

Politechnika Wrocławska

Wydział Informatyki i Telekomunikacji

Kierunek: _____ Informatyka techniczna (ITE)
Specjalność: _____ Systemy informatyki w medycynie (IMT)

PRACA DYPLOMOWA

Inżynierska

**Implementacja systemu optycznego rozpoznawania
znaków oraz opracowanie środowiska eksperymentalnego
opartego o samodzielnie skonstruowany zbiór danych**

Krzysztof Zalewa

Opiekun pracy
Dr inż., Paweł Zyblewski

Słowa kluczowe: 3-6 słów kluczowych

WROCŁAW (2025)

Streszczenie

Dodaj streszczenie pracy w języku polskim. Staraj się uwzględnić wymienione na stronie tytułowej słowa kluczowe. Uwaga przedstawiony rekomendowany szablon dotyczy pracy dyplomowej pisanej w języku angielskim. W przeciwnym wypadku, student powinien samodzielnie zmienić nazwy „Chapter” na „Rozdział” itp stosując odpowiednie pakiety systemu L^AT_EX oraz ustawienia w pliku *latex-settings.tex*.

Abstract

Streszczenie w języku angielskim.

Spis treści

1	Wstęp	1
1.1	Opis problemu	1
1.2	Cel pracy	1
2	Przegląd literatury	3
2.1	Narzędzia OCR	3
2.1.1	Tesseract	3
2.1.2	Easy OCR	3
2.1.3	DocTR OCR	3
2.1.4	Paddle OCR	3
2.2	Zbiory danych	3
2.2.1	IAM	3
2.2.2	oldbooksdataset	5
2.3	Metryki	7
2.3.1	CER	7
2.3.2	WER	7
3	Aspekt inżynierski	9
3.1	Wybór algorytmów OCR	9
3.2	Akwizycja danych	9
4	Aspekt badawczy	11
4.1	Opis problemu	11
5	Podsumowanie	13

1. Wstęp

1.1. Opis problemu

1.2. Cel pracy

Test

Celem pracy jest porównanie, w oparciu o eksperymenty komputerowe, działania wybranych algorytmów optycznego rozpoznawania znaków. Kluczowym elementem pracy jest pozyskanie i opracowanie autorskiego zbioru danych, składającego się z treści dostępnych za pośrednictwem publicznie dostępnego API serwisu wolnelektury.pl, który pozwoli na rzetelną ewaluację znanych z literatury algorytmów OSR. Opracowany zbiór zawierał będzie dokumenty o zróżnicowanej charakterystyce, uwzględniając m.in. różne kroje i stopnie pisma, a także modyfikacje utrudniające poprawne odczytanie treści.

2. Przegląd literatury

2.1. Narzędzia OCR

2.1.1. Tesseract

Tesseract OCR to najstarszy z wybranych algorytmów optycznego rozpoznawania znaków. Został on stworzony przez firmę HP w latach 1984 - 1994. Algorytm ten działa w kilku fazach:

1. Analiza komponentów, gdzie zarys tych komponentów jest przechowywany. Takie podejście mimo że nakłada dodatkowe koszty obliczeniowe pozwala na łatwiejsze rozpoznawanie tekstu w odwróconych kolorach (biały tekst na czarnym tle) [3].
2. Wyszukiwanie linii w komponentach. Celem tego kroku była eliminacja potrzeby korekty przekrzywienia.
3. Podział linii na słowa.
4. Pierwsza iteracja rozpoznawania. Zaczynając na górze strony algorytm próbuje rozpoznać każde kolejne słowo. Jeżeli jest duże prawdopodobieństwo że słowo jest poprawne jest ono wykorzystywane do douczenia klasyfikatora. W ten sposób z każdym kolejnym słowem celność klasyfikatora powinna rosnąć.
5. Druga iteracja rozpoznawania. Po wykonaniu pierwszej iteracji jest duże prawdopodobieństwo że klasyfikator uzyskałby lepsze wyniki. Więc po raz drugi algorytm próbuje rozpoznać tekst na stronie i aktualizuje słowa które były mniej celnie rozpoznane.

2.1.2. Easy OCR

2.1.3. DocTR OCR

2.1.4. Paddle OCR

2.2. Zbiory danych

2.2.1. IAM

IAM to zbiór ręcznie zapisanych tekstów w języku angielskim. Wykonany przez Instytut matematyki i informatyki na Uniwersytecie Breńskim. [2] Zbiór zawiera obrazy w rozdzielczości 300dpi zapisane w formacie PNG w 256 odcieniach szarości. Każdy pod katalog zawiera teksty zapisane przez jedną osobę.

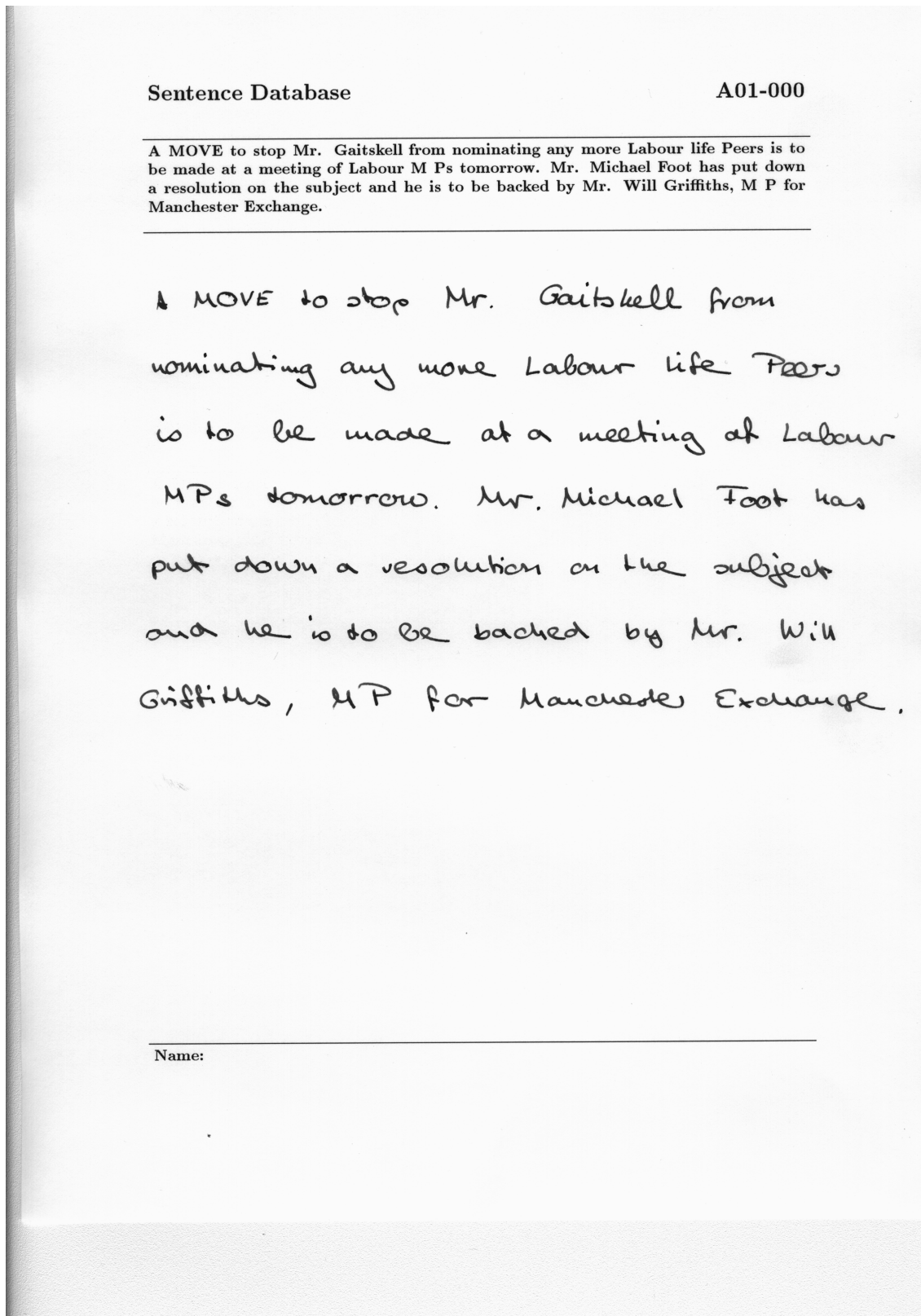


Figure 2.1: Przykładowy obraz ze zbioru danych IAM

2.2.2. oldbooksdataset

Zbiór udostępniony na platformie git hub zawierający skany książek w języku angielski. Książki zapisane są w formacie .tiff w rozdzielczości 300dpi oraz 500dpi [1].

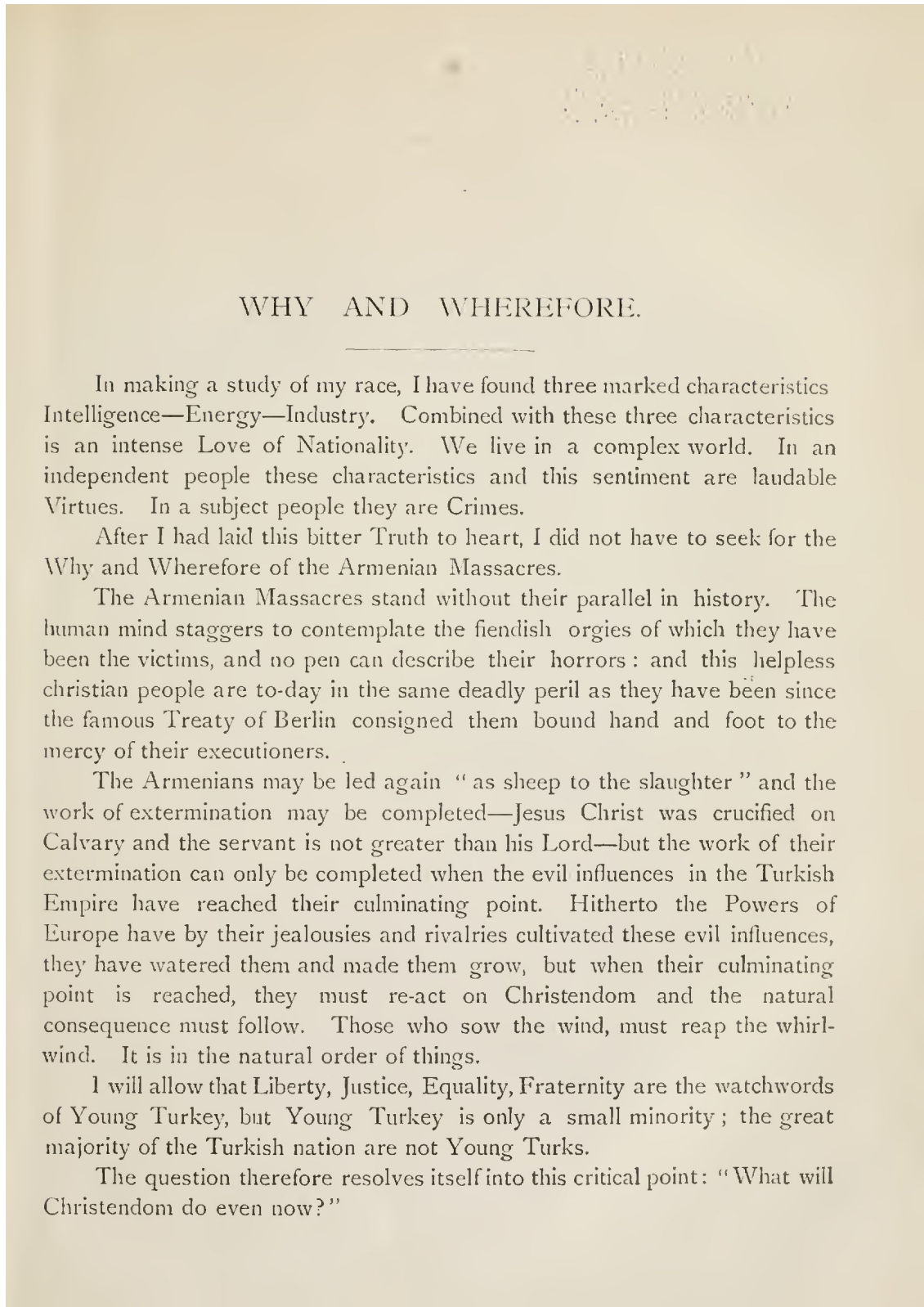


Figure 2.2: Przykładowy obraz ze zbioru danych old-books-dataset

2.3. Metryki

Do oceny wyników zastosowano dwie główne metryki. **CER** i **WER** zasadniczą różnicą między nimi jest celność porównań. CER porównuje na poziomie pojedynczych znaków natomiast, z kolei WER porównuje na poziomie poszczególnych słów. Z tego powodu CER jest przydatniejsze w kontekstach w których nawet pojedynczy znak może zmienić znaczenie słowa (np. w medycynie). Natomiast WER jest lepsze do porównywania spójności zdań itp.

2.3.1. CER

CER (Character Error Rate z ang. częstotliwość błędnych znaków) to metryka dzięki której możliwa jest ocena różnic między tekstem wytworzonym poprzez model OCR a tekstem rzeczywistym. W tym przypadku CER obliczane jest poprzez zsumowanie operacji (wstawień, usunięć oraz zamian znaków) potrzebnych do uzyskania tekstu rzeczywistego.

$$CER = \frac{S + D + I}{N_c}$$

Gdzie:

- S - Liczba zamian znaków (ang. Substitutions)
- D - Liczba usunięć znaków (ang. Deletions)
- I - Liczba wstawień znaków (ang. Inserts)
- N_c - Liczba znaków w tekście (ang. Number of characters)

Na przykład

Tekst oryginalny: Życiem wschód, śmiercią południe;

Tekst wygenerowany przez model: Życiem wschod, siercia poudniex;

Aby przekształcić tekst wygenerowany do tekstu oryginalnego należy wykonać 4 zamiany (Brakujące znaki polskie), 1 wstawienie (Brakujące 'm' w tekście wygenerowanym) oraz 1 usunięcie ('x' nie występuje w tekście oryginalnym). Więc $CER = 6/28 = 0.2141 \approx 21.4\%$

2.3.2. WER

WER (Word Error Rate z ang. częstotliwość błędnych słów) podobnie jak CER jest to metryka dzięki której możliwa jest ocena różnic między tekstem wytworzonym poprzez model OCR a tekstem rzeczywistym. Jak sama nazwa wskazuje WER porównuje tekst na poziomie poszczególnych słów.

$$WER = \frac{S + D + I}{N_w}$$

Gdzie:

- S - Liczba zamian słów (ang. Substitutions), czyli słowa które występują w obu tekstach ale te w tekście są różne od tych w tekście oryginalnym.
- D - Liczba usunięć słów (ang. Deletions), czyli słowa które występują w tekście oryginalnym jednakże nie ma ich w tekście wygenerowanym.
- I - Liczba wstawień słów (ang. Inserts), czyli słowa nadmiarowe których nie ma w tekście oryginalnym.
- N_w - Liczba słów w tekście (ang. Number of words)

Na przykład

Tekst oryginalny: Życiem wschód, śmiercią południe;

Tekst wygenerowany przez model: Życiem wschod, śmiercia poudniex;

Aby przekształcić tekst wygenerowany do tekstu oryginalnego należy wykonać 4 zamiany (Słowa zbliżone do oryginału ale nie takie same). Więc $WER = 4/4 = 1 = 100\%$

3. Aspekt inżynierski

Implementacja wybranych, znanych z literatury metod optycznego rozpoznawania znaków (ang. optical character recognition <OSR>). Opracowanie autorskiego zbioru danych – na podstawie treści samodzielnie pozyskanych z serwisu wolnelektury.pl z wykorzystaniem dostępnego publicznie API – na potrzeby ewaluacji zaimplementowanych algorytmów. Całość implementacji zostanie wykonana w języku Python z wykorzystaniem właściwie dobranych bibliotek programistycznych. Dodatkowo, na potrzeby ewaluacji, opracowane oraz zaimplementowane zostanie odpowiednie środowisko eksperymentalne.

3.1. Wybór algorytmów OCR

3.2. Akwizycja danych

4. Aspekt badawczy

4.1. Opis problemu

5. Podsumowanie

Bibliografia

- [1] P. Barcha. Old books dataset. <https://github.com/PedroBarcha/old-books-dataset>, 2024. Accessed: 2024.
- [2] U. Marti and H. Bunke. The iam-database: An english sentence database for off-line handwriting recognition. *International Journal on Document Analysis and Recognition*, 5:39–46, 2002.
- [3] R. Smith. An overview of the tesseract ocr engine. In *Ninth International Conference on Document Analysis and Recognition (ICDAR)*, volume 2, pages 629–633, 2007.

Spis ilustracji

2.1	Przykładowy obraz ze zbioru danych IAM	4
2.2	Przykładowy obraz ze zbioru danych old-books-dataset	6

Spis tabel