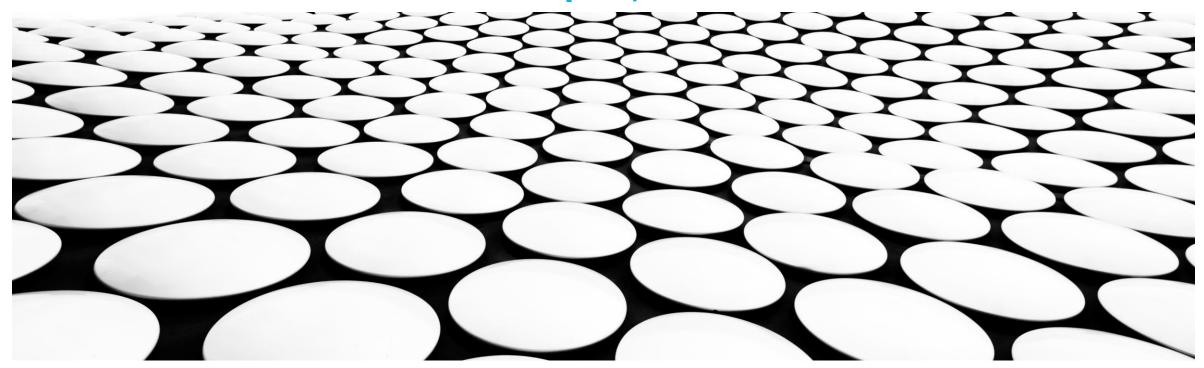
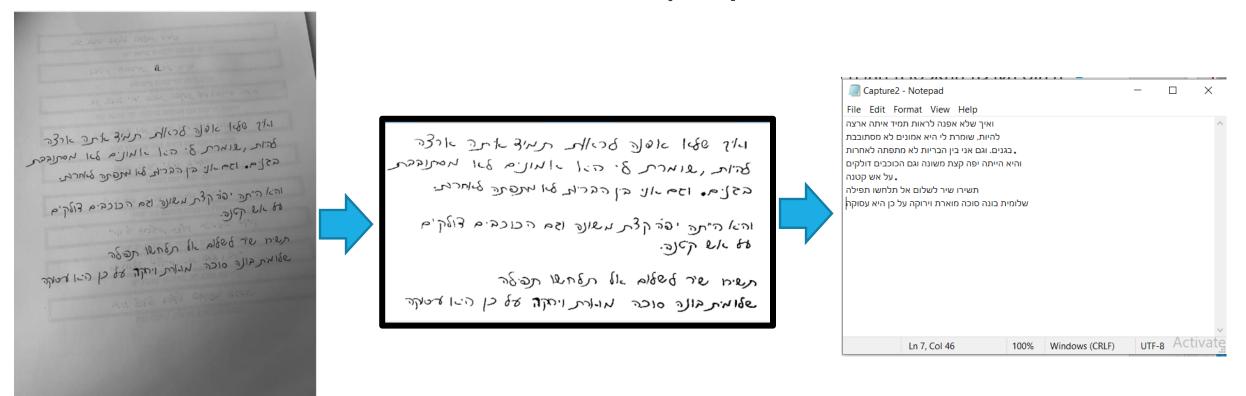
מתמונה לטקסט - חילוץ כתב יד בעברית SHECODES

עדי רוזנטל, אוקטובר 2020



מטרת הפרוייקט

מימוש מערכת המאפשרת המרה של תמונה המכילה טקסט בעברית (בכתב יד / דפוס) לכדי כתב מחשב אותו ניתן לערוך ולעצב בצורה נוחה.



2



















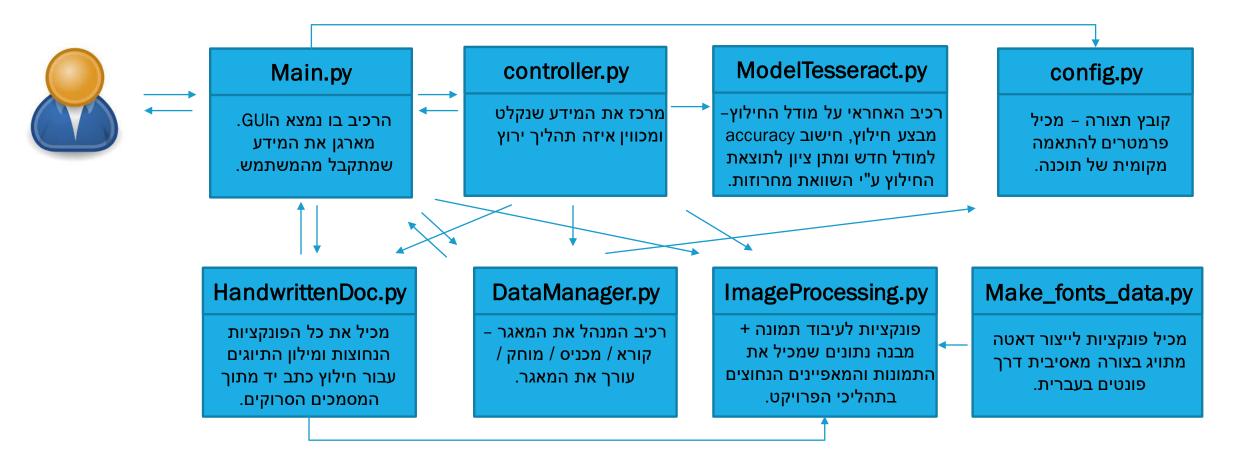


שלושה תהליכים עיקריים בפרוייקט:

- איסוף וחילוץ של דאטה מתוייג בכתב יד בעברית.1
- OCR המבוססת tesseract אימון רשת נוירונים.2
- 3. יצירת ממשק נוח למשתמש לחילוץ טקסט מתוך תמונה ואיסוף דאטה עבור אימון הרשת

4

ארכיטקטורה



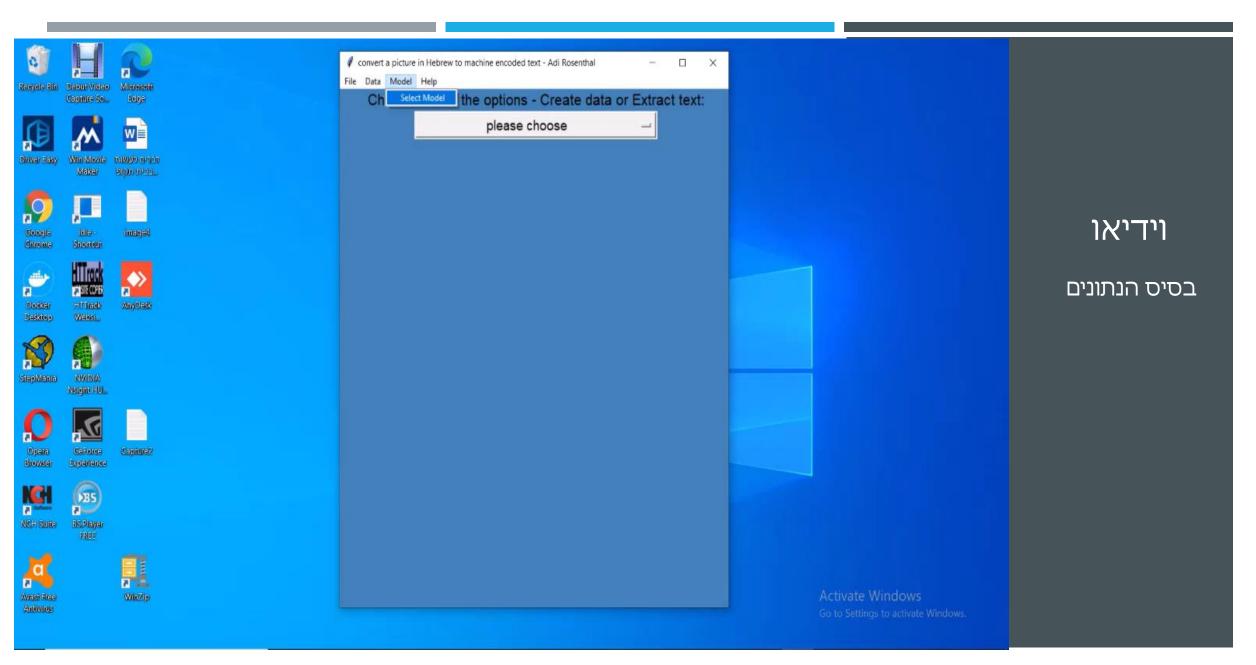
מבנה נתונים (DB)

תהליך טיוב המכונה לחילוץ טקסט דורש אימון של הרשת נוירונים עם כמות גדולה של דאטה.

- תצורת הדאטה: קובץ תמונה (TIF) + קובץ טקסט (.txt).
- הדאטה צריך להיות מאורגן כך שבכל תמונה מופיע שורה אחת של כתב יד + קובץ טקסט בעל
 שם זהה המכיל את טקסט שכתוב בתמונה

6

- הדאטה צריך להיות מדוייק ובאיכות גבוהה
- (overfitting כדאי לבצע גיוון בכתבי היד (כדי למנוע -

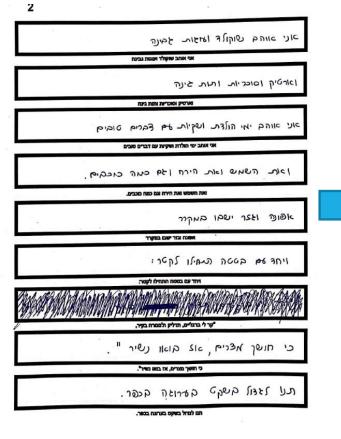


יצירה והבניה של מידע

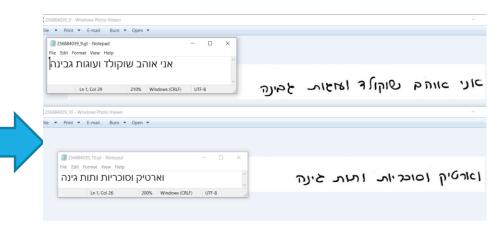
שלוש שיטות עיקריות בהן השתמשתי ליצירת דאטה בפרויקט:

- 1. סריקת וחילוץ של שורות טקסט מתוך דפים מוכנים ("Template") אליהם היו צריכים המשתמשים להעתיק שורות של שירים בכתב ידם לפי ההנחיות בדף.
 - 2. סריקה של תמונה עם טקסט ותיוג התמונה בתוך ממשק המערכת.
 - .3 איסוף פונטים בכתב יד בעברית וחילוץ שורות טקסט ותיוגם.

סריקת וחילוץ של שורות טקסט מתוך דפים מוכנים



9

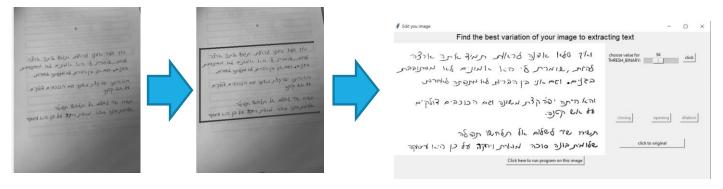


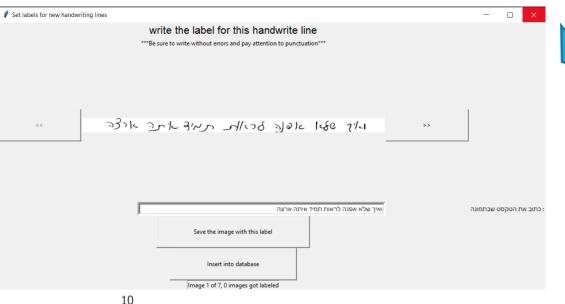
: תהליך

- קבלת קובץ PDG / קובץ תמונה
- במידה ומתקבל PDF מתבצעת המרה לPNG.
 - עיבוד קבצי התמונה ע"י המשתמש – יישור וניקוי רעשים.
 - חילוץ מלבנים מהתמונה •
 - ניקוי "רעשים" בתיחומי המלבנים לפי גודל
- ניקוי של מלבנים מושחרים 🕟
 - חיתוך של התמונה לפיתיחום המלבן
 - תיוג התמונה לפי מילון מוגדר :

(page, line) -> label

תיוג תמונות על ידי המשתמש

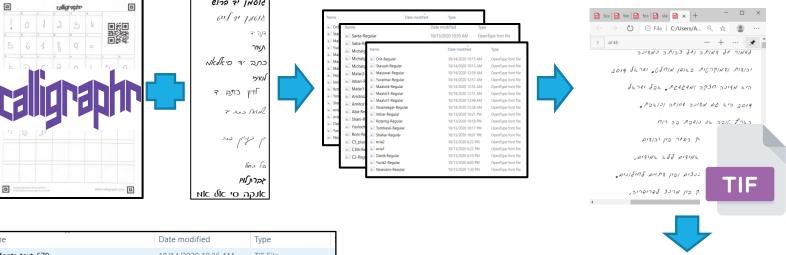


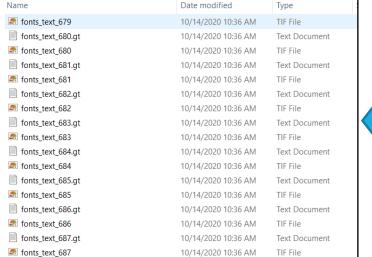


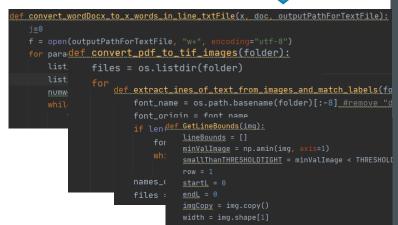
: תהליך

- קבלה כקלט תמונה עם טקסט בכתב יד
- עיבוד קבצי התמונה ע"י המשתמש – יישור וניקוי רעשים.
- הפרדה של הטקסט לשורות על ידי ערך מינימלי של צבע בכל שורה
 - הצגה של השורות בממשק תיוג עבור המשתמש

איסוף דאטה בעזרת פונטים

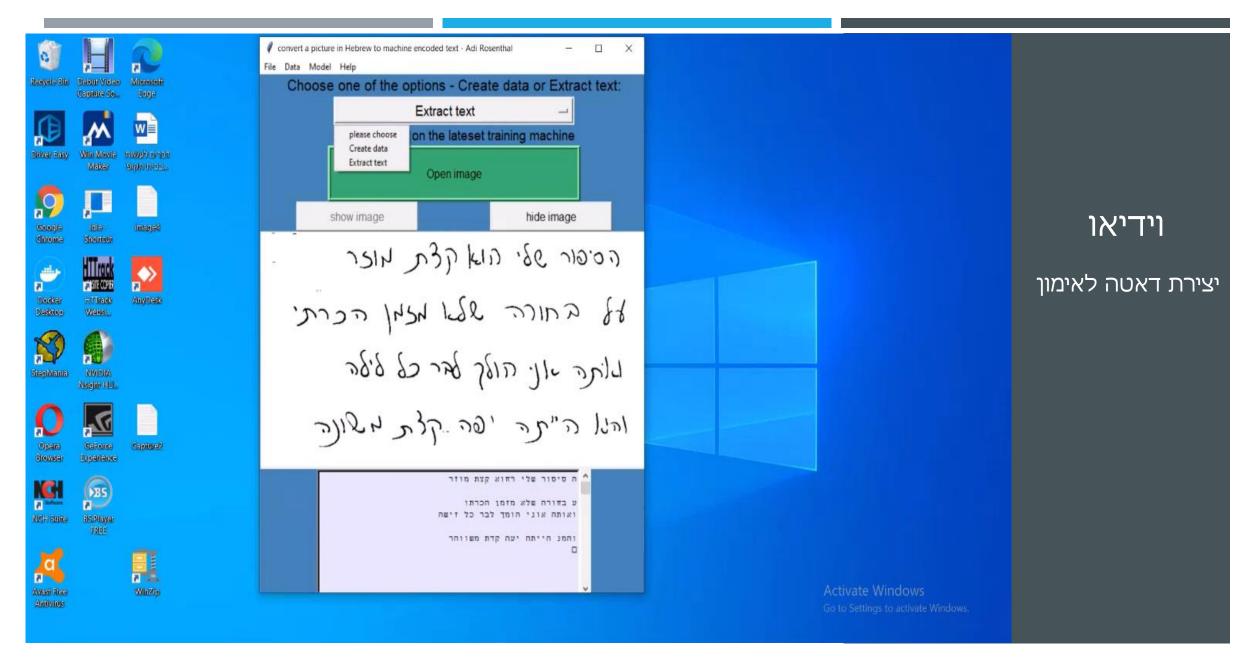






: תהליך

- ריכוז פונטים של כתב יד: •
- חיפוש פונטים שקיימים באינטרנט
- איסוף כ-30 כתבי יד מאנשים שונים ויצירה של פונטים ע"י אתר caligraphr
- ריכוז טקסט ארוך בעברית •
- הפרדה של הטקסט לשורות של כ-6 מילים בשורה
- השמה של הפונטים במסמך וורד, המרה לPDF ובסוף לקובץ TIF + יצירת קובץ טקסט מותאם לתמונות
 - חילוץ שורות טקסט מתוך התמונות
 - תיוג השורות לפי מספר שורה במסמך טקסט



TESSERACT - אימון הרשת

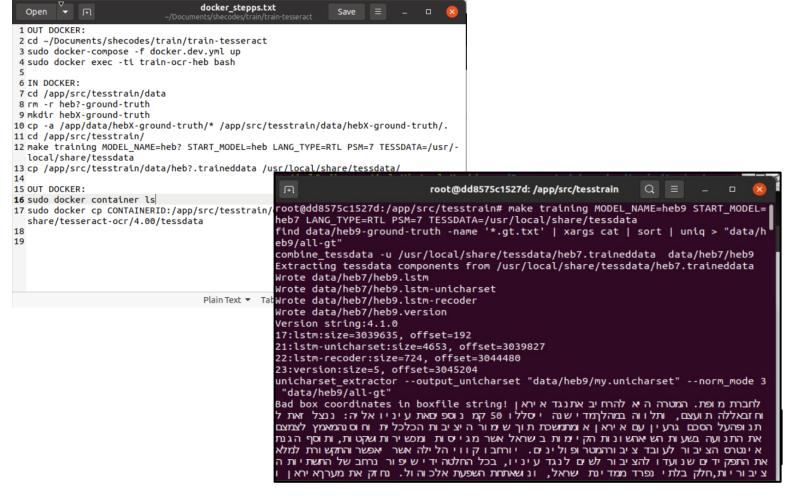
שלושה שלבים עיקריים במהלך הפרוייקט עבור תהליך אימון הרשת:

1. מחקר כלים ואפשרויות לאימון הרשת + יצירת סביבה מתאימה לאימון

.2 תהליך האימון

.3 בדיקה של הרשת המאומנת

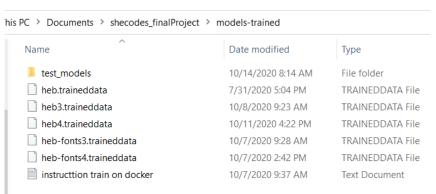
תהליך האימון



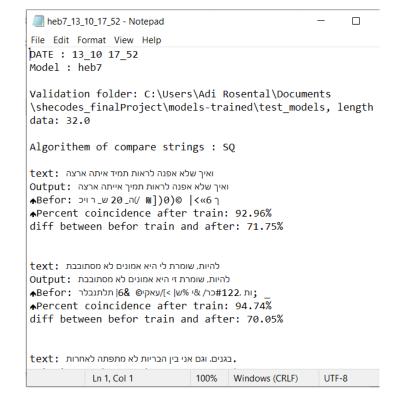
- ארגון של הדאטה והעברתו לסביבה המתאימה
 - **DOCKER**
 - UBUNTO 20
 - Virtual machine
 - שימוש בקוד פתוח דesstrain שעוזר להכין את הדאטה עבור אימון
- .txt + tif -> .box -> .lstm
 - Tesseract training tools
- העברה של הרשת המאומנת (hebX.traineddata) לתקית התקנה של Tesseract

בדיקת הרשת

תוצאות של אימוני הרשת



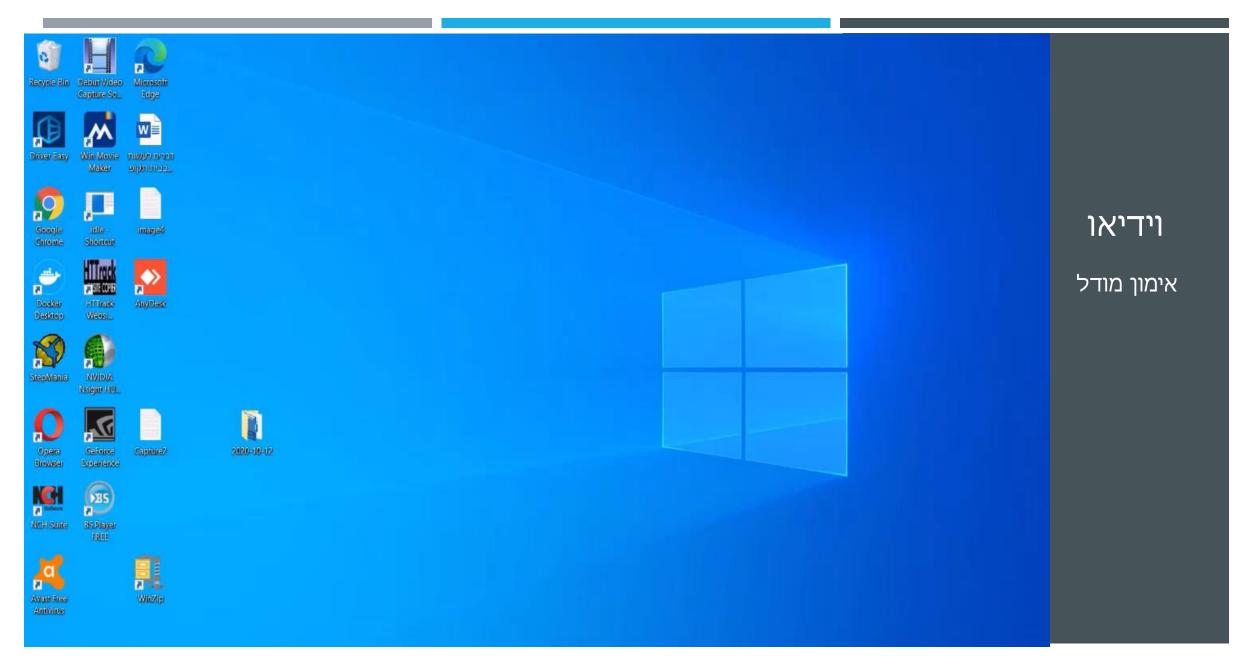
דו"ח בדיקות רשת



- הרצה של הרשת החדשה על הדאטה לבדיקות (בתקיית validation)
- הדאטה מתויג אך הרשתלא אומנה עליו
- השוואה של תוצאת הרשת לתיוג "אמת"
- שימוש בפונקציית
 SequenceMatcher
 על אלגוריתמים להשוואת
 LCS longest מחרוזות
 contiguous matching
- האלגוריתם לא מתייחס למשקלים של אותיות – כמו למשל אותיות דומות וכו'
 - יצירת דוח
- השוואה בין תוצאות הרשת המאומנת על כל שורה בvalidation לתוצאת הרשת על אותם המשפטים לפני האימון והערכה של אחוז שיפור

from difflib import SequenceMatcher as <u>SQ</u>

def Check_model_tesseract(self, folder_validation, folder_output_txtfile, psm=7, compare_methods = "SQ"):

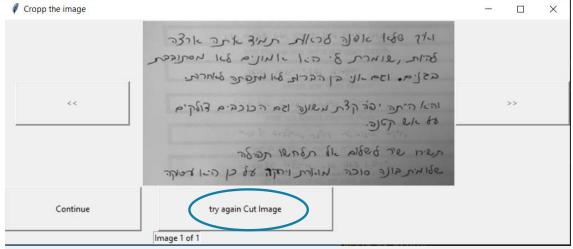


חילוץ טקסט מתוך תמונה

- עיבוד תמונה לפני חילוץ
 - Pytesseract •
- "השוואת תוצאות התוצאה הרשת ל



עיבוד תמונה לפני חילוץ

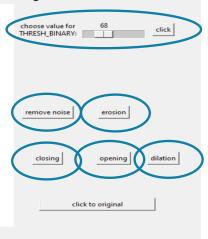


Find the best variation of your image to extracting text

Click here to run program on this image

אבשברונחטירטענסבבקרשת ללן ך מ אבשברונחטינטענסבבקר שתללן ך מ אבשברונחטינטענסבבקרשתלן ך מ

18



כדי שהחילוץ יהיה אופטימלי כדאי לבצע עיבוד תמונה לפני שימוש ברשת נוירונים

בממשק ישנן מספר אפשרויות להטיב עם תוצאת הרשת :

- מרכוז של הטקסט על ידי סימון האזור בו מצוי הטקסט אותו אנחנו רוצים לחלץ.
- ניקוי רעשים בעזרת פונקציות (2) לעיבוד תמונה (OpenCV):

Binary – thresholding

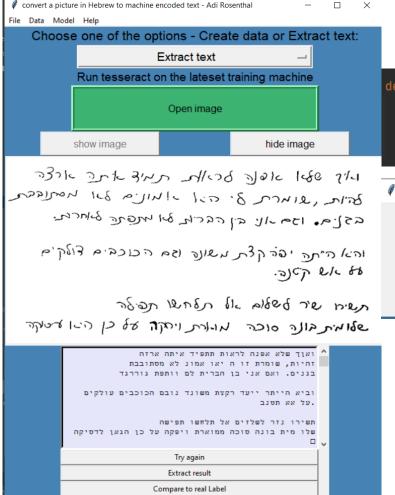
medianBlur

Dilate

Erosion

dilate

חילוץ טקסט

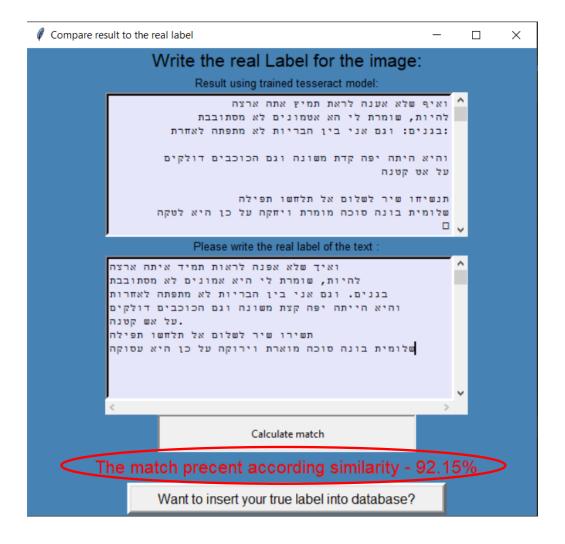


19

תהליך חילוץ הטקסט:

- אופציה לבחירת רשת נוירונים בעזרתה נרצה לחלץ
 - הרשת הנבחרת צריכה להיות שמורה בתקיית tessdata
- באופן דיפולטי הרשת מוגדרת להיות הרשת הכי איכותית
- החילוץ עצמו מתבצע ע"י שימוש בספריה pytesseract שמתממשקת לתוכנה של tesseract – מקבלת תמונה ושם של רשת
 - שמירת תוצאות הרשת כקובץ טקסט (Exstract (result

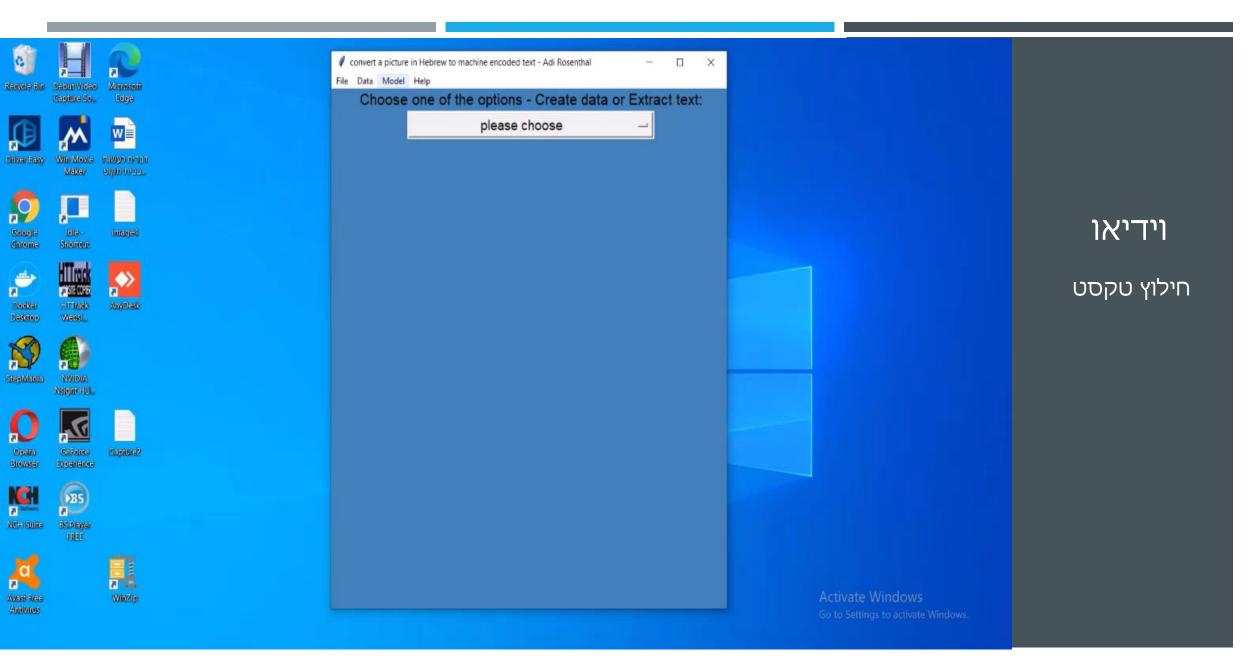
השוואת תוצאות



כדי לבחון את טיב התוצאה שהתקבלה בחילוץ קיימת אופציה עבור המשתמש להכניס את הטקסט הנכון של התמונה – ה"אמת".

> ההשואה בין הטקסטים מתקבלת באחוזים על ידי פונקציית "SequenceMatcher"

(לאחר הכנסת התוצאות אמת ניתן גם להכניס את תמונת הקלט עם הטקסט שהמשתמש הכניס בבדיקה לתוך בסיס הנתונים)



מה למדתי בפרוייקט?

- Creating Virtual machine
- Using Linux
- Work with Dockers
- Git + Github
- Tesseract / OCR
- Design + Architecture
- Using open source

- OpenCV
- Tkinter
- PIL (Image)
- Difflib
- Training tesseract
- Similarity between strings
- Designing and planning a GUI

22

צעדים להמשך...

- שיפור תוצאות הרשת בשיטות נוספות כמו:
 - NLP •
 - שימוש בתיקון על פי מילון -
- סינון של מילים ואותיות- בשילוב עם התניות בשפה העברית
 - https://github.com/NLPH/NLPH Resources
 - בניית אתר אינטרנט / אפליקציה נוחה למשתמש
- ותהליך האימון כדי להשפיע יותר על האימון tesseract בנה מעמיקה יותר של אלגוריתם
 - בדיקה יותר מעמיקה של תוצאות הרשת.

קישורים

- Tesseract-ocr : https://github.com/tesseract-ocr/tesstrain
- https://medium.com/@quiem/how-to-train-tesseract-4-ebe5881ff3b7
- Train-tesseract and docker files for train : https://github.com/guiem/train-tesseract
- Tesstrain : https://github.com/tesseract-ocr/tesstrain
- https://tesseract-ocr.github.io/tessdoc/ImproveQuality
- https://www.makeuseof.com/tag/create-virtual-machine-using-windows-10-hyper-v/
- Learn OpenCV: https://www.youtube.com/watch?v=N81PCpADwKQ&t=6764s
- Learn TKinter: https://www.youtube.com/watch?v=YXPyB4XeYLA
- Learn docker: https://www.youtube.com/watch?v=i7ABlHngi1Q
- SequenceMarcher in python: https://towardsdatascience.com/sequencematcher-in-python-6b1e6f3915fc

אלות?