

Translating Old Persian cuneiform by artificial intelligence (AI)

Author: Shaghayegh Rahmani (Melanee)

Email: melaneepython@gmail.com

GitHub: <https://github.com/Electronic-Old-Persian-Library>

Iran, Tehran

24 July 2024

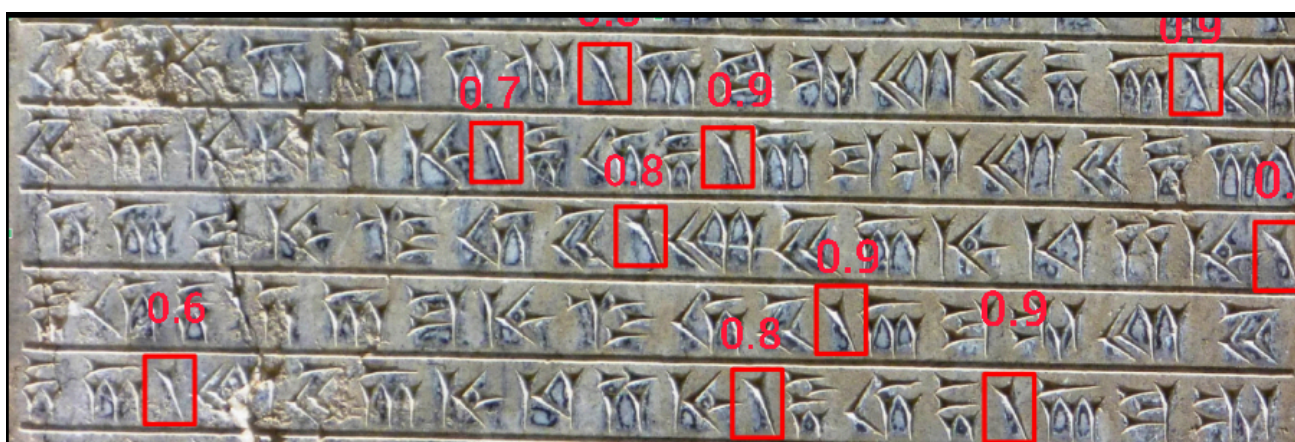


Image by author

Keywords

Old Persian, translating, cuneiform, artificial intelligence, OCR, NLP

Abstract

In this project, I developed an optical character recognition (OCR) model to decipher Old Persian language to modern languages. The process of converting an image text to machine-readable text is called OCR. To the best of my knowledge, I am the first one who is developing this language in this scale in the world. The primary raw data is collected from all over the world. For instance, the [British museum collection](#), cuneiform digital library Initiative ([CDLI](#)), [Livius](#) and my personal photography from the national museum of Iran and Takht-e-Jamshid (Persepolis). I have trained my data with some OCR models, the best one is based on the EasyOCR model.

Introduction

Saving the history of my country, Iran, is a very important and valuable issue that we have to care about. The history that we have nowadays is our ancestor's heritage that we must transfer to the next generations. In this regard, deciphering ancient cuneiforms has been one of the difficult challenges during the span of time to figure out how our ancestors could live in this world. One day, I was thinking, if there is no one in this world who can translate the Old Persian language, what will happen to this language?! How will the next generation decipher new found inscriptions or tablets?! So, this new brainstorm crosses my mind to develop a new AI model to achieve this goal and keep this ancient language alive forever.

Previous research

In June 2024, I joined [Professor Enrique Jiménez](#) team. He and his team at Ludwig Maximilian University (LMU) was utilising advanced artificial intelligence (AI) and computer vision techniques to analyse and decipher Babylonian and Akkadian inscriptions, including the Epic of Gilgamesh. They have developed several tools and algorithms to facilitate this process, which are available in their [GitHub](#) repositories and [electronic babylonian library](#) (eBL). My project is inspired by the eBL project. Professor Enrique Jiménez has developed models for Babylonian cuneiform but I am going to develop my models for Old Persian cuneiform.

Developing models

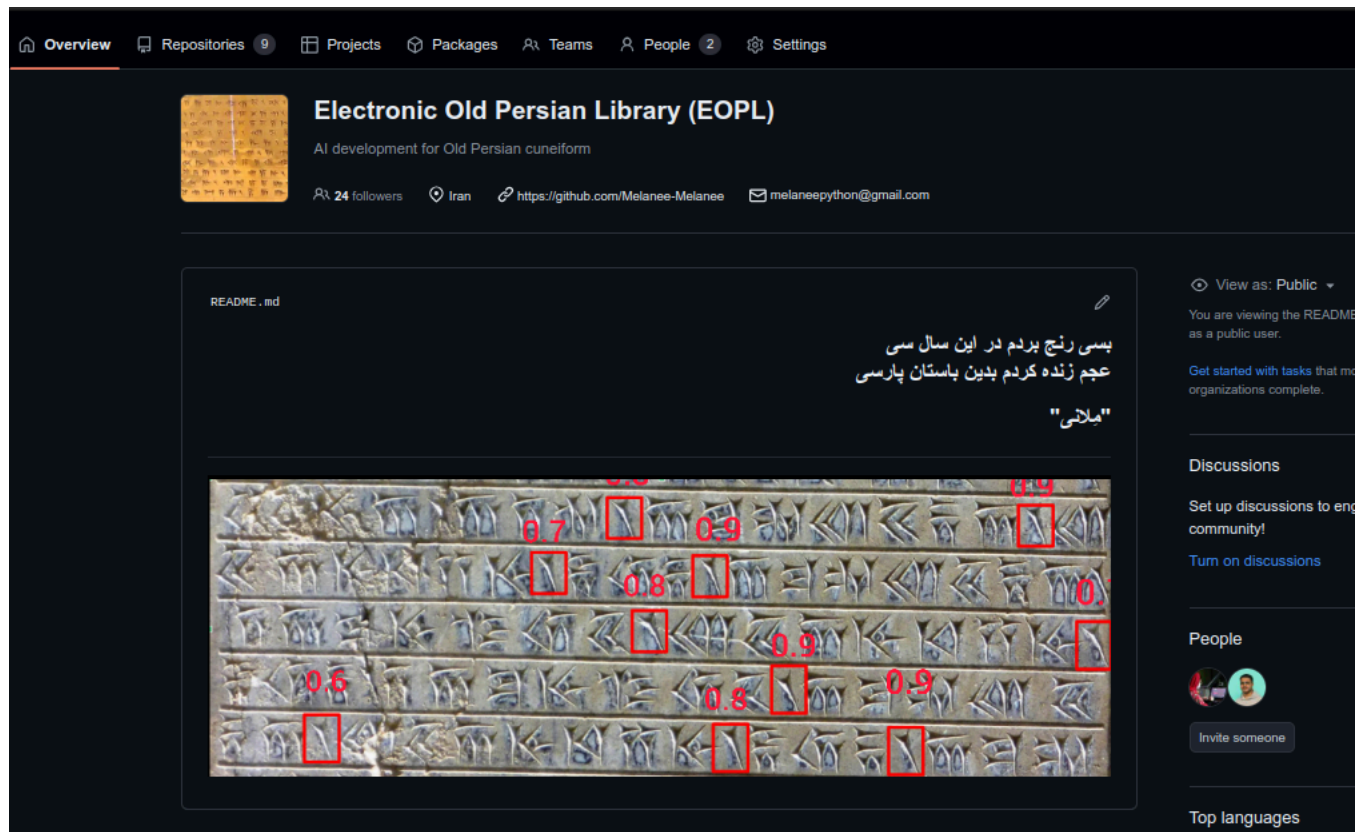
I have developed three OCR models in this project:

- yolo_cnn_old_persian
- tesseract_old_persian
- easyocr_old_persian

I have developed an “easyocr_old_persian” model from scratch and trained it with 42 images of the last 12 lines of the great Darius's inscription in Persepolis, [DPd inscription](#) as primary data. This model deciphers Old Persian cuneiform to English transcription and is published on the [Old-Persian-Cuneiform-OCR](#) repository on Github. This repository is a part of [Electronic Old Persian Library](#) (EOPL) organisation. Since I was looking for an OCR model for the Old Persian language, I have not implemented image pre-processing for my models yet and they work on just black and white images. Additionally, my models are still under developing. In the future, I will train the “easyocr_old_persian” model with more real data to get better results.

The “tessearct_old_persian” model is a pre-trained model developed by [S. Muhammad Hossein Mousavi](#). I just developed a code for evaluating this pre-trained model in Python programming language.

The “yolo_cnn_old_persian” model is not completed yet.



Electronic Old Persian Library (EOPL) organisation on GitHub

Results and discussion

Evaluating “tessearct_old_persian” model code:

```
!pip install pytesseract
!apt-get install tesseract-ocr
!pip install pillow
%cd /usr/share/tesseract-ocr/4.00/tessdata/
!mv/content/drive/MyDrive/Persiancuneiform1/tesseract/myLang.tr
aineddata /usr/share/tesseract-ocr/4.00/tessdata

from PIL import Image
import pytesseract

pytesseract.pytesseract.tesseract_cmd = r'/usr/bin/tesseract'

img_path = '/content/drive/MyDrive/Persiancuneiform1/DPd.png'
img = Image.open(img_path)
text = pytesseract.image_to_string(img, lang='myLang')
print(text)
```

Input:



The image displays a sample of Old Persian cuneiform script. The text is arranged in approximately 10 horizontal lines. The characters are stylized, wedge-shaped symbols. The script is a form of the Old Persian language, which was used by the Achaemenid Empire. The text is likely a fragment from a larger inscription, such as the Behistun inscription, which is a significant historical document in the study of ancient Iran and the Achaemenid Empire.

Output:

Zittiy ; iaryvuS ; xrSayZiy;

mnc;aurmzia;upstam; rlauv;

hia ; ViZiriS ; rgiriS ; uta;

im am ; i h yaum ; au lm z i a ;

pitTucs;hca;hinaya; hca;

QuSiyala ; hca;iruga;ariy;

imam ;ihyaum;ma; ajMiya; ait;

aim ;yanm;jDiyaMiy;

aitmiy ; iiaTuv

In the next step. I translate that Old Persian transcription to modern languages by [Chat-GPT](#) (GPT-4 version):

Translate to Modern Persian:

"این منم داریوش شاهنشاه؛ به لطف اهورامزدا، من این را بنا کردم؛ من این امپراتوری را بنیان نهادم و آن را نیرومند ساختم. باشد که اهورامزدا من و پادشاهی مرا محافظت کند؛ باشد که برای همیشه پایدار بماند؛ و باشد که از دروغ در امان باشد؛ این است آنچه من انجام دادم؛ این است آنچه من می‌گویم."

Translate to Modern English:

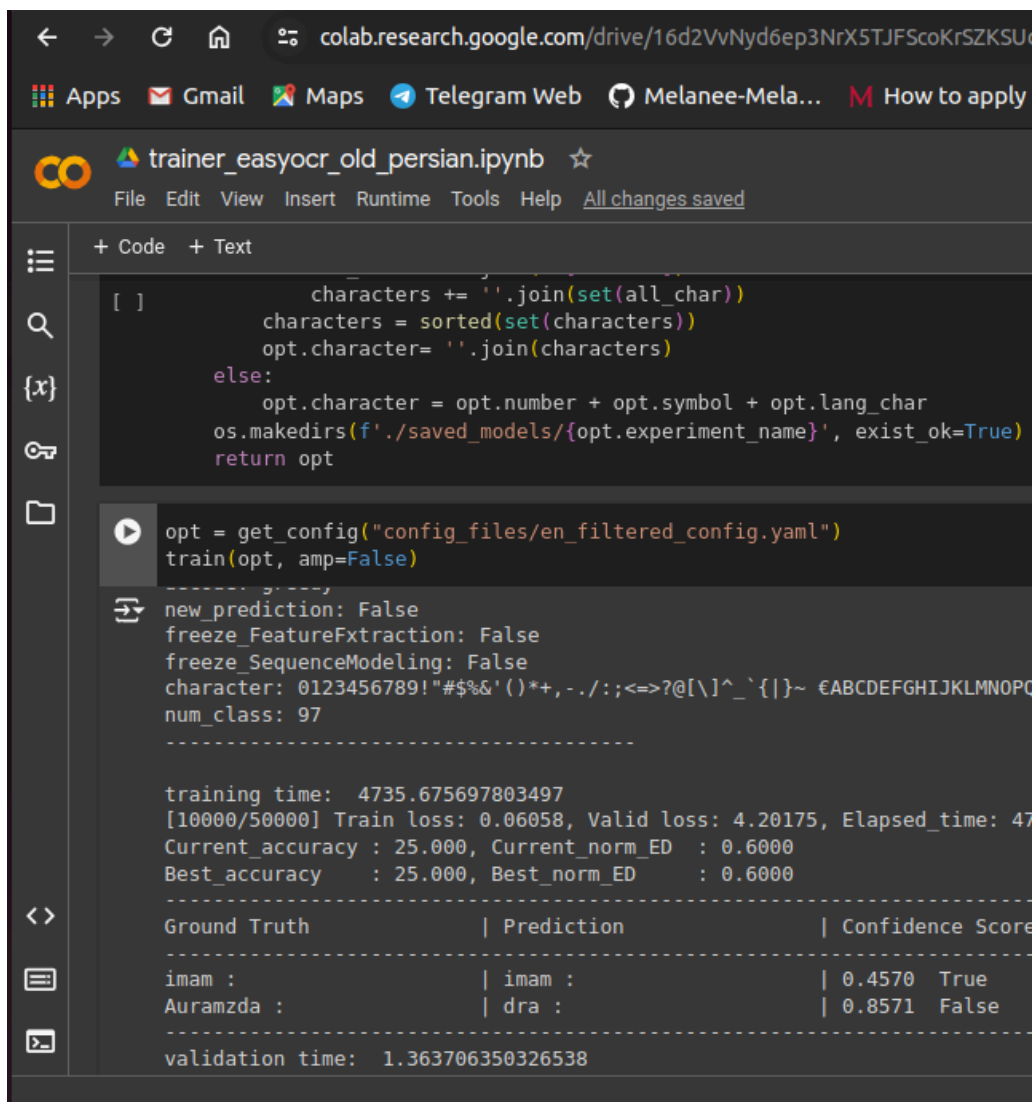
"This is me, Dariush king; By the grace of Ahura Mazda, I have built this; I founded this empire and made it strong. May Ahuramazda protect me and my kingdom; may it last forever; and it would be safe from lies; that is what I did; That is what I am saying."

But in the future, I will develop my own natural language processing (NLP) model to translate the output instead of using ChatGPT.

Results for “easyocr_old_persian” model:

Step	num_iter	valInterval	Confidence score	Execution time(min)
1	300	200	0.0032	5
2	3000	2000	0.5	15
3	10000	50000	0.45	102

Table 1. Training results for “easyocr_old_persian” model



```

[ ]      characters += ''.join(set(all_char))
          characters = sorted(set(characters))
          opt.character= ''.join(characters)
        else:
          opt.character = opt.number + opt.symbol + opt.lang_char
          os.makedirs(f'./saved_models/{opt.experiment_name}', exist_ok=True)
          return opt

opt = get_config("config_files/en_filtered_config.yaml")
train(opt, amp=False)
-----
new_prediction: False
freeze_FeatureExtraction: False
freeze_SequenceModeling: False
character: 0123456789!"#$%&'()*+,-./:;<=>?@[\\]^_`{|}~ €ABCFGHIJKLMNOPQ
num_class: 97
-----

training time: 4735.675697803497
[10000/50000] Train loss: 0.06058, Valid loss: 4.20175, Elapsed_time: 47
Current_accuracy : 25.000, Current_norm_ED : 0.6000
Best_accuracy   : 25.000, Best_norm_ED   : 0.6000
-----
Ground Truth      | Prediction      | Confidence Score
-----
imam :            | imam :         | 0.4570 True
Auramazda :      | dra :          | 0.8571 False
-----
validation time: 1.363706350326538

```

Training results for “easyocr_old_persian” model in my notebook

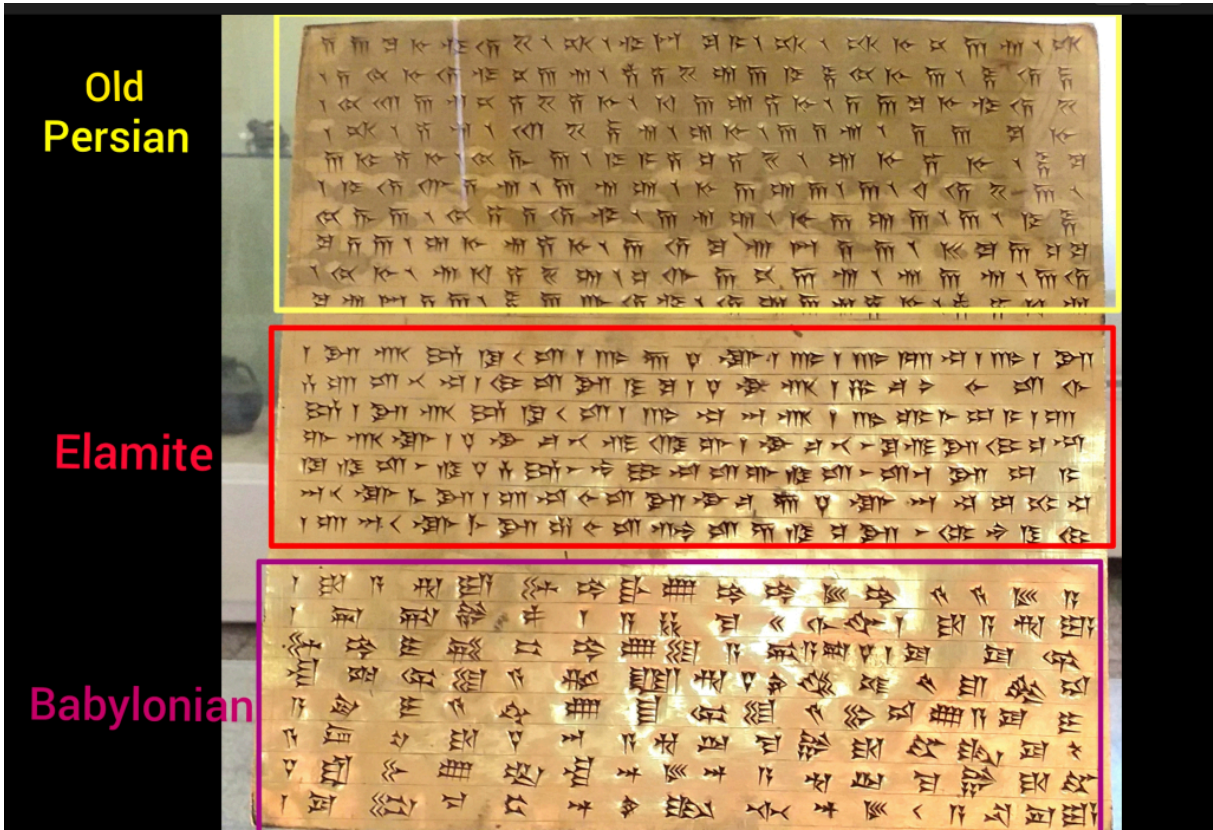
All my models have been developed on Google Colab notebook infrastructure:

Operating system	version	GPU	CPU	RAM
Linux, Ubuntu	22.04 LTS	T4GPU	-	16G

Table 2. Infrastructure of my Google Colab notebook

Future Research

- This “easyocr_old_persian” model can be implemented for other languages, for example, Akkadian, Pahlavi, Avestan, Sogdian and so on.
- By deciphering Old Persian cuneiform in trilingual inscriptions or tablets (e.g, Old Persian, Babylonian and Elamite), we can decipher Babylonian and Elamite cuneiform because they have same meaning.



Trilingual gold plate of king Darius, photo is taken from national museum of Iran

- Matching broken tablets or inscriptions, we can use OCR to convert the cuneiform signs from images into machine-readable text. Then we will use Prof Enrique Jiménez's NLP project to apply [algorithms](#) to detect and match segments of different tablets or inscriptions, aiding in the reconstruction of fragmented texts.

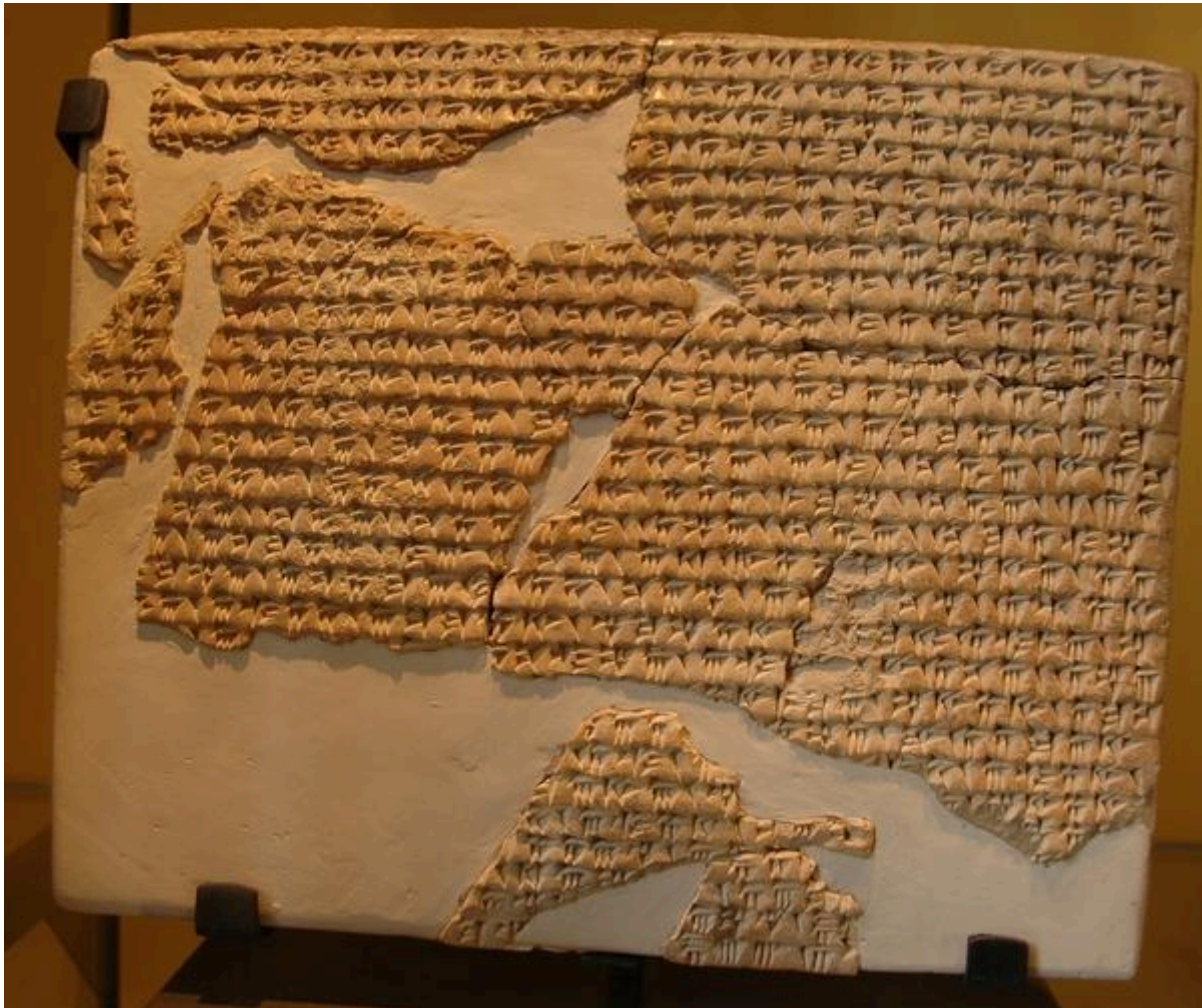


Photo is from Apadana Castle Shush, [DSf inscription](#)

Future work

- Developing pre-processing for the raw image dataset
- Training “easyocr_old_persian” model with huge images
- Developing NLP models to translate English transcriptions to modern Persian

Conclusion

Acquired results of the evaluation indicate that my models will be able to properly translate Old Persian cuneiform. The acquired results are promising that they are able to make and improve NLP in this area.

References

https://www.researchgate.net/publication/370482228_Translating_Akkadian_to_English_with_neural_machine_translation

https://www.researchgate.net/publication/324725453_Extracting_old_persian_cuneiform_font_out_of_noisy_images_handwritten_or_inscription

https://www.researchgate.net/publication/267867570_Intelligent_Recognition_of_Ancient_Persian_Cuneiform_Characters

The pre-print of this paper is archived on “arxiv.org” with submission number 5778688 on 07 Aug 2024 to avoid any plagiarism.

All source code of this paper is under CC-BY-NC license and any commercial use is prohibited.