

# Der Weg zum nutzbaren Volltext. Werkspezifisches Training als Baustein der OCR-Volltexterkennung für Alte Drucke

Personen mit 2 337 812 074 *M.* versichert blieben, davon 552 246  
Personen mit 1 847 622 742 *M.* bei den 35 Se im  
Deutschen Reich. Der gesammte Zuwachs im Jahre 1877 stellt sich  
auf 15 777 Personen mit 110 873 820 *M.* d. h. 2,14 % der  
Zahl der Versicherten und 4,98 % der Wachen,  
summe. Die Berücksichtigung der Unfälle, welche die  
Liquidation einiger Gesellschafenschaften im Jahre 1877 zur Folge gehabt  
hat vermindert sich der Zuwachs auf 8542 (1,16%) der Versicherten

Jan Kamlah und Thomas Schmidt

# Der Weg zum nutzbaren Volltext. Werkspezifisches Training als Baustein der OCR-Volltexterkennung für Alte Drucke



1. OCR und werkspezifisches Training an der UB Mannheim
2. Einführung
3. Training mit synthetischer Ground Truth (Ansatz A)
4. Training mit realer Ground Truth (Ansatz B)
5. Fazit

# Der Weg zum nutzbaren Volltext. Werkspezifisches Training als Baustein der OCR-Volltexterkennung für Alte Drucke



1. OCR und werkspezifisches Training an der UB Mannheim
2. Einführung
3. Training mit synthetischer Ground Truth (Ansatz A)
4. Training mit realer Ground Truth (Ansatz B)
5. Fazit

# 1. OCR und werkspezifisches Training an der UB Mannheim

- 2022: werkspezifisches Training (eScriptorium + Kraken) mit theologischen Werken (DigiTheo Ground Truth: <https://github.com/UB-Mannheim/DTGT>)
- Seit 2021: Modulprojekt „Werkspezifisches Training“ (für Tesseract + Calamari) im Rahmen der 3. DFG-Förderphase von OCR-D
- 2013–2019: Digitalisierung und OCR „Aktienführer-Datenarchiv“
- 2019: Werkspezifisches Training für „Weisthümer“ (von Jacob Grimm, 1840) <https://github.com/UB-Mannheim/Weisthuemer>
- 2011–2018: Digitalisierung und OCR der Sammlung „Desbillons“ (Reiseberichte & juristische Texte)



# Der Weg zum nutzbaren Volltext. Werkspezifisches Training als Baustein der OCR-Volltexterkennung für Alte Drucke



1. OCR und werkspezifisches Training an der UB Mannheim
- 2. Einführung**
3. Training mit synthetischer Ground Truth (Ansatz A)
4. Training mit realer Ground Truth (Ansatz B)
5. Fazit

## 2. Einführung

Im Ganzen waren im Jahre 1877 815272 Personen mit 2 511 478 651 *M.* versichert; hiervon schieden 13 109 Personen (35 906 127 *M.*) durch Tod und 48 754 Personen (137 760 445 *M.*) anderweitig aus, so daß am Jahresluß 753 409 Personen mit 2 337 812 074 *M.* versichert blieben, davon 552 246 Personen mit 1 847 622 742 *M.* bei den 35 Gesellschaften im Deutschen Reich. Der gesammte Zuwachs im Jahre 1877 stellt sich auf 15 777 Personen mit 110 873 820 *M.*, d. h. 2,14 % der Zahl der Versicherten und 4,98 % der Versicherungssumme. Die Berücksichtigung der Unfälle, welche die Liquidation einiger Gesellschaften im Jahre 1877 zur Folge gehabt hat, vermindert sich der Zuwachs auf 8542 (1,16 %) der Versicherten und 93 396 677 *M.* (4,19 %) der Summe (gegen 4,35 bzw. 6,57 % in 1876; 5,30 bzw. 8,09 % in 1875; 7,12 bzw. 9,29 % in 1874; 8,84 bzw. 11,44 % in 1873). Der Durchschnitt der Summe belief sich pro Kopf im Jahre 1877 auf 3103 *M.*, gegen 3021 *M.* in 1876. Den stärksten Zuwachs hatten im Jahre 1877 von den 35 Gesellschaften im Deutschen Reich diejenigen zu Karlsruhe (2 810 616 *M.*) und Gotha (1 167 800 *M.*), außerdem hatten noch 4 dieser Gesellschaften einen Zuwachs, alle übrigen eine Abnahme der Versicherungssumme erlitten.

*Deutscher Reichsanzeiger und Preußischer  
Staatsanzeiger, Jg. 1878, Ausgabe 248*

## 2. Einführung

Im Ganzen waren im Jahre 1877 815272 Personen mit 2511478651 M. versichert; hiervon schieden 13109 Personen (35906127 M.) durch Tod und 48754 Personen (137760445 M.) anderweitig aus, so daß am Jahreschluß 753409 Personen mit 2337812074 M. versichert blieben, davon 552246 Personen mit 1847622742 M. bei den 35 Gesellschaften im Deutschen Reich. Der gesammte Zuwachs im Jahre 1877 stellt sich auf 15777 Personen mit 110873820 M., d. h. 2,14% der Zahl der Versicherten und 4,98% der Versicherungssumme. Die Berücksichtigung der Unfälle, welche die Liquidation einiger Gesellschaften im Jahre 1877 zur Folge gehabt hat, vermindert sich der Zuwachs auf 8542 (1,16%) der Versicherten und 93396677 M. (4,19%) der Summe (gegen 4,35 bzw. 6,57% in 1876; 5,30 bzw. 8,09% in 1875; 7,12 bzw. 9,29% in 1874; 8,84 bzw. 11,44% in 1873). Der Durchschnitt der Summe belief sich pro Kopf im Jahre 1877 auf 3103 M., gegen 3021 M. in 1876. Den stärksten Zuwachs hatten im Jahre 1877 von den 35 Gesellschaften im Deutschen Reich diejenigen zu Karlsruhe (2810616 M.) und Gotha (1167800 M.), außerdem hatten noch 4 dieser Gesellschaften einen Zuwachs, alle übrigen eine Abnahme der Versicherungssumme erlitten.

*Deutscher Reichsanzeiger und Preußischer Staatsanzeiger, Jg. 1878, Ausgabe 248*

Im Ganzen waren im Jahre 1877 815272 Personen mit 2511478651 -ÄZ erfichert; hiervon schieden 13109 Personen (35906127 ii) durch Tod und 48754 11 (137760445 kii) anderweitig aus, so daß am Jahreschluß 753409 Per= fonen mit 2337812074 &a versichert blieben, davon 552246 Personen mit 1847622742 bei den 35 re 157 fteh im Deutschen Reich. Der gesammte Zuwachs im Jahre 1877 ftellt sich auf 15777 Perionen mit 110873820 , d. h. 2,14 %0 der Zahl der Versicherten und 4,98 % der Wunche die Summe. Die Berücksichtigung der Unfälle, welche die Liquidation einiger Gesellschaften im Jahre 1877 zur Folge gehabt hat, vermindert sich der Zuwachs auf 8542 (1,16%) der Versicherten und 93396677 a (4,19%) der Summe (gegen 4,35 bzw. 6,57 / in 1876; 5,30 bzw. 8,09 % in 1875; 7,12 bzw. 9,29 % in 1874; 8,84 bzw. 11,44% in 187, Der Durchschnitt der Summe belief sich pro Kopf im Jahre 1877 auf 3103 a, gegen 302 s in 1876. Den stärksten Zuwachs hatten im Jahre 1877 von den 35 Gefell= schaften im Deutschen Reich diejenigen zu Karlsruhe (2810616 und Gotha (1167800 ), außerdem hatten noch 4 dieser Gefell= schaften einen Zuwachs, alle übrigen eine Abnahme der Versiche= rungssumme erlitten.

OCR-Ergebnis

Tesseract (tesseract 5.1.0-20)

Tesseract-Model: [frak2021](#) (UB Mannheim)

## 2. Einführung

Im Ganzen waren im Jahre 1877 815272 Personen mit 2 511 478 651 *M.* versichert; hiervon schieden 13 109 Personen (35 906 127 *M.*) durch Tod und 48 754 Personen (137 760 445 *M.*) anderweitig aus, so daß am Jahreschluß 753 409 Personen mit 2 337 812 074 *M.* versichert blieben, davon 552 246 Personen mit 1 847 622 742 *M.* bei den 35 Gesellschaften im Deutschen Reich. Der gesammte Zuwachs im Jahre 1877 stellt sich auf 15 777 Personen mit 110 873 820 *M.*, d. h. 2,14 % der Zahl der Versicherten und 4,98 % der Versicherungssumme. Die Berücksichtigung der Unfälle, welche die Liquidation einiger Gesellschaften im Jahre 1877 zur Folge gehabt hat, vermindert sich der Zuwachs auf 8542 (1,16 %) der Versicherten und 93 396 677 *M.* (4,19 %) der Summe (gegen 4,35 bzw. 6,57 % in 1876; 5,30 bzw. 8,09 % in 1875; 7,12 bzw. 9,29 % in 1874; 8,84 bzw. 11,44 % in 1873). Der Durchschnitt der Summe belief sich pro Kopf im Jahre 1877 auf 3103 *M.* gegen 3021 *M.* in 1876. Den stärksten Zuwachs hatten im Jahre 1877 von den 35 Gesellschaften im Deutschen Reich diejenigen zu Karlsruhe (2 810 616 *M.*) und Gotha (1 167 800 *M.*), außerdem hatten noch 4 dieser Gesellschaften einen Zuwachs, alle übrigen eine Abnahme der Versicherungssumme erlitten.

*Deutscher Reichsanzeiger und Preußischer Staatsanzeiger*, Jg. 1878, Ausgabe 248

Im Ganzen waren im Jahre 1877 815272 Personen mit 2511478 651 *-ÄZ* erfichert; hiervon schieden 13 109 Personen (35 906 127 *ii*) durch Tod und 48 754 11 (137 760 445 *κii*) anderweitig aus, so daß am Jahreschluß 753 409 Personen mit 2337 812 074 *&* erfichert blieben, davon 552 246 Personen mit 1847 622 742 bei den 35 re 157 fteh im Deutschen Reich. Der gesammte Zuwachs im Jahre 1877 ftellt sich auf 15 777 Perionen mit 110 873 820 , d. h. 2,14 % der Zahl der Versicherten und 4,98 % der Wunche die Summe. Die Berücksichtigung der Unfälle, welche die Liquidation einiger Gesellschaften im Jahre 1877 zur Folge gehabt hat, vermindert sich der Zuwachs auf 8542 (1,16%) der Versicherten und 93 396 677 *α* (4,19%) der Summe (gegen 4,35 bzw. 6,57 / in 1876; 5,30 bzw. 8,09 % in 1875; 7,12 bzw. 9,29 % in 1874; 8,84 bzw. 11,44% in 187, Der Durchschnitt der Summe belief sich pro Kopf im Jahre 1877 auf 3103 *α*, gegen 302 *σ* in 1876. Den stärksten Zuwachs hatten im Jahre 1877 von den 35 Gesellschaften im Deutschen Reich diejenigen zu Karlsruhe (2810 616 *α*) und Gotha (1 167 800 *α*), außerdem hatten noch 4 dieser Gesellschaften einen Zuwachs, alle übrigen eine Abnahme der Versicherungssumme erlitten.

OCR-Ergebnis

Tesseract (tesseract 5.1.0-20)

Tesseract-Model: *frak2021* (UB Mannheim)



## 2. Einführung

Im Ganzen waren im Jahre 1877 815272 Personen mit 2 511 478 651 *M.* versichert; hiervon schieden 13 109 Personen (35 906 127 *M.*) durch Tod und 48 754 Personen (137 760 445 *M.*) anderweitig aus, so daß am Jahreschluß 753 409 Personen mit 2 337 812 074 *M.* versichert blieben, davon 552 246 Personen mit 1 847 622 742 *M.* bei den 35 Gesellschaften im Deutschen Reich. Der gesammte Zuwachs im Jahre 1877 stellt sich auf 15 777 Personen mit 110 873 820 *M.*, d. h. 2,14 % der Zahl der Versicherten und 4,98 % der Versicherungssumme. Die Berücksichtigung der Unfälle, welche die Liquidation einiger Gesellschaften im Jahre 1877 zur Folge gehabt hat, vermindert sich der Zuwachs auf 8542 (1,16 %) der Versicherten und 93 396 677 *M.* (4,19 %) der Summe (gegen 4,35 bzw. 6,57 % in 1876; 5,30 bzw. 8,09 % in 1875; 7,12 bzw. 9,29 % in 1874; 8,84 bzw. 11,44 % in 1873). Der Durchschnitt der Summe belief sich pro Kopf im Jahre 1877 auf 3103 *M.* gegen 3021 *M.* in 1876. Den stärksten Zuwachs hatten im Jahre 1877 von den 35 Gesellschaften im Deutschen Reich diejenigen zu Karlsruhe (2 810 616 *M.*) und Gotha (1 167 800 *M.*), außerdem hatten noch 4 dieser Gesellschaften einen Zuwachs, alle übrigen eine Abnahme der Versicherungssumme erlitten.

*Deutscher Reichsanzeiger und Preußischer Staatsanzeiger*, Jg. 1878, Ausgabe 248

Im Ganzen waren im Jahre 1877 815272 Personen mit 2511 478 651 *M.* versichert; hiervon schieden 13 109 Personen (35 906 127 *M.*) durch Tod und 48 754 14. (137 760 445 *M.*) anderweitig aus, so daß am Jahreschluß 753 409 Personen mit 2 337 812 074 *M.* versichert blieben, davon 552 246 Personen mit 1 847 622 742 *M.* bei den 35 Se im Deutschen Reich. Der gesammte Zuwachs im Jahre 1877 stellt sich auf 15777 Personen mit 110 873 820 *M.*, d. h. 2,14 % der Zahl der Versicherten und 4,98 % der Versicherungssumme. Die Berücksichtigung der Unfälle, welche die Liquidation einiger Gesellschaften im Jahre 1877 zur Folge gehabt hat, vermindert sich der Zuwachs auf 8542 (1,16 %) der Versicherten und 93 396 677 *M.* (4,19 %) der Summe (gegen 4,35 bzw. 6,57 % in 1876; 5,30 bzw. 8,09 % in 1875; 7,12 bzw. 9,29 in 1874; 8,84 bzw. 11,44 % in 71. Der Durchschnitt der Summe belief sich pro Kopf im Jahre 1877 auf 3103 *M.*, gegen 3021 *M.* in 1876. Den stärksten Zuwachs hatten im Jahre 1877 von den 35 Gesellschaften im Deutschen Reich diejenigen zu Karlsruhe (2 810 616 *M.*) und Gotha (1 167 800 *M.*), außerdem hatten noch 4 dieser Gesellschaften einen Zuwachs, alle übrigen eine Abnahme der Versicherungssumme erlitten.

OCR-Ergebnis

Tesseract (tesseract 5.1.0-20)

[Tesseract-Model: frak2021 \(UB Mannheim\) nach werkspezifischem Training](#)

## 2. Einführung



Gründe für ein werkspezifisches Training ...

## 2. Einführung

Gründe für ein werkspezifisches Training ...

- Verbesserung der Texterkennungsqualität

## 2. Einführung

Gründe für ein werkspezifisches Training ...

- Verbesserung der Texterkennungsqualität
- Erweiterung des zu erkennenden Zeichenvorrates → historische Währungssymbole, astronomische oder mathematische Zeichen etc.

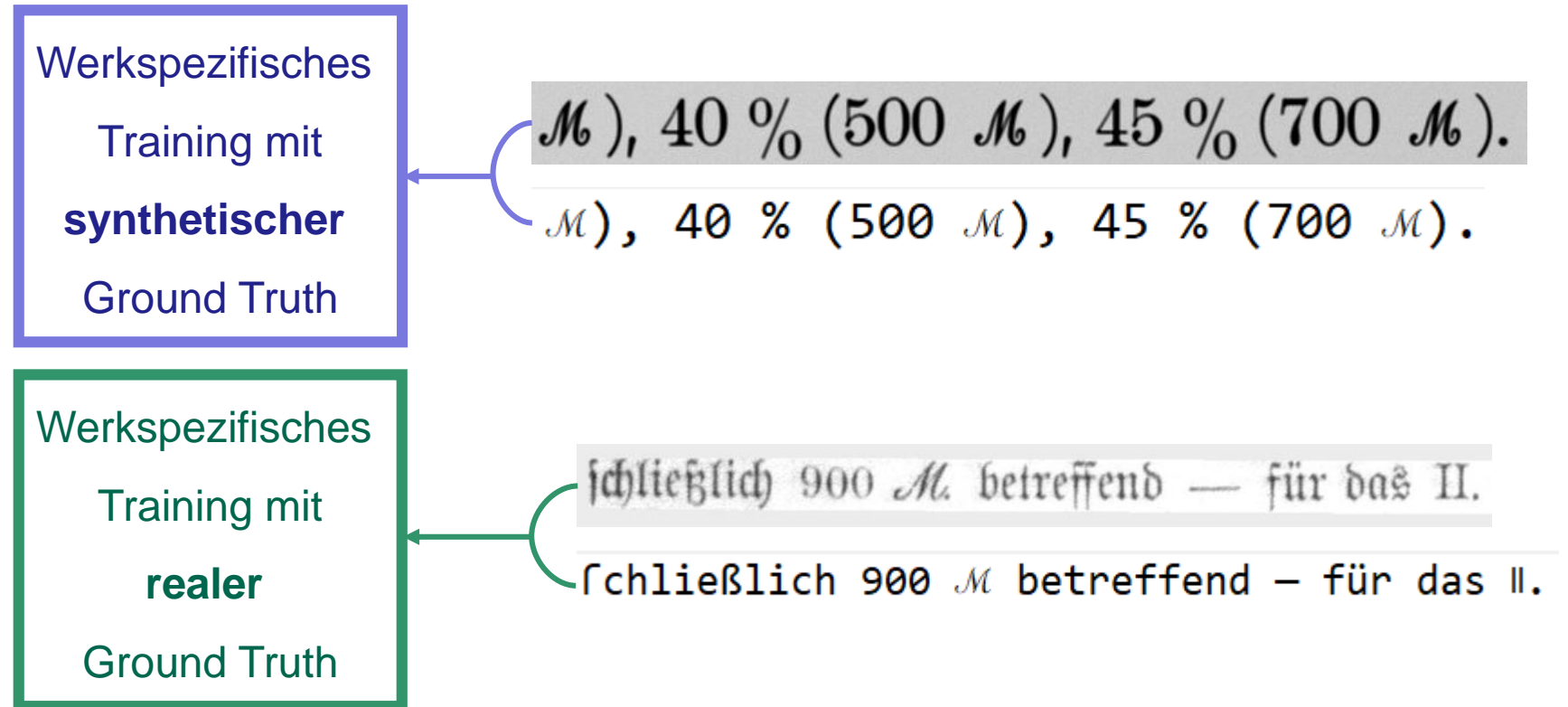
## 2. Einführung

Gründe für ein werkspezifisches Training ...

- Verbesserung der Texterkennungsqualität
- Erweiterung des zu erkennenden Zeichenvorrates → historische Währungssymbole, astronomische oder mathematische Zeichen etc.
- Spezielle Mappings → Normalisierungen, Auflösung von Kürzeln, Umgang mit Leerräumen (gesperrte Schrift) etc.

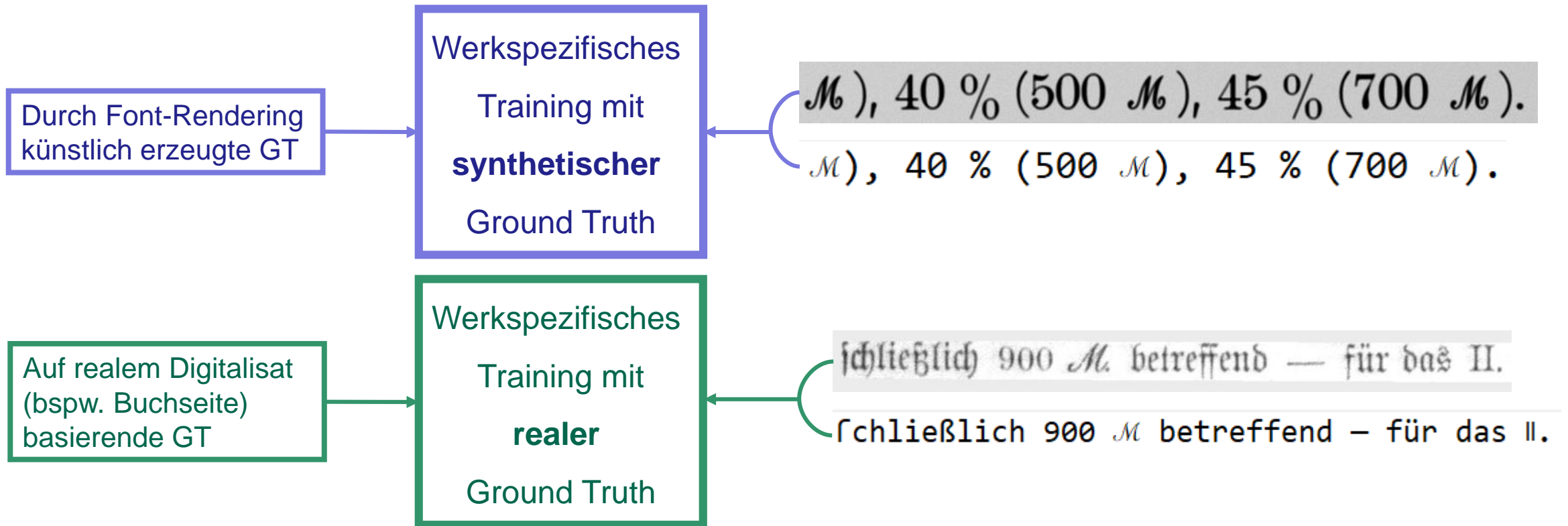
## 2. Einführung

Zwei getestete Ansätze für werkspezifisches Training:



## 2. Einführung

Zwei getestete Ansätze für werkspezifisches Training:



# Der Weg zum nutzbaren Volltext. Werkspezifisches Training als Baustein der OCR-Volltexterkennung für Alte Drucke



1. OCR und werkspezifisches Training an der UB Mannheim
2. Einführung
3. Training mit synthetischer Ground Truth (Ansatz A)
4. Training mit realer Ground Truth (Ansatz B)
5. Fazit



### 3. Training mit synthetischer Ground Truth (Ansatz A)

- Tesseract-Tool **Text2Image** erzeugt aus Textinput und Font eine Bilddatei

### 3. Training mit synthetischer Ground Truth (Ansatz A)

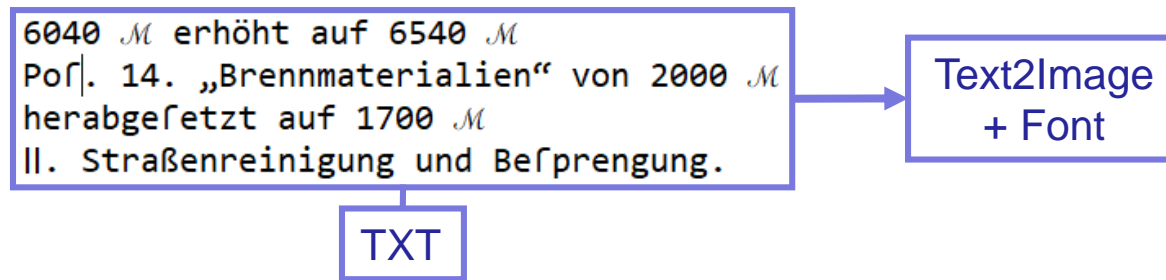
- Tesseract-Tool **Text2Image** erzeugt aus Textinput und Font eine Bilddatei

6040 *ℳ* erhöht auf 6540 *ℳ*  
Poſ|. 14. „Brennmaterialien“ von 2000 *ℳ*  
herabgefetzt auf 1700 *ℳ*  
II. Straßenreinigung und Befprengung.

TXT

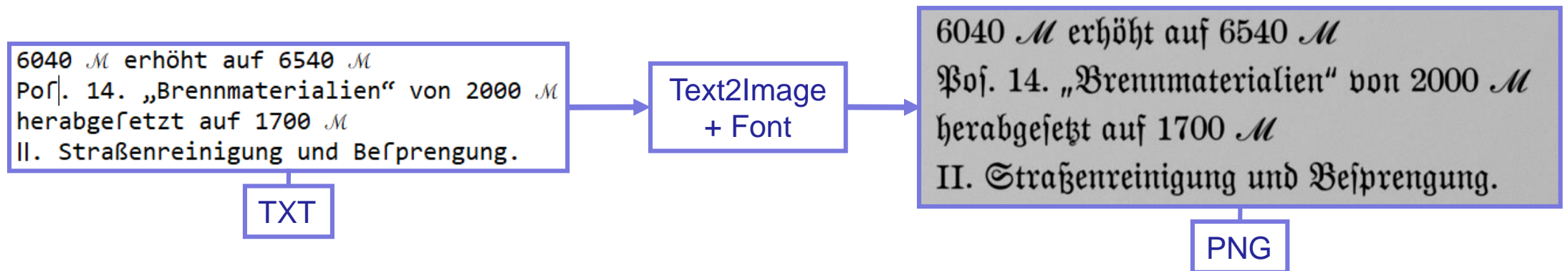
### 3. Training mit synthetischer Ground Truth (Ansatz A)

- Tesseract-Tool **Text2Image** erzeugt aus Textinput und Font eine Bilddatei



### 3. Training mit synthetischer Ground Truth (Ansatz A)

- Tesseract-Tool **Text2Image** erzeugt aus Textinput und Font eine Bilddatei



### 3. Training mit synthetischer Ground Truth (Ansatz A)

2 verwendete Fraktur-Fonts für Rendering:

Gründerfrach=Fraktur	Berthold Mainzer Fraktur
<b>entspricht</b> dem verwendeten Font des zu erkennenden Ausgangsmaterials	<b>ähnelt</b> dem Font des zu erkennenden Ausgangsmaterial
hoher Zeichenvorrat: langes f, rundes r, M etc.	hoher Zeichenvorrat: langes f, rundes r, M etc.

*Tabelle 1: Verwendete Fonts*

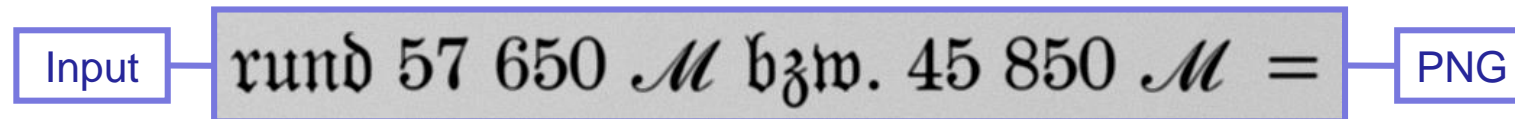
### 3. Training mit synthetischer Ground Truth (Ansatz A)



- Zusätzliches Skript, um Bildzeilen per data augmentation zu verändern und die GT mit Variationen anzureichern

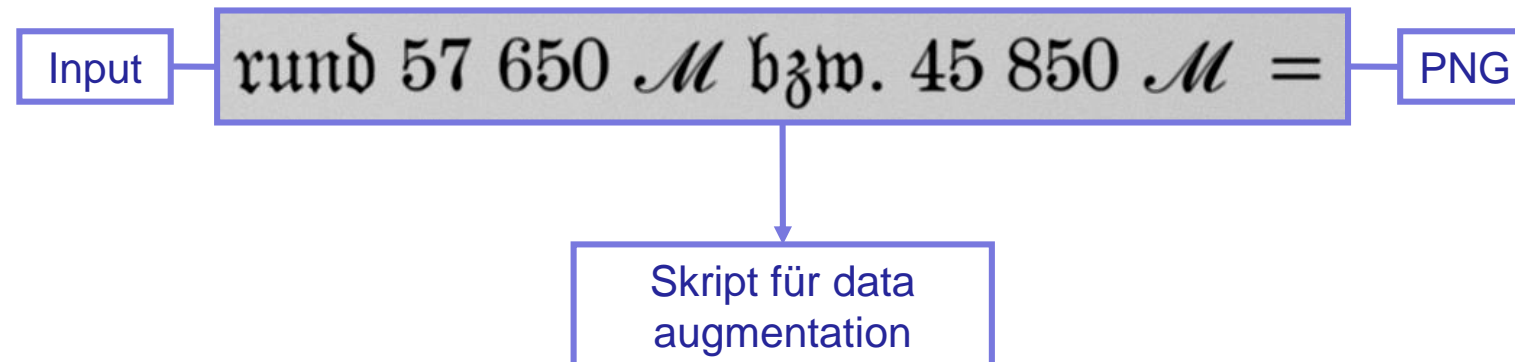
### 3. Training mit synthetischer Ground Truth (Ansatz A)

- Zusätzliches Skript, um Bildzeilen per data augmentation zu verändern und die GT mit Variationen anzureichern



### 3. Training mit synthetischer Ground Truth (Ansatz A)

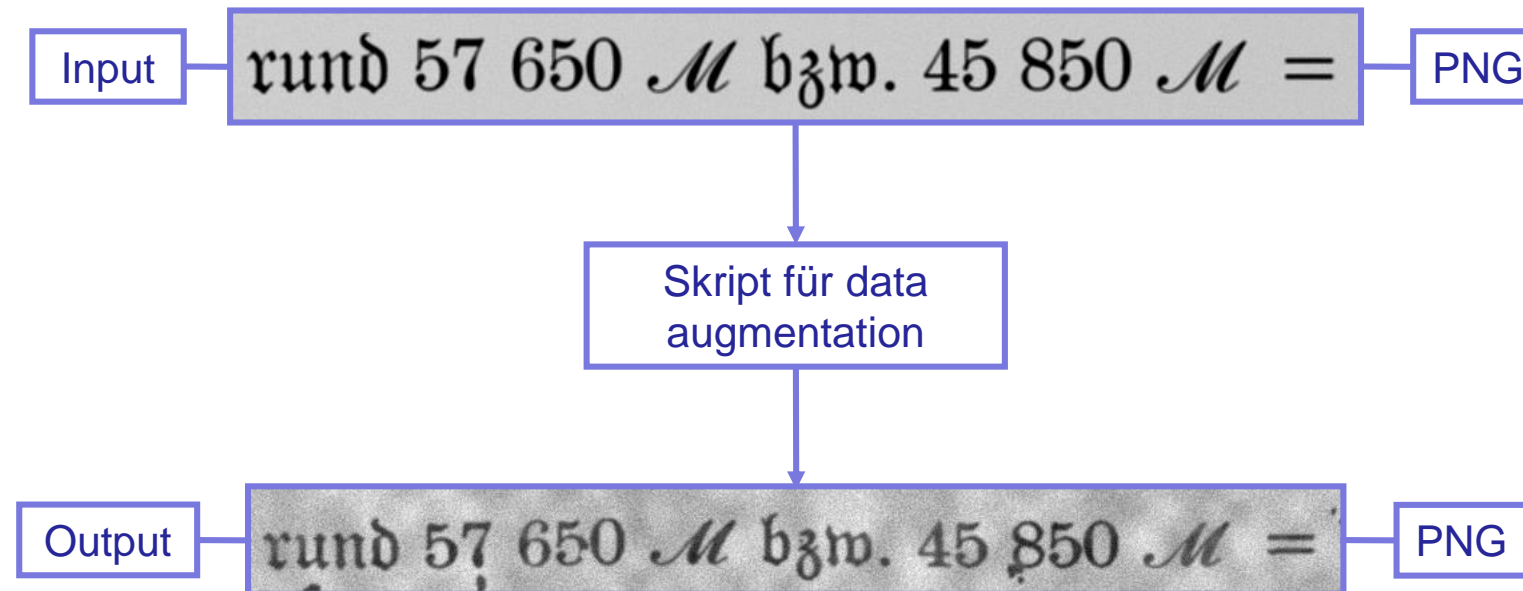
- Zusätzliches Skript, um Bildzeilen per data augmentation zu verändern und die GT mit Variationen anzureichern





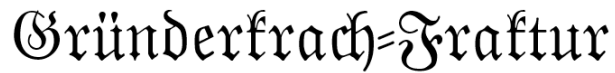

### 3. Training mit synthetischer Ground Truth (Ansatz A)

- Zusätzliches Skript, um Bildzeilen per data augmentation zu verändern und die GT mit Variationen anzureichern



### 3. Training mit synthetischer Ground Truth (Ansatz A)

Tabelle 2: Trainingssets mit synthetischer Ground Truth

Werkspezifisches Training mit synthetischer Ground Truth (2 Fonts / 6 Datasets)			
Ausgangsmodell:	frak2021 (UB Mannheim)		
Textquelle:	„Charlottenburger Amtsschrifttum“, 1879–1919 (ZLB)		
Fonts:	 © 2021 Philipp Poll	 © Peter Wiegel	
Synthetische GT:	<b>LARGE*</b> 3.000 GT lines	<b>MEDIUM†</b> 1.000 GT lines	<b>SMALL†</b> 100 GT lines
	PNG, Graustufe, 50 % per data augmentation bearbeitet		
Trainingsziel:	Erweiterung des Zeichenvorrats um: $\mathcal{M}$		
Trainingslaufzeit:	60 Epochen mit Abbruchkriterium → Tatsächliche Laufzeit ca. 40–50 Epochen		

\* enthält Lines mit und ohne  $\mathcal{M}$ ; † enthält nur Lines mit  $\mathcal{M}$

### 3. Training mit synthetischer Ground Truth (Ansatz A)

Evaluation: *Gründerkrach=Fraktur*

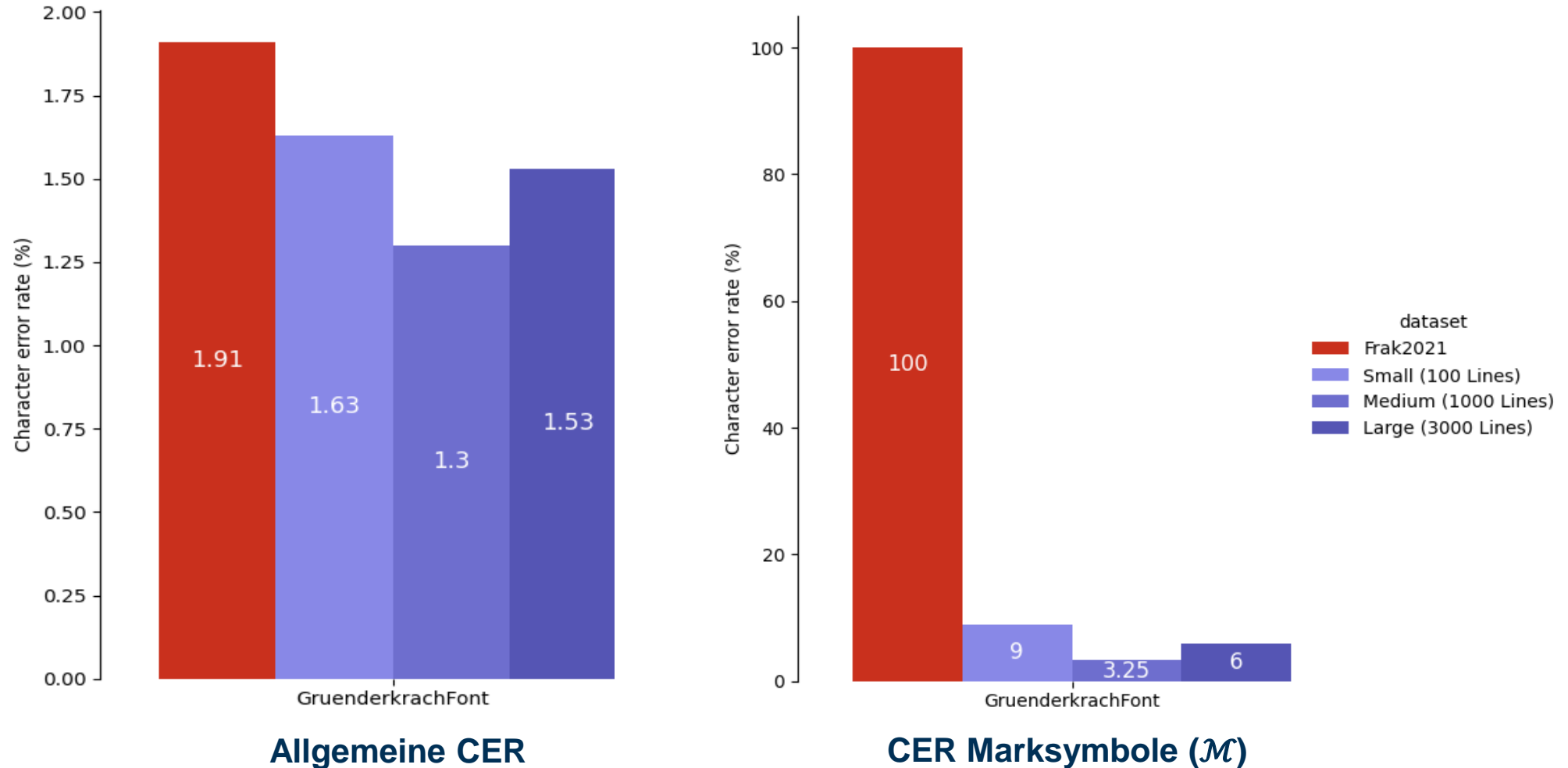


Fig. 1

### 3. Training mit synthetischer Ground Truth (Ansatz A)

#### Evaluation: **Berthold Mainzer Fraktur**

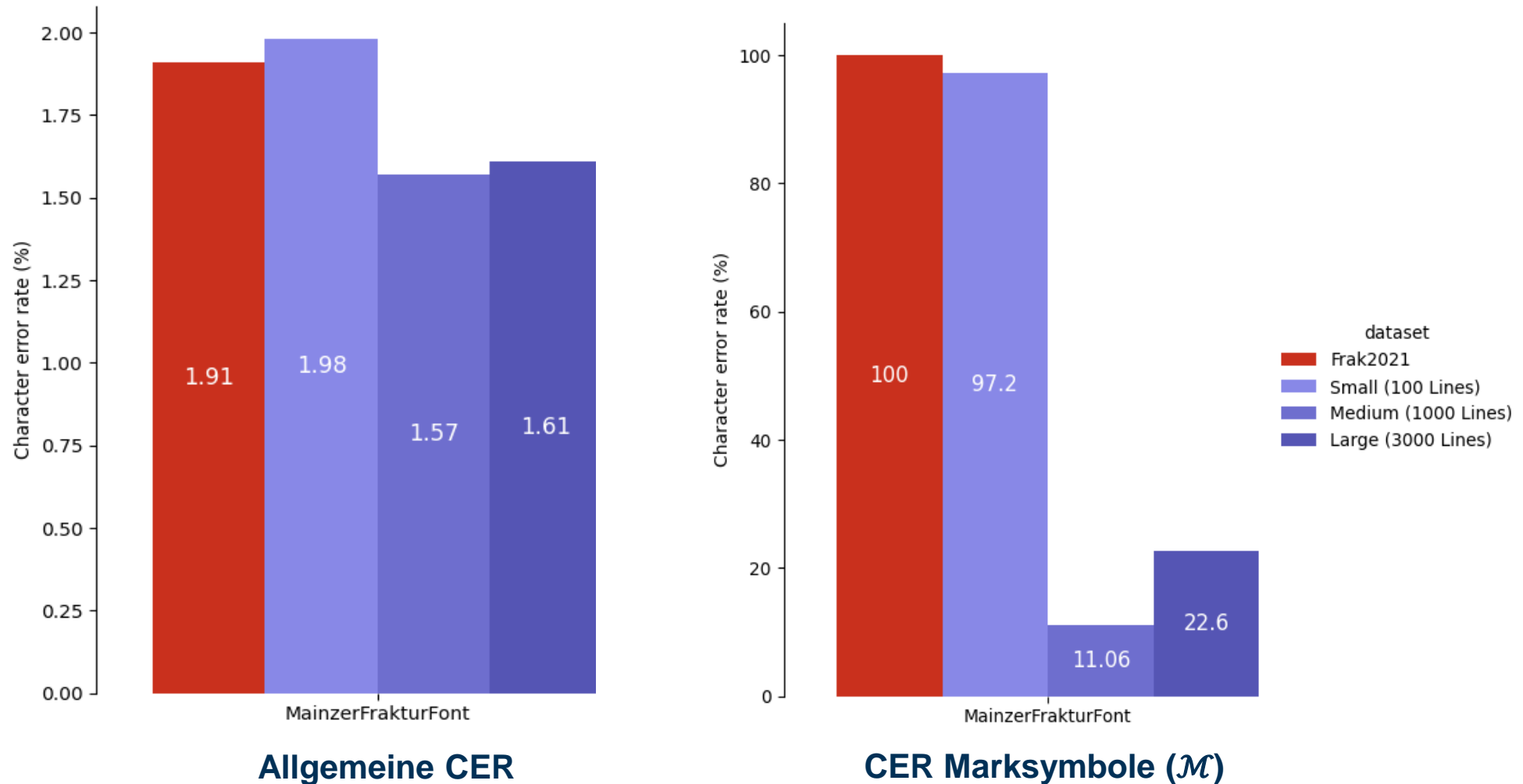


Fig. 2

# Der Weg zum nutzbaren Volltext. Werkspezifisches Training als Baustein der OCR-Volltexterkennung für Alte Drucke



1. OCR und werkspezifisches Training an der UB Mannheim
2. Einführung
3. Training mit synthetischer Ground Truth (Ansatz A)
4. Training mit realer Ground Truth (Ansatz B)
5. Fazit

## 4. Training mit realer Ground Truth (Ansatz B)

- Analog zum **synthetischen** Training wurde als Quelle der GT-Produktion das „Charlottenburger Amtsschrifttum“, 1879–1919 (Zentrale Landesbibliothek Berlin) verwendet.

## 4. Training mit realer Ground Truth (Ansatz B)

- Analog zum **synthetischen** Training wurde als Quelle der GT-Produktion das „Charlottenburger Amtsschrifttum“, 1879–1919 (Zentrale Landesbibliothek Berlin) verwendet.
- Die Transkription des Quellmaterials erfolgte in **Transkribus**.

## 4. Training mit realer Ground Truth (Ansatz B)

- Analog zum **synthetischen** Training wurde als Quelle der GT-Produktion das „Charlottenburger Amtsschrifttum“, 1879–1919 (Zentrale Landesbibliothek Berlin) verwendet.
- Die Transkription des Quellmaterials erfolgte in **Transkribus**.
- Im Gegensatz zum **synthetischen** Training erfolgte als data augmentation nur eine Graustufennormalisierung für die Bildinhalte.



## 4. Training mit realer Ground Truth (Ansatz B)

- Analog zum **synthetischen** Training wurde als Quelle der GT-Produktion das „Charlottenburger Amtsschrifttum“, 1879–1919 (Zentrale Landesbibliothek Berlin) verwendet.
- Die Transkription des Quellmaterials erfolgte in **Transkribus**.
- Im Gegensatz zum **synthetischen** Training erfolgte als data augmentation nur eine Graustufennormalisierung für die Bildinhalte.
- Sowohl Umfang der Ground Truth als auch Trainingslaufzeit (in Epochen) waren mit dem synthetischen Training identisch.

## 4. Training mit realer Ground Truth (Ansatz B)

Tabelle 3: Trainingsaufbau reale Ground Truth

Werkspezifisches Training mit realer Ground Truth (3 Datasets)			
Ausgangsmodell:	frak2021 (UB Mannheim)		
Quelle:	„Charlottenburger Amtsschrifttum“, 1881–1906 (ZLB)		
Ground Truth:	<b>LARGE*</b> 3.000 GT lines	<b>MEDIUM†</b> 1.000 GT lines	<b>SMALL†</b> 100 GT lines
	PNG, Graustufe (+ Graustufennormalisierung)		
Trainingsziel:	Erweiterung des Zeichenvorrats um: $\mathcal{M}$		
Trainingslaufzeit:	60 Epochen mit Abbruchkriterium → Tatsächliche Laufzeit ca. 40–50 Epochen		

\* enthält Lines mit und ohne  $\mathcal{M}$ ; † enthält nur Lines mit  $\mathcal{M}$

## 4. Training mit realer Ground Truth (Ansatz B)

Evaluation: reale Ground Truth

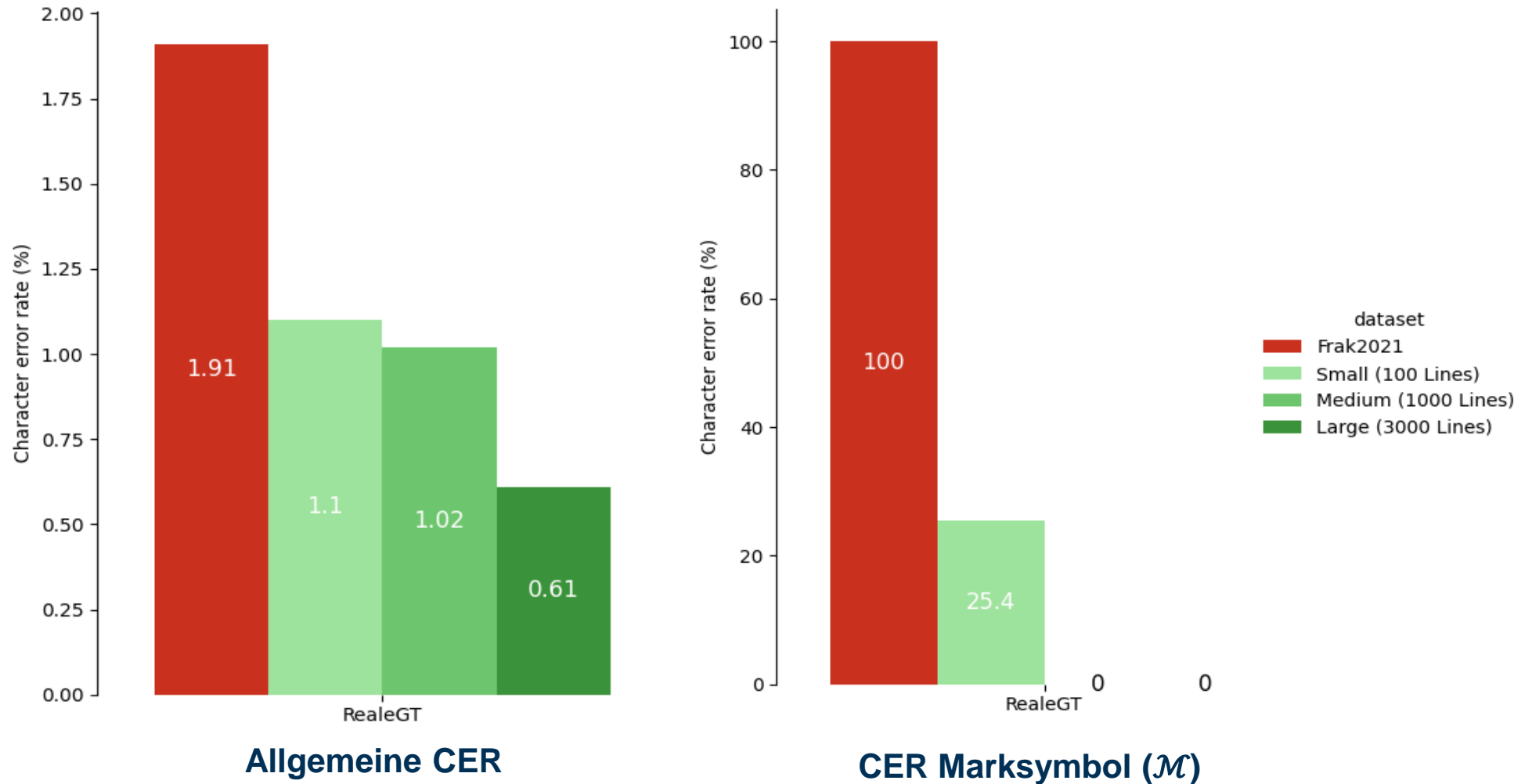


Fig. 3

# Der Weg zum nutzbaren Volltext. Werkspezifisches Training als Baustein der OCR-Volltexterkennung für Alte Drucke



1. OCR und werkspezifisches Training an der UB Mannheim
2. Einführung werkspezifisches Training
3. Training mit synthetischer Ground Truth (Ansatz A)
4. Training mit realer Ground Truth (Ansatz B)
5. **Fazit**

## 5. Fazit

### Allgemeine CER von realer und synthetischer GT

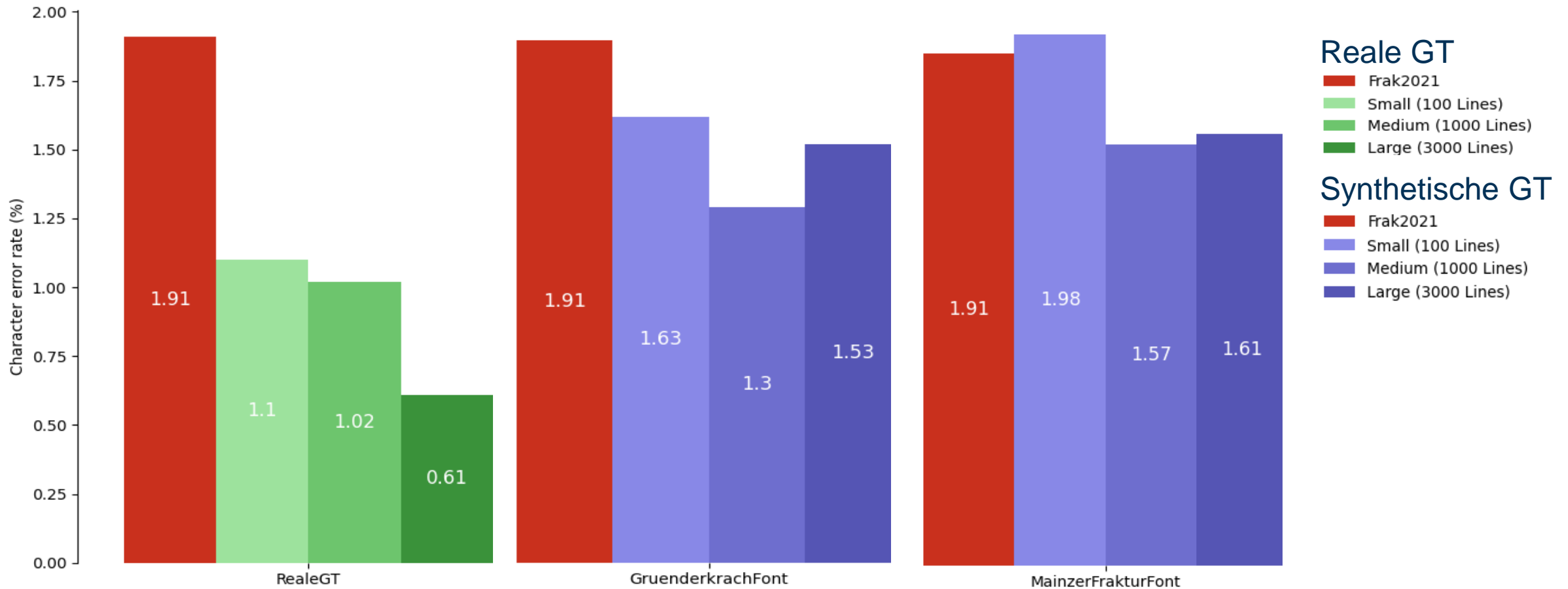


Fig. 4

## 5. Fazit

### CER Marksymbole ( $\mathcal{M}$ ) realer und synthetischer GT

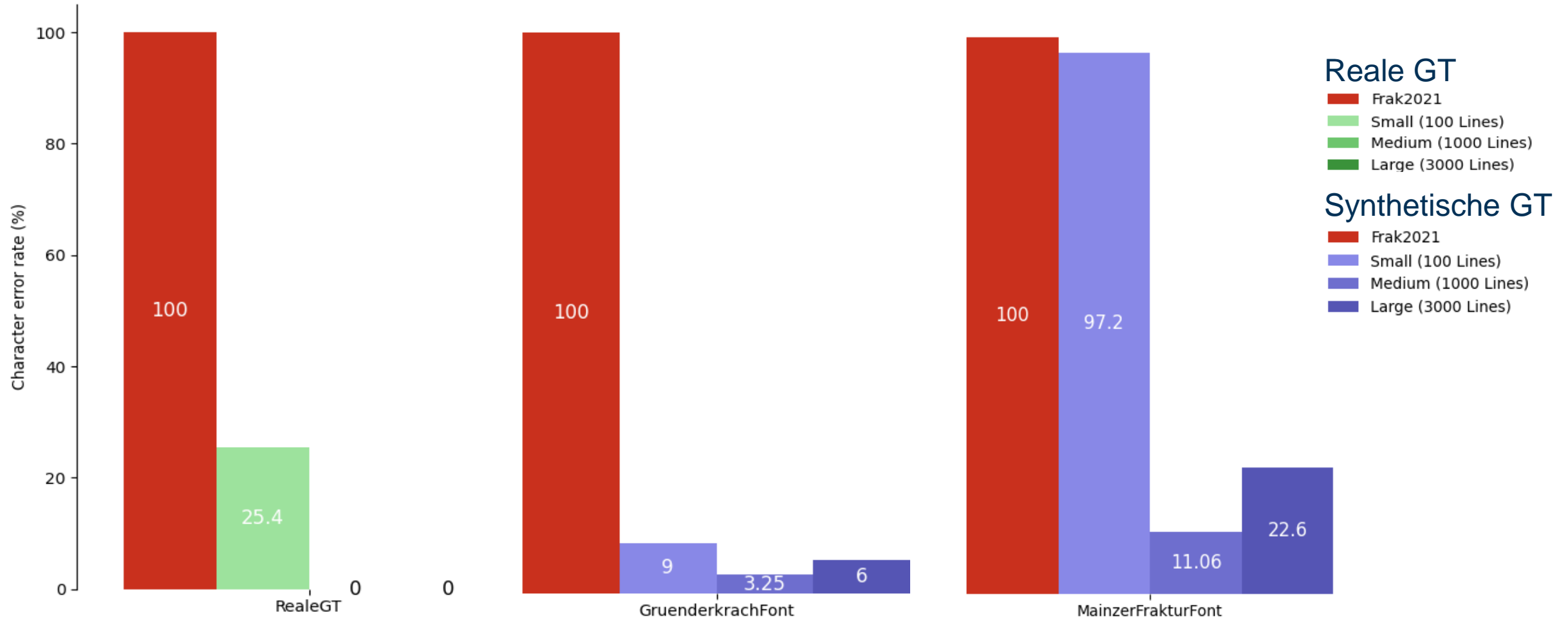


Fig. 5

## 5. Fazit

Training mit synthetischer Ground Truth		Training mit realer Ground Truth	
Vorteile	Nachteile	Vorteile	Nachteile
Geringer zeitlicher Aufwand zur Erzeugung umfangreicher GT	Fonts, die Ausgangsmaterial entsprechen, sind selten verfügbar	Sehr gute Verbesserung des Ausgangsmodells	Relativ hoher Aufwand zur Erstellung der GT (Auswahl entsprechender Textstellen mit infrage kommenden Glyphen + Transkription + QA)
Gute Ergebnisse, wenn der genutzte Font dem Ausgangsmaterial entspricht	Zeichensatz verfügbarer Fonts meist unzureichend (besonders bei seltenen Glyphen)	Nachtrainiertes Modell ist Modellen mit synthetischer GT überlegen	
	Seltene Glyphen in verschiedenen Fonts uneinheitlich codiert		

# Der Weg zum nutzbaren Volltext. Werkspezifisches Training als Baustein der OCR-Volltexterkennung für Alte Drucke



Ground Truth, weitere Auswertungen und zusätzliche Informationen auf Github:  
<https://github.com/UB-Mannheim/charlottenburger-amtsschrifttum>

Vielen Dank!

**Jan Kamlah (Entwicklung):** [jan.kamlah@bib.uni-mannheim.de](mailto:jan.kamlah@bib.uni-mannheim.de)

**Thomas Schmidt (Projektmanagement):** [thomas.schmidt@bib.uni-mannheim.de](mailto:thomas.schmidt@bib.uni-mannheim.de)



## 5. Fazit

- Werkspezifisches Training ist ein effektives Verfahren zur Optimierung der Texterkennungsausgabe
- In der Regel sollte das Training mit realen GT-Daten bevorzugt werden
- Training mit geringem Material starten und ggf. erweitern
- Data augmentation regularisiert das Training und reduziert die Gefahren des Overfittings
- Min. 40 Epochen trainieren
  - besser wären 100 Epochen mit einem Abbruchkriterium wie etwa 0.01 % CER
  - Nachträgliche Evaluation der Modelle um die Verwendung von Overfitting-Modelle zu vermeiden
- Kontrolle des Zeichenvorrates des Trainingsmaterials

## Literatur

- „Charlottenburger Amtsschrifttum“, 1879–1919 (ZLB),  
<https://digital.zlb.de/viewer/metadata/16046633/1/>
- BertholdMainzerFraktur, Peter Wiegel,  
<http://www.peter-wiegel.de/MainzerFraktur.html>
- Weil, Stefan: Neue Frakturmodelle für Tesseract, 2019,  
<https://madoc.bib.uni-mannheim.de/53748>
- Weil Stefan: 126 Jahre Zeitung online - Fundgrube für historisch Interessierte und Motor für die Bibliotheks-IT : 126 years of the newspaper online, 2018,  
<https://madoc.bib.uni-mannheim.de/46507/>
- Weil, Stefan und Kamlah, Jan: Forschungsdaten aus Digitalisaten, 2019,  
<https://madoc.bib.uni-mannheim.de/52204>