleiter: Richard, Martin u. Paul Piechatzek.
Prokuristen: Paul Gräning, Alfred Knop, Otto Kuhwald.
Bankverbindung: Reichskredit-Gesellschaft A.-G., Berlin.
Postscheck-Konto: 4847 Berlin.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 9500 qm, davon 4500 qm bebaut.
Anlagen: Maschinenbau-Werkstätten (Dreherei, Schleiferei, Eiserei, Presserei u. Schlosserei).
Eigene Vertretungen: Im Ausland.
Besondere Angaben: Der Export erstreckt sich auf alle Weltteile.
Gefolgschaft: 350 Mitglieder.

Otto Pieron
siehe Maschinenfabrik

Paul Pietzschmann Wasserwerksbau, Berlin-Spandau, Schönwalder Str. 34.

Fernruf: 37 68 71.
Gründung: 1903.
Fabrikationsprogramm: Entwurf und Bau von Trink- und Gebrauchswasserwerken jeder Größe, Enteisungsanlagen einschl. der erforderlichen Antriebsanlagen (Dampf — Diesel — Wasserkraft, Hoch- und Niederspannung).
Inhaber: Ing. Paul Pietzschmann.
Bankverbindungen: Dresdner Bank, Spandauer Bank, Spandau.

Besondere Angaben: Der Firmeninhaber beschäftigt sich hauptsächlich mit der Erstellung halb- u. vollautomatisch arbeitender Wasserwerke u. besitzt hierüber große u. langjährige Erfahrungen.

Anton Piller Maschinenfabrik, Osterode (Harz), Abgunst 24.

Fernruf: 211. **Drahtanschrift:** apo
Gründung: 1909.
Fabrikationsprogramm: Ventilatoren für Heizungs-, Lüftungs-, Absaug- u. sonstige Zwecke.
Bankverbindung: Reichsbank, Städt. Sparkasse, Osterode a. H.
Postscheck-Konto: 40 278 Hannover. ×

Pilot, G. m. b. H.
siehe Maschinenfabrik

Friedrich Piltz & Sohn K.-G., Heidenheim a. d. Brenz, Friedrichstr. 9.

Fernruf: 637. **Drahtanschrift:** piltz
Gründung: 1863.
Fabrikationsprogramm: Genauigkeitswerkzeuge für Herstellung u. Kontrolle von Gewinden; Gewindeschleifeinrichtungen; Drehbank-Schleifapparate.
Gesellschafter: Otto u. Walther Piltz.
Geschäftsführer: Die Gesellschafter.
Bankverbindungen: Reichsbank, Deutsche Bank u. Disconto-Ges., Heidenheim.
Postscheck-Konto: 2178 Stuttgart.
Geschäftsjahr: Kalenderjahr.

Grundbesitz: 2600 qm, davon 2000 qm bebaut; gepachtet sind 1200 qm mit 700 qm bebauter Fläche.
Anlagen: Verwaltungs- u. Fabrikationsräume in Heidenheim u. München.
Besondere Angaben: In München befindet sich ein Zweigwerk der Gesellschaft.

Eduard Pincuss Armaturenfabrik, Sanitäre Einrichtungen, Berlin O 17, Gr. Frankfurter Str. 13.

Fernruf: 59 13 18. **Drahtanschrift:** epal
Gründung: 1859.
Fabrikationsprogramm: Wasserleitungs-Armaturen.
Inhaber: Arthur Landsberger, Bln.-Charlottenburg; Ernst Reichenbach, Bln.-Grünwald.
Prokuristen: Paul Derpsch, Frieda Thierschmann.

462

1. Einführung & Digitalisierung



Herausforderungen (Infrastruktur, Personal):

- Mehrteiliger Workflow mit spezifischen infrastrukturellen und personellen Anforderungen
- Digitalisierung (1 Projektkoordinator, 1 Hiwi)
- OCR (1 Projektkoordinator, 1 Entwickler)
- Datenstrukturierung (1 Entwickler)
- OCR-Server (eScriptorium)

Phönix — Pincuss

Phönix-Werk G. m. b. H., Spezialfabrik moderner Trocken- Apparate, Meerane (Sa.).

Fernruf: 2424. **Drahtanschrift:** phönixwerk
Gründung: 1907.
Fabrikationsprogramm: Trockenapparate; Holzbearbeitungs-Maschinen.
Kapital: RM 17 500.—.
Anteilhaber: G. R. Nestmann, Meerane (100%).
Geschäftsführer: Obering. A. Wackermann.
Prokurist: J. Frenzel.
Bankverbindungen: Meeraner Bank A.-G., Reichsbank, Meerane.
Postscheck-Konto: 115 485 Leipzig.
Geschäftsjahr: 1./1.—31./12.

Grundbesitz: 3700 qm, davon 1403 qm bebaut.
Anlagen: Fabrikationsräume mit Montagehalle, elektr. Schweißanlage; kaufm. u. techn. Büro.
Eigene Vertretung in Berlin: Willy Böckel, W 30, Martin-Luther-Str. 12.
Besondere Angaben: Gegründet als Spezialfirma moderner Trocknungsanlagen. Hergestellt wurden Schlangensammel-Trockenapparate System Otto und Getreidetrockner nach dem Zellsystem neben anderen allgemeinen Trocknungsanlagen. 1921 erfolgte neben dem Trocknerbau die Aufnahme der Fabrikation von Holzbearbeitungsmaschinen. 1932 übernahm das Werk zusätzlich die Herstellung von Bewöld-Apparaturen für die Papierindustrie nach den Patenten Dr. Bruno Wiegler, Berlin.

Piccolo-Automaten G. m. b. H., Berlin W 35, Kurfürstenstraße 146.

Fernruf: 21 20 95. **Drahtanschrift:** piccoloautomat
Gründung: 1932.
Fabrikationsprogramm: Schokoladen-Verkaufsapparate (Tischautomaten, Kugelschapparate).
Kapital: RM 20 000.—.
Anteilhaber: Dr. Richard Schönthal, Wien.
Geschäftsführer: Erwin Hantsch, Bln.-Steglitz.
Bankverbindung: Deutsche Bank u. Disconto-Ges., Berlin.
Postscheck-Konto: 8012 Berlin.
Geschäftsjahr: 1./1.—31./12.
Anlagen: Mechanische Werkstätten.
Gefolgschaft: 13 Arbeiter u. 6 Angestellte.

F. Piechatzek, Kran- u. Aufzug-Werke, Berlin N 65, Seestraße 51-56.

Fernruf: 46 43 11. **Drahtanschrift:** lüderszug
Gründung: 1885.
Fabrikationsprogramm: Krane u. Aufzüge, Hebezeuge u. Hebemaschinen (Flaschenzüge, Laufkatzen, Winden, Elektro-Flaschenzüge).
Geschäftsleiter: Richard, Martin u. Paul Piechatzek.
Prokuristen: Paul Gräning, Alfred Knop, Otto Kuhwald.
Bankverbindung: Reichskredit-Gesellschaft A.-G., Berlin.
Postscheck-Konto: 4847 Berlin.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 9500 qm, davon 4500 qm bebaut.
Anlagen: Maschinenbau-Werkstätten (Dreherei, Schleiferei, Eiserei, Presserei u. Schlosserei).
Eigene Vertretungen: Im Ausland.
Besondere Angaben: Der Export erstreckt sich auf alle Weltteile.
Gefolgschaft: 350 Mitglieder.

Otto Pieron
siehe Maschinenfabrik

Paul Pietzschmann Wasserwerksbau, Berlin-Spandau, Schönwalder Str. 34.

Fernruf: 37 68 71.
Gründung: 1903.
Fabrikationsprogramm: Entwurf und Bau von Trink- und Gebrauchswasserwerken jeder Größe, Enteisungsanlagen einschl. der erforderlichen Antriebsanlagen (Dampf — Diesel — Wasserkraft, Hoch- und Niederspannung).
Inhaber: Ing. Paul Pietzschmann.
Bankverbindungen: Dresdner Bank, Spandauer Bank, Spandau.

Besondere Angaben: Der Firmeninhaber beschäftigt sich hauptsächlich mit der Erstellung halb- u. vollautomatisch arbeitender Wasserwerke u. besitzt hierüber große u. langjährige Erfahrungen.

Anton Piller Maschinenfabrik, Osterode (Harz), Abgunst 24.

Fernruf: 211. **Drahtanschrift:** apo
Gründung: 1909.
Fabrikationsprogramm: Ventilatoren für Heizungs-, Lüftungs-, Absauge u. sonstige Zwecke.
Bankverbindung: Reichsbank, Städt. Sparkasse, Osterode a. H.
Postscheck-Konto: 40 278 Hannover. ×

Pilot, G. m. b. H.
siehe Maschinenfabrik

Friedrich Piltz & Sohn K.-G., Heidenheim a. d. Brenz, Friedrichstr. 9.

Fernruf: 637. **Drahtanschrift:** piltz
Gründung: 1863.
Fabrikationsprogramm: Genauigkeitswerkzeuge für Herstellung u. Kontrolle von Gewinden; Gewindeschleifeinrichtungen; Drehbank-Schleifapparate.
Gesellschafter: Otto u. Walther Piltz.
Geschäftsführer: Die Gesellschafter.
Bankverbindungen: Reichsbank, Deutsche Bank u. Disconto-Ges., Heidenheim.
Postscheck-Konto: 2178 Stuttgart.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 2600 qm, davon 2000 qm bebaut; gepachtet sind 1200 qm mit 700 qm bebauter Fläche.
Anlagen: Verwaltungs- u. Fabrikationsräume in Heidenheim u. München.
Besondere Angaben: In München befindet sich ein Zweigwerk der Gesellschaft.

Eduard Pincuss Armaturenfabrik, Sanitäre Einrichtungen, Berlin O 17, Gr. Frankfurter Str. 13.

Fernruf: 59 13 18. **Drahtanschrift:** epal
Gründung: 1859.
Fabrikationsprogramm: Wasserleitungs-Armaturen.
Inhaber: Arthur Landsberger, Bln.-Charlottenburg; Ernst Reichenbach, Bln.-Grünwald.
Prokuristen: Paul Derpsch, Frieda Thierschmann.

2. Layoutsegmentierung und OCR via eScriptorium

2. Layoutsegmentierung und OCR via eScriptorium

eScriptorium

Transkriptionssoftware

Open-Source-Plattform zur manuellen oder automatisierten Segmentierung und Texterkennung von historischen Handschriften und Drucken.



Entwickler:	PSL (Paris)
Erscheinungsjahr:	2018
Aktuelle Version:	0.10.5 (2022)
Betriebssystem:	plattformunabhängig
Programmiersprache:	Python, JavaScript, Hypertext Markup Language
OCR-Engine:	Kraken
Lizenz:	MIT license
Modelle:	https://ub-backup.bib.uni-mannheim.de/~stweil/eScriptorium/ (Fraktur)
Weitere Infos:	https://ocr-bw.bib.uni-mannheim.de/tag/escriptorium/

2. Layoutsegmentierung und OCR via eScriptorium

Kraken

OCR-Engine für Layoutsegmentierung und Texterkennung

Kraken ist eine auf historische Dokumente und nicht-lateinische Schriftsysteme optimierte "all-in-one"-OpenSource-OCR-Softwarelösung mit vollständig trainierbarer Layout-Analyse und Texterkennung.

Entwickler: Benjamin Kiessling

Erscheinungsjahr: 2015

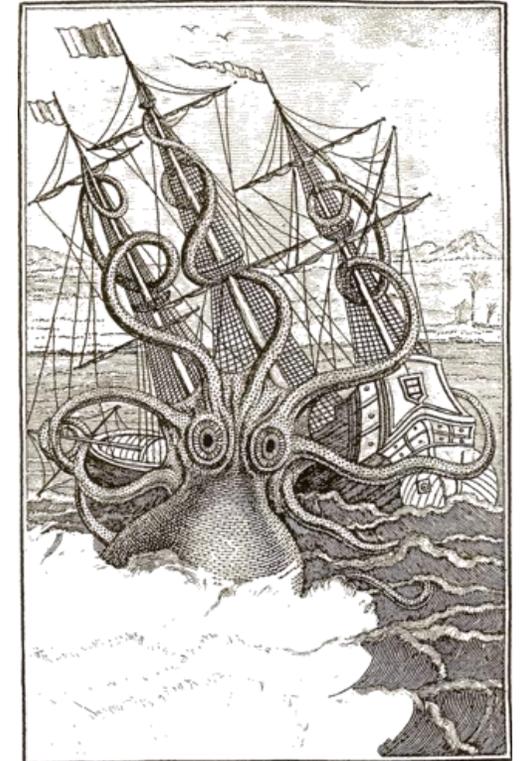
OCR-Technik: bounding box und baseline-based (eScriptorium)

Ausgabeformate: ALTO, PageXML, abbyXML and hOCR

Programmiersprache: Python

Modelle: <https://ub-backup.bib.uni-mannheim.de/~stweil/eScriptorium/> (Fraktur)

Lizenz: MIT License



eScriptorium-Workflow

Erkennung des Layouts und von Textinhalten

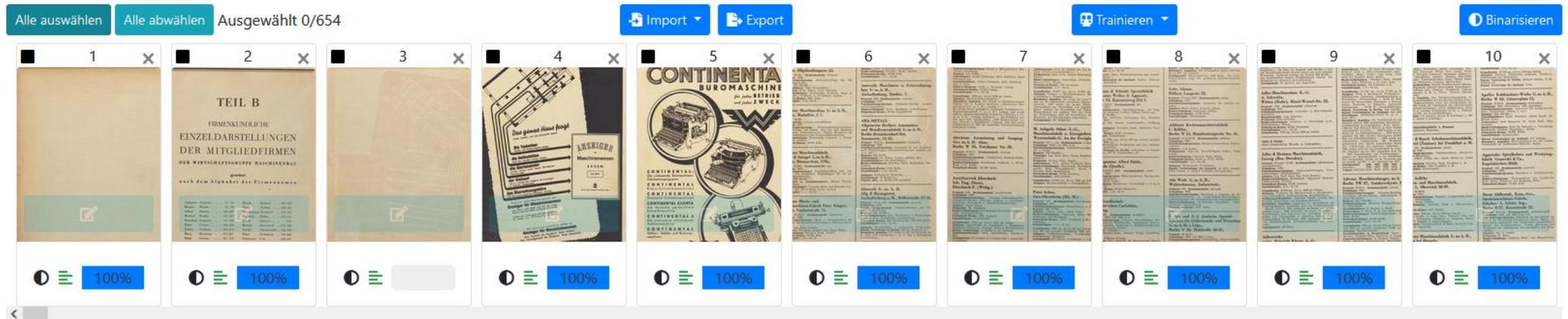
1. Upload der Digitalisate
2. Auswahl eines geeigneten Basismodells für Layoutsegmentierung und Texterkennung
3. Ermittlung von Schwachstellen der Modelle
4. Erstellung von Ground-Truth-Daten
5. Werkspezifisches Nachtraining der Modelle
6. Erkennung von Layout und Textinhalten

Wenn kein Basismodell gefunden werden kann, sollte auf ein anderes System ausgewichen werden oder es muss ein „from-scratch“ Modell trainiert werden (wesentlich höherer Trainingsaufwand)

2. Layoutsegmentierung und OCR via eScriptorium

Upload der Seiten

Digitalisate können über die GUI oder die API einzeln oder in einem Batch hochgeladen werden



Übersicht über die 654 Digitalisate in dem neu angelegten Projekt

Layoutsegmentierung

Auswahl eines geeigneten Basismodells

Die Segmentierung besteht aus 2 Schritten:

1. Segmentierung der Textregionen
2. Segmentierung der Textzeilen (Baseline und Polygonzug)



Geeignetes Basismodell:
cbad_1800_compensated_50

2. Layoutsegmentierung und OCR via eScriptorium



Schwachstelle: Segmentierung der Profile

Abwärme — Ackermann

Fabrikationsprogramm: Beton- u. Mörtelmischer; Bauwinden; Basenofen.
Kapital: RM 20.000.—
Anteilhaber: Arthur Schumann, Rich. Findeisen, Ernst Schumann sen.
Geschäftsführer: Arthur Schumann, Rich. Findeisen, Ing. O. Selz.
Bankverbindungen: Stadt- u. Girobank, Leipzig.
Postcheck-Konto: 70.550 Leipzig.
Geschäftsjahr: 1./1.—31./12.

Grundbesitz: 6000 qm, davon 3800 qm bebaut.
Anlagen: Maschinenwerkstätten.
Besondere Angaben: Die Fabrik ist Fabrikationsnachfolgerin der in Konkurs gerathenen „Allgemeinen Baumaschinen-Ges. m. b. H.“, Leipzig G. A.
Die Firma hat sich anfangl. Schumann, Findeisen & Co., Baumaschinenfabrik G. m. b. H., Leipzig, und wurde 1904 in „Baumaschinenfabrik Schumann, Findeisen & Co. G. m. b. H.“, Leipzig, umfirmirt. — Die Fabrikate werden unter dem Schutzzeichen „ABO-Baumaschinen“ so wie „Neoroll- u. Rib-Mischer“ verkauft. Aus dem ersten Schutzzeichen ist 1907 die endgültige Firmenbezeichnung entwickelt worden.

Abwärme Ausnutzung und Saugzug
Ges. m. b. H. Abas,
Berlin W 35, Potsdamer Str. 28.
Fernruf: 22 63 17. **Drachenschrift:** abwasch.
Gründung: 1921.
Fabrikationsprogramm: Ventilatoren, Staubabscheider, Gasanhalterbühnen.
Bankverbindungen: Reichsbank, Commerz- u. Privatbank, A.-G., Berlin.
Postcheck-Konto: 118.120 Berlin.

Acetylenwerk Ebersbach
Inh. Eug. Zinser,
Ebersbach-F. (Witthg.)
Fernruf: 216. **Drachenschrift:** acetylenwerk.
Gründung: 1898.
Fabrikationsprogramm: Antogen-Schweißapparate; Antogen-Werkzeuge (Schweiß-, Schneid- u. Lötlöffel); Sauerstoff-, Wasserstoff- u. Dünnsäureventile).
Inhaber: Eugen Zinser.
Präkurist: Th. Friedrich.
Bankverbindungen: Gewerbank, Ebersbach-F.; Reichsbank, Göttingen-Witthg.
Postcheck-Konto: 3228 Stuttgart.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 1500 qm, davon 1500 qm bebaut.
Anlagen: Apparatebau u. Schweißerei u. Klempnerei, Schneid-, Schlosserei, Montage, Eisen- u. Metalldreherei.
Eigene Vertretung: in Berlin u. im Ausland.
Tochtergesellschaft: 8844 „Spindelwerke, Ebersbach-F.“
Gefolgschaft: 30 Arbeiter (4 Lehrlinge) und 7 Angestellte (1 Lehrling).

Achenbach Söhne G. m. b. H.,
Buschhütten (Kr. Siegen).
Fernruf: Siegen 5011. **Drachenschrift:** achenbachsöhne.
Gründung: 1846.
Fabrikationsprogramm: Walzwerkzeug- und -walzenfabrik.
Kapital: RM 1.564.000.—

Anteilhaber: Frau Dr. Barten, Dr. Ing. Ernst Barten, Ernst Gieseler (Geschäftsführer).
Präkuristen: Karl Roth, Eduard Reinschmidt, Heinrich Bester.
Bankverbindungen: Reichsbank, Deutsche Bank und Disconto-Ges., Siegen.
Postcheck-Konto: 873 Dortmund.
Geschäftsjahr: 1./7.—30./6.

Grundbesitz: 16.600 qm, davon 16.000 qm bebaut.
Anlagen: Gießerei u. Modellschreineri; Maschinenbauwerkstätten u. Walzendreherei.
Besondere Angaben: Die Firma hat ihren Ursprung in dem im 15. Jahrhundert errichteten Hütten- und Hammerwerk. — 1846 wurde die jetzige Firma gegründet (Klempnerei zur Herstellung von gelohlenen Ofen).
Später nahm die Herstellung von Walzen und der Walzwerkzeug den größten Teil der Produktion ein.
Die Werkanlagen sind mit neuesten Bearbeitungs-Maschinen und Einrichtungen ausgestattet.
Gefolgschaft: 503 Mitglieder.

M. Achgelis Söhne A.-G.,
Maschinenfabrik u. Eisengießerei,
Wesermünde-G, An der Zweighahn 1.
Fernruf: 101 u. 146. **Drachenschrift:** achgeliswerke.
Gründung: 1883; seit 1918 A.-G.
Fabrikationsprogramm: Schälhülsmaschinen in jeder Art u. Größe.
Kapital: RM 225.000.—
Vorstand: Ing. Karl Boos, Georg Brinkmann, Werner Sander.
Präkurist: Abt. Ing. Wilh. Barth.
Aufsichtsrat: Vorn: Arthur Friedrichs, Bremerhaven, Bankverbindungen: Reichsbank, Wesermünde-G.; Gesellschafter-Bank, Wesermünde-G.; Nordl.-Kreditbank, Bremer Bank, Bremerhaven.
Postcheck-Konto: 19.028 Hamburg.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 18.500 qm, davon 7200 qm bebaut.
Anlagen: Maschinenfabrik u. Eisengießerei.

Peter Acker,
Gau-Odernheim (Rh. H.)
Fernruf: 226. **Drachenschrift:** maschinenacker.
Gründung: 1878.
Fabrikationsprogramm: Höhenförderer u. Pöbge.
Kapital: RM 70.000.—
Inhaber: Johann Acker, Jakob Acker.
Bankverbindungen: Volksbank, Alzey; Spars- u. Darlehnskass. Gau-Odernheim.
Postcheck-Konto: 23.577 Frankfurt a/M.; 8617 Ludwigshafen.
Geschäftsjahr: 1./1.—31./12.
Grundbesitz: 1200 qm, davon 800 qm bebaut; gepachtet sind 390 qm; gesamte Nutzfläche 1500 qm.
Anlagen: Fabrikationsanlage, Verfahrungs- u. Anstellstraßen, Treiberei.
Gefolgschaft: 29 Arbeiter, 5 Lehrlinge u. 1 Angestellter.

Friedrich Ackermann,
Werkzeug- und Maschinenfabrik,
Wuppertal-Barmen, Oberdenkmalstr. 89.
Fernruf: 54 282.
Gründung: 1912.
Fabrikationsprogramm: Maschinenschraubstöcke, Zahnrad-, Gewindestift-, Fräskutter-, Vorrichtungen- und Drehstöße.

Segmentierung mit cbad_1800

Optimale Segmentierung

Abwärme — Ackermann

Fabrikationsprogramm: Beton- u. Mörtelmischer; Bauwinden; Basenofen.
Kapital: RM 20.000.—
Anteilhaber: Arthur Schumann, Rich. Findeisen, Ernst Schumann sen.
Geschäftsführer: Arthur Schumann, Rich. Findeisen, Ing. O. Selz.
Bankverbindungen: Stadt- u. Girobank, Leipzig.
Postcheck-Konto: 70.550 Leipzig.
Geschäftsjahr: 1./1.—31./12.

Grundbesitz: 6000 qm, davon 3800 qm bebaut.
Anlagen: Maschinenwerkstätten.
Besondere Angaben: Die Fabrik ist Fabrikationsnachfolgerin der in Konkurs gerathenen „Allgemeinen Baumaschinen-Ges. m. b. H.“, Leipzig G. A.
Die Firma hat sich anfangl. Schumann, Findeisen & Co., Baumaschinenfabrik G. m. b. H., Leipzig, und wurde 1904 in „Baumaschinenfabrik Schumann, Findeisen & Co. G. m. b. H.“, Leipzig, umfirmirt. — Die Fabrikate werden unter dem Schutzzeichen „ABO-Baumaschinen“ so wie „Neoroll- u. Rib-Mischer“ verkauft. Aus dem ersten Schutzzeichen ist 1907 die endgültige Firmenbezeichnung entwickelt worden.

Abwärme Ausnutzung und Saugzug
Ges. m. b. H. Abas,
Berlin W 35, Potsdamer Str. 28.
Fernruf: 22 63 17. **Drachenschrift:** abwasch.
Gründung: 1921.
Fabrikationsprogramm: Ventilatoren, Staubabscheider, Gasanhalterbühnen.
Bankverbindungen: Reichsbank, Commerz- u. Privatbank, A.-G., Berlin.
Postcheck-Konto: 118.120 Berlin.

Acetylenwerk Ebersbach
Inh. Eug. Zinser,
Ebersbach-F. (Witthg.)
Fernruf: 216. **Drachenschrift:** acetylenwerk.
Gründung: 1898.
Fabrikationsprogramm: Antogen-Schweißapparate; Antogen-Werkzeuge (Schweiß-, Schneid- u. Lötlöffel); Sauerstoff-, Wasserstoff- u. Dünnsäureventile).
Inhaber: Eugen Zinser.
Präkurist: Th. Friedrich.
Bankverbindungen: Gewerbank, Ebersbach-F.; Reichsbank, Göttingen-Witthg.
Postcheck-Konto: 3228 Stuttgart.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 1500 qm, davon 1500 qm bebaut.
Anlagen: Apparatebau u. Schweißerei u. Klempnerei, Schneid-, Schlosserei, Montage, Eisen- u. Metalldreherei.
Eigene Vertretung: in Berlin u. im Ausland.
Tochtergesellschaft: 8844 „Spindelwerke, Ebersbach-F.“
Gefolgschaft: 30 Arbeiter (4 Lehrlinge) und 7 Angestellte (1 Lehrling).

Achenbach Söhne G. m. b. H.,
Buschhütten (Kr. Siegen).
Fernruf: Siegen 5011. **Drachenschrift:** achenbachsöhne.
Gründung: 1846.
Fabrikationsprogramm: Walzwerkzeug- und -walzenfabrik.
Kapital: RM 1.564.000.—

Anteilhaber: Frau Dr. Barten, Dr. Ing. Ernst Barten, Ernst Gieseler (Geschäftsführer).
Präkuristen: Karl Roth, Eduard Reinschmidt, Heinrich Bester.
Bankverbindungen: Reichsbank, Deutsche Bank und Disconto-Ges., Siegen.
Postcheck-Konto: 873 Dortmund.
Geschäftsjahr: 1./7.—30./6.

Grundbesitz: 16.600 qm, davon 16.000 qm bebaut.
Anlagen: Gießerei u. Modellschreineri; Maschinenbauwerkstätten u. Walzendreherei.
Besondere Angaben: Die Firma hat ihren Ursprung in dem im 15. Jahrhundert errichteten Hütten- und Hammerwerk. — 1846 wurde die jetzige Firma gegründet (Klempnerei zur Herstellung von gelohlenen Ofen).
Später nahm die Herstellung von Walzen und der Walzwerkzeug den größten Teil der Produktion ein.
Die Werkanlagen sind mit neuesten Bearbeitungs-Maschinen und Einrichtungen ausgestattet.
Gefolgschaft: 503 Mitglieder.

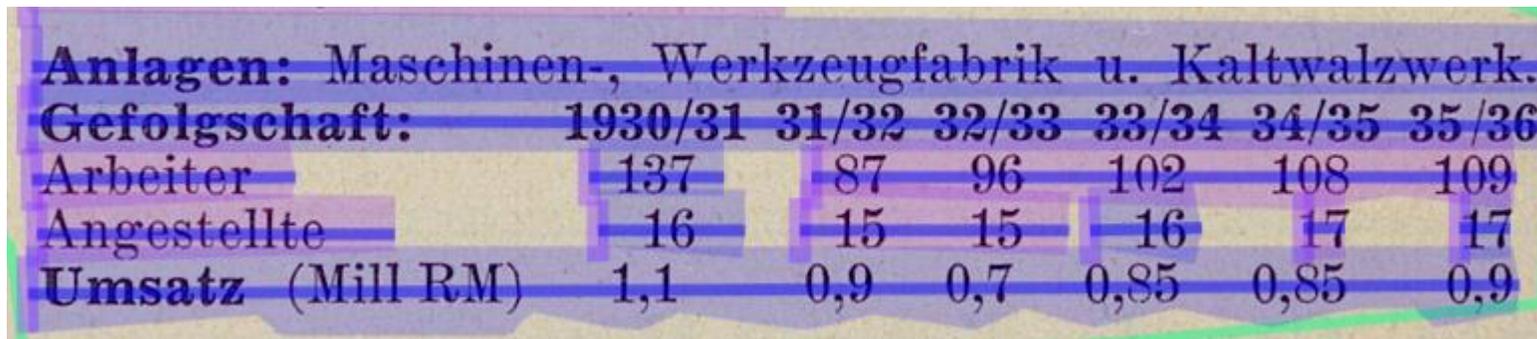
M. Achgelis Söhne A.-G.,
Maschinenfabrik u. Eisengießerei,
Wesermünde-G, An der Zweighahn 1.
Fernruf: 101 u. 146. **Drachenschrift:** achgeliswerke.
Gründung: 1883; seit 1918 A.-G.
Fabrikationsprogramm: Schälhülsmaschinen in jeder Art u. Größe.
Kapital: RM 225.000.—
Vorstand: Ing. Karl Boos, Georg Brinkmann, Werner Sander.
Präkurist: Abt. Ing. Wilh. Barth.
Aufsichtsrat: Vorn: Arthur Friedrichs, Bremerhaven, Bankverbindungen: Reichsbank, Wesermünde-G.; Gesellschafter-Bank, Wesermünde-G.; Nordl.-Kreditbank, Bremer Bank, Bremerhaven.
Postcheck-Konto: 19.028 Hamburg.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 18.500 qm, davon 7200 qm bebaut.
Anlagen: Maschinenfabrik u. Eisengießerei.

Peter Acker,
Gau-Odernheim (Rh. H.)
Fernruf: 226. **Drachenschrift:** maschinenacker.
Gründung: 1878.
Fabrikationsprogramm: Höhenförderer u. Pöbge.
Kapital: RM 70.000.—
Inhaber: Johann Acker, Jakob Acker.
Bankverbindungen: Volksbank, Alzey; Spars- u. Darlehnskass. Gau-Odernheim.
Postcheck-Konto: 23.577 Frankfurt a/M.; 8617 Ludwigshafen.
Geschäftsjahr: 1./1.—31./12.
Grundbesitz: 1200 qm, davon 800 qm bebaut; gepachtet sind 390 qm; gesamte Nutzfläche 1500 qm.
Anlagen: Fabrikationsanlage, Verfahrungs- u. Anstellstraßen, Treiberei.
Gefolgschaft: 29 Arbeiter, 5 Lehrlinge u. 1 Angestellter.

Friedrich Ackermann,
Werkzeug- und Maschinenfabrik,
Wuppertal-Barmen, Oberdenkmalstr. 89.
Fernruf: 54 282.
Gründung: 1912.
Fabrikationsprogramm: Maschinenschraubstöcke, Zahnrad-, Gewindestift-, Fräskutter-, Vorrichtungen- und Drehstöße.

Schwachstelle: Textzeilen

Bei der Textzeilenerkennung ist die Baseline nicht immer korrekt und teilweise werden Zeilen in kleine Bereiche unterteilt



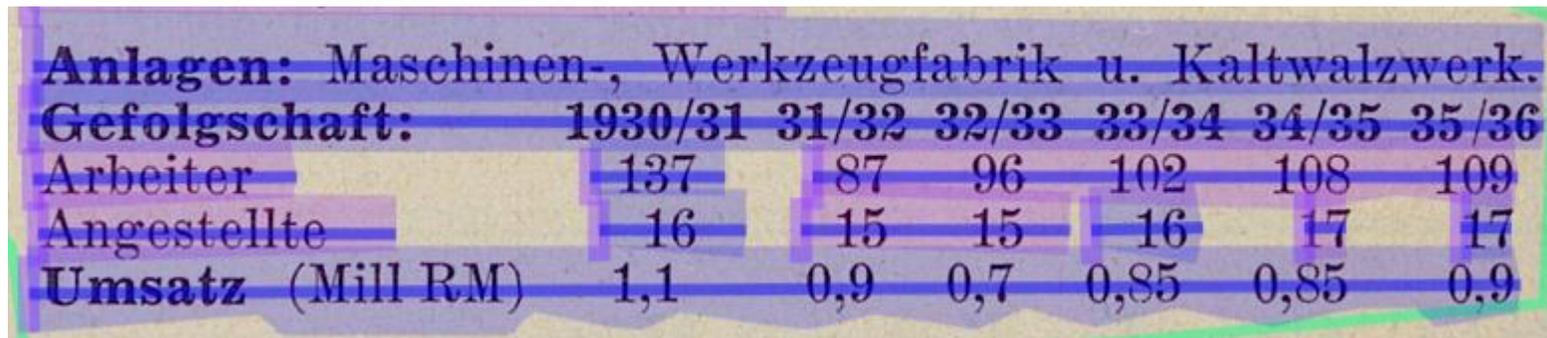
Anlagen:	Maschinen-, Werkzeugfabrik u. Kaltwalzwerk.					
Gefolgschaft:	1930/31	31/32	32/33	33/34	34/35	35/36
Arbeiter	137	87	96	102	108	109
Angestellte	16	15	15	16	17	17
Umsatz (Mill RM)	1,1	0,9	0,7	0,85	0,85	0,9

Beispiel: Textzeilensegmentierung einer Tabelle

Werkspezifisches Training

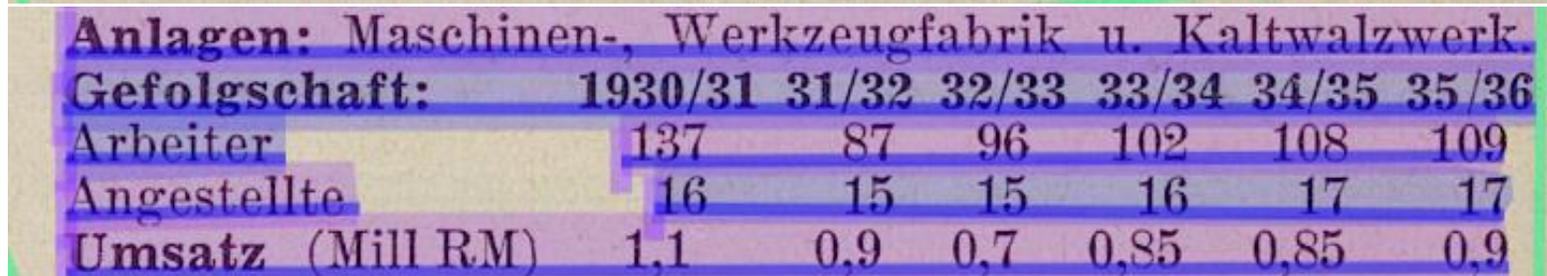
- Ground Truth:** 47 Seiten
Ergebnisse Textregionen: Unbefriedigend (Lösung: Verwendung einer anderen Vor- und Nachverarbeitung)
Ergebnisse Textzeilen: Deutliche Verbesserung

cbad_1800



Anlagen: Maschinen-, Werkzeugfabrik u. Kaltwalzwerk.						
Gefolgschaft:	1930/31	31/32	32/33	33/34	34/35	35/36
Arbeiter	137	87	96	102	108	109
Angestellte	16	15	15	16	17	17
Umsatz (Mill RM)	1,1	0,9	0,7	0,85	0,85	0,9

cbad_1800_trained



Anlagen: Maschinen-, Werkzeugfabrik u. Kaltwalzwerk.						
Gefolgschaft:	1930/31	31/32	32/33	33/34	34/35	35/36
Arbeiter	137	87	96	102	108	109
Angestellte	16	15	15	16	17	17
Umsatz (Mill RM)	1,1	0,9	0,7	0,85	0,85	0,9

OCR

Auswahl eines geeigneten Basismodells:

Ein geeignetes Basismodell sollte

- 1) den Zeichenvorrat des Quellmaterials (möglichst) abbilden
- 2) bereits eine gute Erkennungsrate auf das Quellmaterial garantieren

Geeignetes OCR-Modell: digitue



Schwachstellen in der Texterkennung

Das Modell liefert bereits sehr gute Ergebnisse mit Texterkennungsraten > 98 %

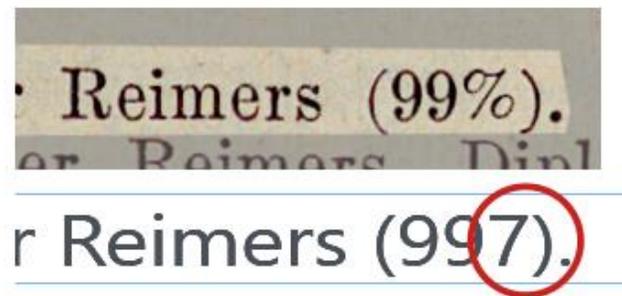
Einige spezifische Fehler lassen sich identifizieren:

Prozentzeichen



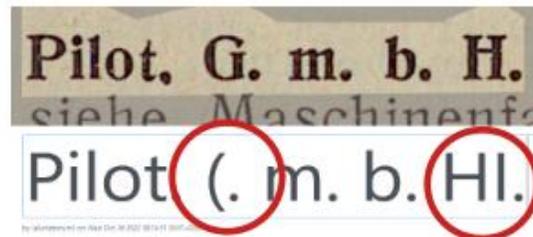
RM 24 000.—.
igner: Oscar Klätte (80%), Ilse
innemann (10%).
igner: Oscar Klätte (8025), Ilse

Oct 26 2022 09:14:22 GMT+0200

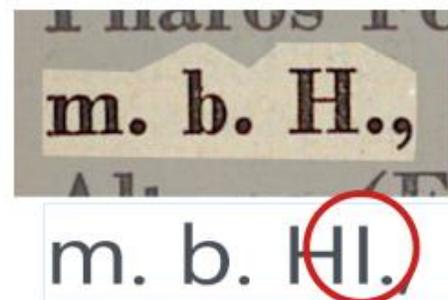


Reimers (99%).
er Reimers Dipl
r Reimers (997).

Zeichenkette „G.m.b.H.“

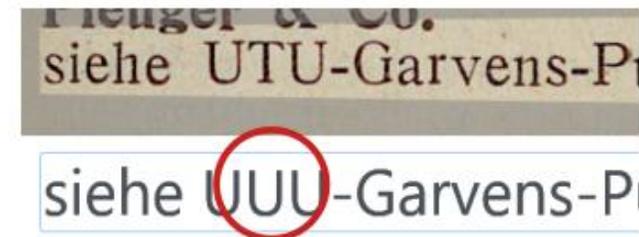


Pilot, G. m. b. H.
siehe Maschinenfa
Pilot (. m. b. Hl.)

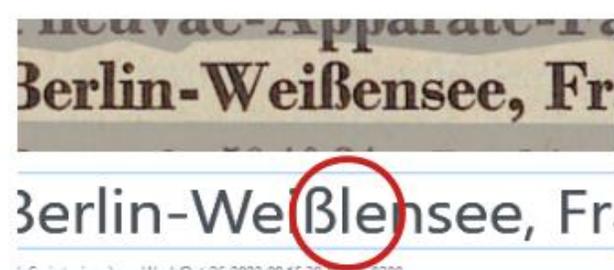


m. b. H.,
Al.
m. b. Hl.)

Seltene Zeichen / Zeichenkombinationen



Heuger & Co.
siehe UTU-Garvens-P
siehe UUU-Garvens-P



neue Apparate-
Berlin-Weißensee, Fr
Berlin-Weißensee, Fr

2. Layoutsegmentierung und OCR via eScriptorium

Werkspezifisches Training

Ground Truth:

26 Seiten

Ergebnisse Texterkennung:

Die Erkennungsrate konnte auf > 99.85 % angehoben werden

digitue

RM 24 000.—.
igner: Oscar Klatte (80%), Ilse
innemann (10%).
igner: Oscar Klatte (8025), Ilse

Pilot, G. m. b. H.
siehe Maschinenfa
Pilot. (. m. b. Hl.

Pleuger & Co.
siehe UTU-Garvens-P
siehe UUU-Garvens-P

digitue_trained

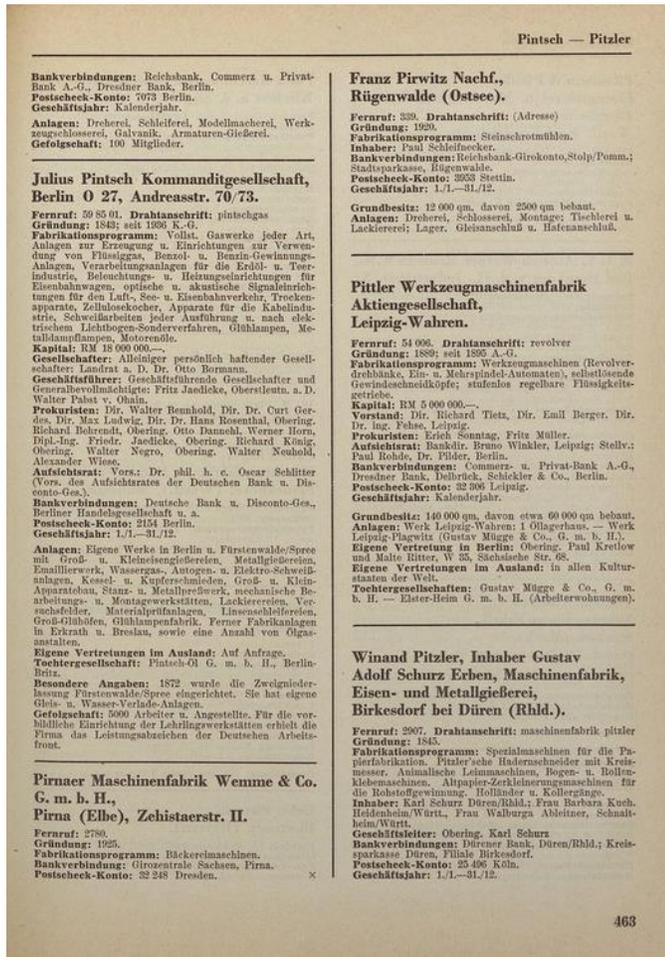
M 24 000.—.
er: Oscar Klatte (80%), Ilse
emann (10%).
er: Oscar Klatte (80%), Ilse

Pilot, G. m. b. H.
siehe Maschinenfa
Pilot. G. m. b. H.

Pleuger & Co.
siehe UTU-Garvens
siehe UTU-Garvens

3. Datenextraktion und -strukturierung via Python

3. Datenextraktion und -strukturierung via Python: Dateninput



```
1 <?xml version="1.0" encoding="UTF-8" standalone="yes"?>
2 <PcGts xmlns="http://schema.primaresearch.org/PAGE/gts/pagecontent/2019-07-15" xmlns
3 xsi:schemaLocation="http://schema.primaresearch.org/PAGE/gts/pagecontent/2019-07-15
4 http://schema.primaresearch.org/PAGE/gts/pagecontent/2019-07-15/pagecontent.xsd">
5 <Metadata>
6 <Creator>escriptorium</Creator>
7 <Created>2022-09-19T12:10:54.541269+00:00</Created>
8 <LastChange>2022-09-19T12:10:54.541290+00:00</LastChange>
9 </Metadata>
10 <Page imageFilename="maschinenindustrie_1937_0452.jpg" imageWidth="1000" imageHeight="1000">
11 <TextRegion id="r_2_1" >
12 <Coords points="127,354 1290,354 1290,685 127,685"/>
13 <TextLine id="eSc_line_835f5cfa" >
14 <Coords points="1269,350 1148,348 1132,348 1109,348 1086
15 880,348 857,348 834,348 825,350 818,350 815,350 811,348 811,348
16 673,348 650,348 632,348 627,348 605,348 582,348 559,348 536,348
17 307,348 284,348 261,348 238,348 215,348 192,348 169,348 146,348
18 <Baseline points="126,382 1292,387"/>
19 <TextEquiv>
20 <Unicode>Bankverbindungen: Reichsbank, Commerz u. Priv
21 </TextEquiv>
22 </TextLine>
23
24
25
26 <TextLine id="eSc_line_eff01337" >
27 <Coords points="128,426 126,449 877,449 880,430 877,396 128,389 128,426"/>
28 <Baseline points="128,426 880,430"/>
29 <TextEquiv>
30 <Unicode>Bank A.-G., Dresdner Bank, Berlin.</Unicode>
31 </TextEquiv>
32 </TextLine>
```

Reicht PAGE-XML aus, um Forschungsfrage zu beantworten?

JPG

PAGE-XML

Blatt Funktionalitäten:

1. liest PAGE-XML
2. entfernt Bindestriche
3. konvertiert PAGE-XML in TXT- und TSV-Formate
4. bietet ein Modul zur Datenstrukturierung

Blatt

pypi package 0.1.6

NLP-helper for OCR-ed pages in [PAGE XML](#) format.

Installation

```
pip install blatt
```

Open Source + CLI



<https://pypi.org/project/blatt>

<https://github.com/UB-Mannheim/blatt>

OCR-Segmentierung der Textregionen mit Problemen

blatt



1. sortiert Textzeilen
2. führt eine geometrische Segmentierung auf Basis der Koordinaten durch
3. Datenbereinigung
4. führt Segmente nachfolgender Spalten und Seiten zusammen

JPG



Transformation
via blatt

Ansatz: Doppelpunkt als
Separator nutzen

Problem: inkonsistente
Attributbezeichnung,
Druckfehler, OCR Fehler

e.g.: OCR problems:

*{'Postscheck-Konto', 'Postseheck-Konto', 'Postscheckkonto', 'Postscheck-Konnto', 'Fostscheck-Konto', 'Postcheck-Konto', 'Postscheek-Konto',
'Potscheck-Konto', 'Postsbeck-Konto', 'Postscheckkonto', 'Postschek-Konto', 'Postscheck-Konten', 'Ponstscheck-Konto', 'Postscheck-Konto'}*

Strukturierte Forschungsdaten

Firma	Piccolo-Automaten G.m.b.H., Berlin W 35, Kurfürstenstraße 146
Rechtsform	G.m.b.H.
Fernruf	212095
Drahtanschrift	piccoloautomat
Gründung	1932
Produkt	Schokoladen- Verkaufsapparate (Tischautomaten, Kugelstechapparate)

...

Sortierung und Gruppierung

Postscheck-Konto

```
{'Postscheck-Konto', 'Postseheck-Konto', 'PostscheckKonto', 'Postscheck-Konnto', 'Fostscheck-Konto', 'Postcheck-Konto',  
'Postscheek-Konto', 'Potscheck-Konto', 'Postscheck-Konto', 'Postscheckkonto', 'Postschek-Konto', 'Postscheck-Konten',  
'Ponstscheck-Konto', 'Postscheck-Konto'}
```

Geschäftsjahr

```
{'Geschäftjahr', 'Geschäftsjahr', 'Gescbäftsjahr', 'Geschäftsjahr', '.Geschäftsjahr'}
```

Fabrikationsprogramm

```
{'Fabfikationsprogramm', 'Fabrikationprogramm', 'Fabrikationsprorammm', 'Fabrikationsprogramm:',  
'Fabrikstionsprogramm', 'Fabrkkationsprogramm', 'Fabrikationsprogramm'}
```

Extraktion *Drahtanschrift* aus *Fernruf*:

Fernruf: 212095. Drahtanschrift: piccoloautomat



Fernruf	212095
Drahtanschrift	piccoloautomat

Extraktion der Rechtsformen:

{'m. b. H.', 'G. m. b. H.', 'Gesellschaft m. b. H.', 'G.m.b. H.', 'Ges. m.b. H.', 'GmbH', 'G. m. b.H.', 'mit beschr. Haftung', 'Ges.m.b.H.', 'G.m.b.H.', 'GmbH.', 'GmbH'}

{'A.-G.', 'Aktien-Gesellschaft', 'Actien-Gesellschaft', 'Aktiengesellschaft', 'Akt.-Ges.', 'Aktien-Gesellsch.', 'A. G.'}

{'K.-G.', 'Kom.-Ges.', 'Komm.-Ges.', 'KG.', 'Kommanditgesellschaft', 'Kom. Ges.', 'Komm.-Ges.'}

3. Datenextraktion und -strukturierung via Python: CSV und XLSX

Strukturierte Forschungsdaten

Firma	Piccolo-Automaten G.m.b.H., Berlin W 35, Kurfürstenstraße 146
Rechtsform	G.m.b.H.
Fernruf	212095
Drahtanschrift	piccoloautomat
Gründung	1932
Produkt	Schokoladen-Verkaufsapparate (Tischautomaten, Kugelstechapparate)

...

CSV und XLSX

	Company	RAW_TEXT	W_TEXT_1	LE_SEGME	FABRIKATIONSPROGRAMM	POSTSCHECK-KONTO	FERNRUF	DRAHTANSCHRIFT	BANKVERBINDUNGEN	ANLAGEN
0	Aachener Kratz	Aachener K Cassalette f Aachen, Oli Fernruf: 34 Gründung: Fabrikation Verwendun Kapital: RM Geschäftsfü Prokuristen Bankverbin Disconto-G Postscheck: Geschäftsj	Aachener K	/MI1937/n	Kratzenbeschläge für alle Ver	2952 Köln	34 041	kratzena	Reichsbank, Deutsche Bank u. Disconto-Ges	
1	Aachener Mas	Aachener M Aachen, Ru Fernruf: 25 Drahtansch Gründung: Fabrikation zur Herstell Drahtartike die Kratzen Bankverbin bank A.-G.,	Aachener N	/MI1937/n	Drahtbearbeitungsmaschinen	16 987 Köln	25 205	aachener maschinenbau	Reichsbank, Commerz- u. Privatbank A.-G.,	
2		Aachener N Rothe & Ste Aachen, Bis Fernruf: 25 Gründung: Fabrikation Lieferung v samte Kratz Kapital: RM Anteileigne Peter Rump								

4. Fazit

Alle Inhalte dieser Präsentation stehen unter der [Lizenz Creative Commons BY 4.0 International](https://creativecommons.org/licenses/by/4.0/), sofern nicht anders angegeben.





Ergebnisse:

- Sehr gute Texterkennungsgenauigkeit (> 99,85%) durch werkspezifisches Training
- eScriptorium eignet sich sehr gut als Plattform für Ground-Truth-Produktion und Training:
 - schnelle, benutzerfreundliche, plattformunabhängige Erstellung von Transkriptionen
 - benutzerfreundlicher Trainingsprozess (GUI)
- Datenstrukturierung:
 - keine All-in-One-Lösung → wir haben das Open-Source-Tool *blatt* entwickelt → es kann in ähnlichen Projekten wiederverwendet werden
 - wir haben alle verfügbaren Daten strukturiert (nicht nur Rechtsformen) → die Daten können für andere Forschungsfragen wiederverwendet werden
- Services der UB Mannheim können flexibel und kooperativ in Projekten mit anderen Stakeholdern der Universität genutzt werden

Phönix — Pincuss

Phönix-Werk G. m. b. H., Spezialfabrik moderner Trocken- Apparate, Meerane (Sa.).

Fernruf: 2424. Drahtanschrift: phönixwerk
Gründung: 1907.
Fabrikationsprogramm: Trockenapparate; Holzbearbeitungs-Maschinen.
Kapital: RM 17 500.—.
Anteilseigner: G. E. Nestmann, Meerane (100%).
Geschäftsführer: Obering. A. Wackernann.
Prokurist: J. Frenzel.
Bankverbindungen: Meeraner Bank A.-G., Reichsbank, Meerane.
Postscheck-Konto: 115 485 Leipzig.
Geschäftsjahr: 1./1.—31./12.

Grundbesitz: 3700 qm, davon 1403 qm bebaut.
Anlagen: Fabrikationsräume mit Montagehalle, elektr. Schweißanlage; kaufm. u. techn. Büro.
Eigene Vertretung in Berlin: Willy Böckel, W 30, Martin-Luther-Str. 12.
Besondere Angaben: Gegründet als Spezialfirma moderner Trocknungsanlagen. Hergestellt wurden Schlangensammel-Trockenapparate System Otto und Getreidetrockner nach dem Zellen-system neben anderen allgemeinen Trocknungsanlagen. 1921 erfolgte neben dem Trocknerbau die Aufnahme der Fabrikation von Holzbearbeitungsmaschinen. 1932 übernahm das Werk zusätzlich die Herstellung von Bewoilt-Apparaturen für die Papierindustrie nach den Patenten Dr. Bruno Wiegler, Berlin.

Piccolo-Automaten G. m. b. H., Berlin W 35, Kurfürstenstraße 146.

Fernruf: 21 20 95. Drahtanschrift: piccoloautomat
Gründung: 1932.
Fabrikationsprogramm: Schokoladen-Verkaufsapparate (Tischautomaten, Kugelschapparate).
Kapital: RM 20 000.—.
Anteilseigner: Dr. Richard Schöenthal, Wien.
Geschäftsführer: Erwin Hantsch, Bln.-Steglitz.
Bankverbindung: Deutsche Bank u. Disconto-Ges., Berlin.
Postscheck-Konto: 8012 Berlin.
Geschäftsjahr: 1./1.—31./12.
Anlagen: Mechanische Werkstätten.
Gefolgschaft: 13 Arbeiter u. 6 Angestellte.

F. Piechatzek, Kran- u. Aufzug-Werke, Berlin N 65, Seestraße 51-56.

Fernruf: 46 43 11. Drahtanschrift: lüderszug
Gründung: 1885.
Fabrikationsprogramm: Krane u. Aufzüge, Hebezeuge u. Hebe-maschinen (Flaschenzüge, Laufkatzen, Winden, Elektro-Flaschenzüge).
Geschäftsleiter: Richard, Martin u. Paul Piechatzek.
Prokuristen: Paul Gräning, Alfred Knop, Otto Kuhwald.
Bankverbindung: Reichskredit-Gesellschaft A.-G., Berlin.
Postscheck-Konto: 4847 Berlin.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 9500 qm, davon 4500 qm bebaut.
Anlagen: Maschinenbau-Werkstätten (Dreherei, Schleiferei, Eiserei, Presserei u. Schlosserei).
Eigene Vertretungen: Im Ausland.
Besondere Angaben: Der Export erstreckt sich auf alle Weltteile.
Gefolgschaft: 350 Mitglieder.

Otto Pieron
siehe Maschinenfabrik

Paul Pietzschmann Wasserwerksbau, Berlin-Spandau, Schönwalder Str. 34.

Fernruf: 37 68 71.
Gründung: 1903.
Fabrikationsprogramm: Entwurf und Bau von Trink- und Gebrauchswasserwerken jeder Größe, Enteisungsanlagen einschl. der erforderlichen Antriebsanlagen (Dampf — Diesel — Wasserkraft, Hoch- und Niederspannung).
Inhaber: Ing. Paul Pietzschmann.
Bankverbindungen: Dresdner Bank, Spandauer Bank, Spandau.

Besondere Angaben: Der Firmeninhaber beschäftigt sich hauptsächlich mit der Erstellung halb- u. vollautomatisch arbeitender Wasserwerke u. besitzt hierüber große u. langjährige Erfahrungen.

Anton Piller Maschinenfabrik, Osterode (Harz), Abgunst 24.

Fernruf: 211. Drahtanschrift: apo
Gründung: 1909.
Fabrikationsprogramm: Ventilatoren für Heizungs-, Lüftungs-, Absauge u. sonstige Zwecke.
Bankverbindung: Reichsbank, Städt. Sparkasse, Osterode a. H.
Postscheck-Konto: 40 278 Hannover. ×

Pilot, G. m. b. H.
siehe Maschinenfabrik

Friedrich Piltz & Sohn K.-G., Heidenheim a. d. Brenz, Friedrichstr. 9.

Fernruf: 637. Drahtanschrift: piltz
Gründung: 1863.
Fabrikationsprogramm: Genauigkeitswerkzeuge für Herstellung u. Kontrolle von Gewinden; Gewindeschleif-einrichtungen; Drehbank-Schleifapparate.
Gesellschafter: Otto u. Walther Piltz.
Geschäftsführer: Die Gesellschafter.
Bankverbindungen: Reichsbank, Deutsche Bank u. Disconto-Ges., Heidenheim.
Postscheck-Konto: 2178 Stuttgart.
Geschäftsjahr: Kalenderjahr.
Grundbesitz: 2600 qm, davon 2000 qm bebaut; gepachtet sind 1200 qm mit 700 qm bebauter Fläche.
Anlagen: Verwaltungs- u. Fabrikationsräume in Heidenheim u. München.
Besondere Angaben: In München befindet sich ein Zweigwerk der Gesellschaft.

Eduard Pincuss Armaturenfabrik, Sanitäre Einrichtungen, Berlin O 17, Gr. Frankfurter Str. 13.

Fernruf: 59 13 18. Drahtanschrift: epal
Gründung: 1859.
Fabrikationsprogramm: Wasserleitungs-Armaturen.
Inhaber: Arthur Landsberger, Bln.-Charlottenburg; Ernst Reichenbach, Bln.-Grünewald.
Prokuristen: Paul Derspech, Frieda Thierschmann.

Zeitaufwand:

- **Insgesamt:** ca. 2 Arbeitswochen (verteilt auf 2 Monate)
- **Digitalisierung:** 1 Arbeitstag
- **OCR:** 1 Arbeitswoche:
 - Auswertung vorhandener OCR-Modelle auf dem Material
 - Entwurf eines OCR-Workflows (Layoutsegmentierung + OCR)
 - Ground-Truth-Produktion und anschließendes Training (Layoutsegmentierung + OCR)
 - Endgültige OCR für alle Seiten
- **Strukturierung der Daten:** 1 Arbeitswoche:
 - Segmentierung
 - Nachbearbeitung
 - Erweiterung *blatt*

Team:

- **6 Projektbeteiligte**
 - Prof. Jochen Streb (**Professur für Wirtschaftsgeschichte @ Uni Mannheim**)
 - 1 Hilfskraft für Ground-Truth-Produktion und QA (**Professur für Wirtschaftsgeschichte**)
 - 1 Projektleitung Digitalisierung (**UB Mannheim**)
 - 1 Projektkoordinator OCR (**UB Mannheim**)
 - 1 Entwickler OCR (**UB Mannheim**)
 - 1 Entwickler Datenstrukturierung (**BERD@NFDI, UB Mannheim**)

Feedback, Fragen?

Jan Kamlah (Development): jan.kamlah@uni-mannheim.de
Renat Shigapov (Development): : renat.shigapov@uni-mannheim.de
Thomas Schmidt (Project coordination): thomas.schmidt@uni-mannheim.de

Danke!

<https://github.com/UB-Mannheim/Maschinen-Industrie>

Funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) – 460037581