VALIDAÇÃO DE DOCUMENTOS RG, CPF E DIPLOMAS, COM AUXÍLIO DE OCR-TESSERACT

Gabriel Mancini

Pedro Marcelo Prado



O QUE É OCR?

- OCR OPTICAL CHARACTER RECOGNITION (Reconhecimento Óptico de Caracteres)
- É uma tecnologia que permite a conversão de texto apresentado em imagens ou documentos digitalizados em texto editável e pesquisável por máquinas.



COMO ELE FUNCIONA?

O OCR utiliza algoritmos de processamento de imagem e inteligência artificial para identificar e interpretar caracteres (letras, números, símbolos) contidos em uma imagem.

· Pré-processamento da imagem :

Ajuste de brilho e contraste.

Remoção de ruídos.

Transformação para escala de cinza ou binarização.

· Segmentação:

Identificação de linhas, palavras e caracteres individuais.

· Reconhecimento de padrões :

Comparação de cada característica com padrões predefinidos ou treinamento em modelos baseados em aprendizado de máquina.

· Pós-processamento:

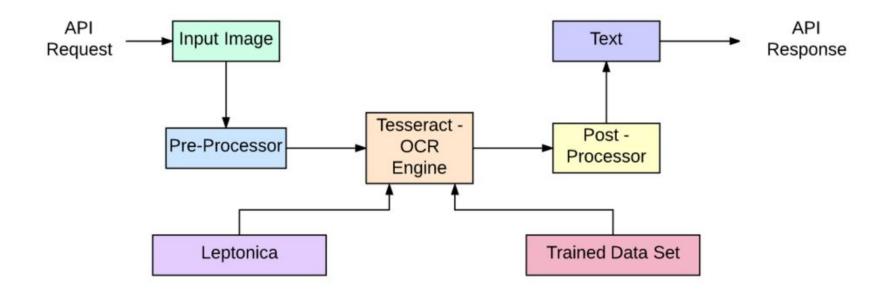
Correção de erros e ajustes no texto reconhecido, utilizando dicionários ou regras gramaticais.







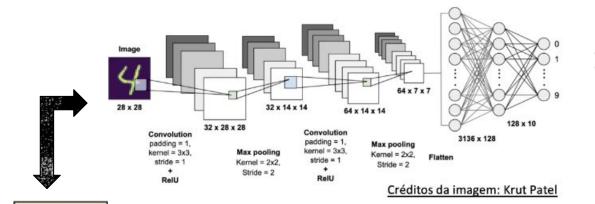
ESTRUTURA DO TESSERACT



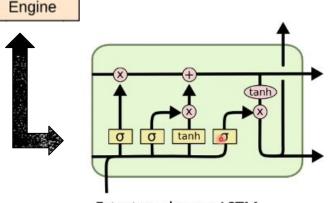
Créditos da imagem: Balaaji Parthasarathy



APRENDIZADO PROFUNDO - TESSERACT

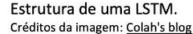


Para reconhecer uma imagem contendo um único caractere, normalmente usamos uma Rede Neural Convolucional (CNN)



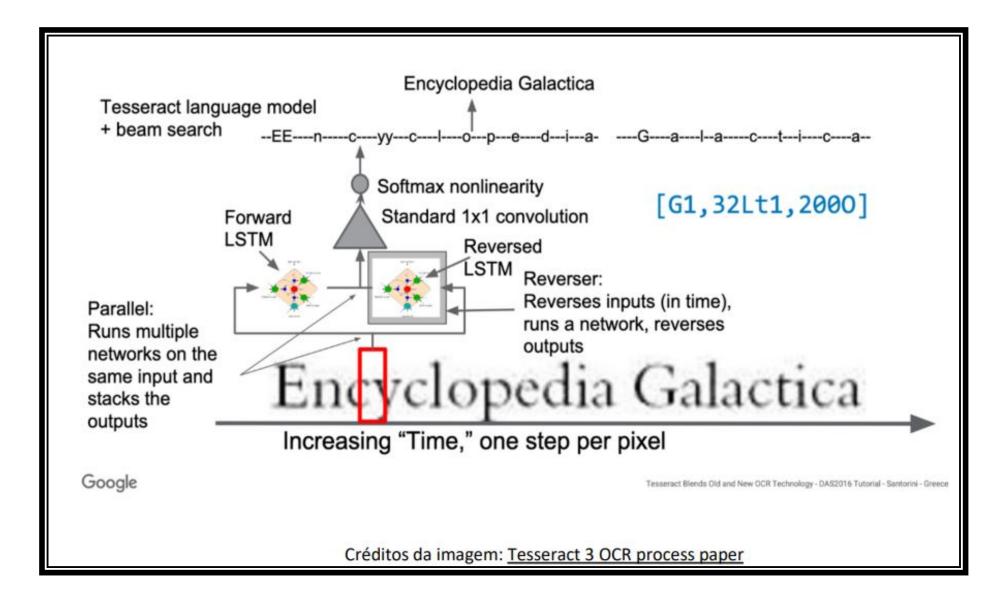
Tesseract -OCR

Um texto de comprimento arbitrário é uma sequência de caracteres e é mais interessante utilizar RNNs (Redes Neurais Recorrentes). LSTM (Long short-term memory) é uma forma popular de RNN.





FUNCIONAMENTO



FUNCIONAMENTO DO CHATBOT RG E CPF

Importação de Bibliotecas Montagem do Google Drive e Cópia de Arquivos

Função para Exibição de Imagens

Pré-processa mento da Imagem Detecção e Ordenação de Contornos

Correção de Perspectiva Reconhecimento Óptico de Caracteres (OCR)

Extração de Dados Interface com
ChatBot



FUNCIONAMENTO DO CHATBOT PARA DIPLOMA

Importação de Bibliotecas

Montagem do Google Drive e Cópia de Arquivos Função para Exibição de Imagens

Reconhecimento Óptico de Caracteres (OCR) Extração de Dados Interface com ChatBot

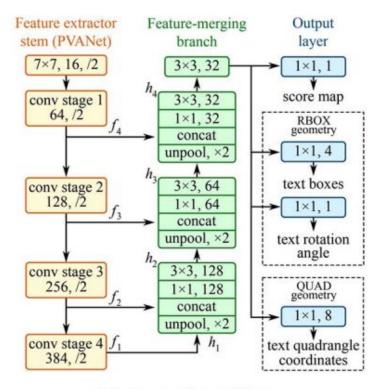


1° TENTATIVA EAST COM TESSERACT

- O EAST (Efficient Accurate Scene Text detector) é um modelo de aprendizagem profunda (deep learning), publicado oficialmente em 2017 por Zhou et al.
 - Usa camadas convolucionais para extrair características das imagens e dessa forma detectar a presença de textos

Estrutura do EAST

(Fully-Convolutional Network)



Créditos: Zhou et al.



RESULTADO DO EAST

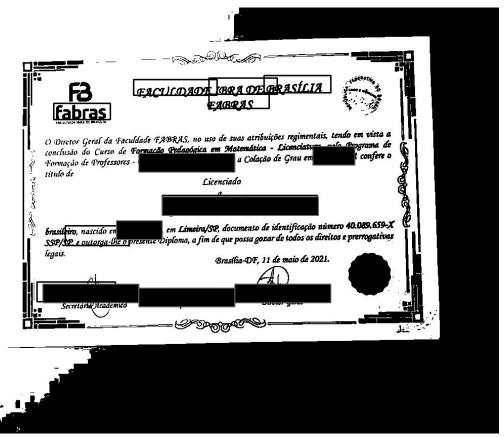






RESULTADO DO EAST



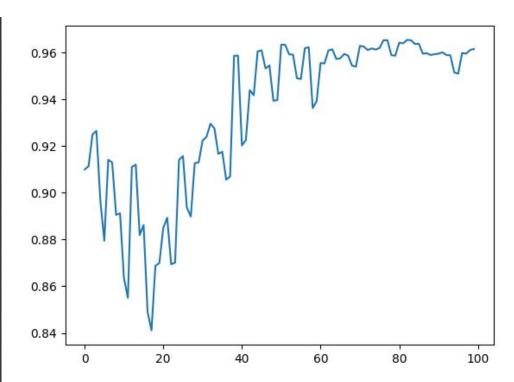




2° TENTATIVA OCR PERSONALIZADO

Layer (type)	Output Shape	Param #
conv2d (Conv2D)	(None, 26, 26, 32)	320
max_pooling2d (MaxPooling2D)	(None, 13, 13, 32)	Ø
conv2d_1 (Conv2D)	(None, 13, 13, 64)	18,496
max_pooling2d_1 (MaxPooling2D)	(None, 6, 6, 64)	0
conv2d_2 (Conv2D)	(None, 4, 4, 128)	73,856
max_pooling2d_2 (MaxPooling2D)	(None, 2, 2, 128)	Ø
flatten (F <mark>latten</mark>)	(None, 512)	Ø
dense (Dense)	(None, 64)	32,832
dense_1 (Dense)	(None, 128)	8,320
dense_2 (Dense)	(None, 36)	4,644

Total params: 138,470 (540.90 KB)
Trainable params: 138,468 (540.89 KB)
Non-trainable params: 0 (0.00 B)
Optimizer params: 2 (12.00 B)



Banco de dados:

MNIST 0-9 Kaggle A-Z



2° TENTATIVA OCR PERSONALIZADO

T -> 99.89%

E -> 100.00%

S -> 98.41%

T -> 99.34%

A -> 100.00%

N -> 99.99%

D -> 100.00%

0 -> 54.03%

D -> 93.65%

M -> 99.94%

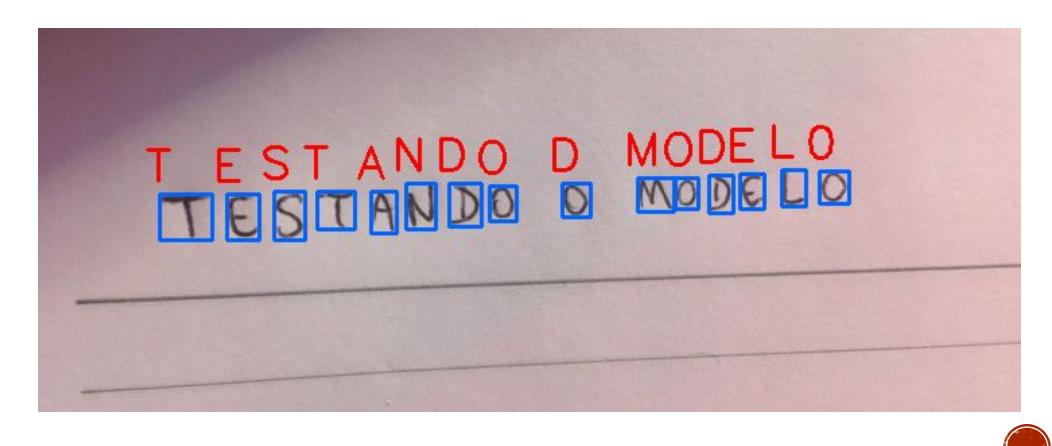
0 -> 44.77%

D -> 99.97%

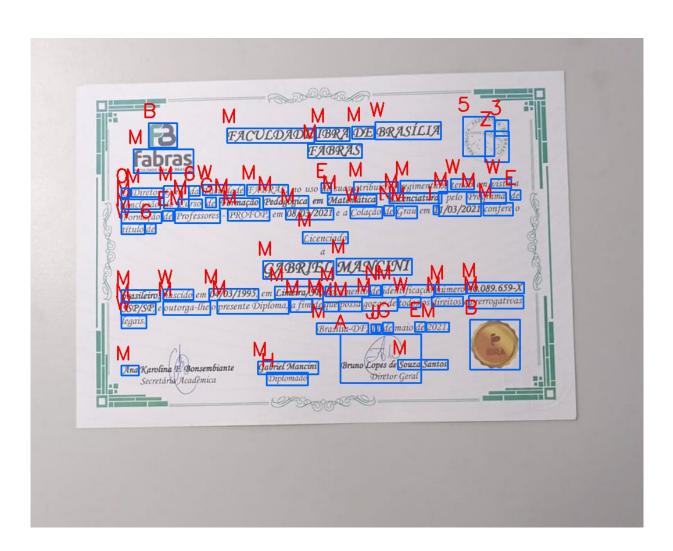
E -> 99.98%

L -> 100.00%

0 -> 77 17%



2° TENTATIVA OCR PERSONALIZADO

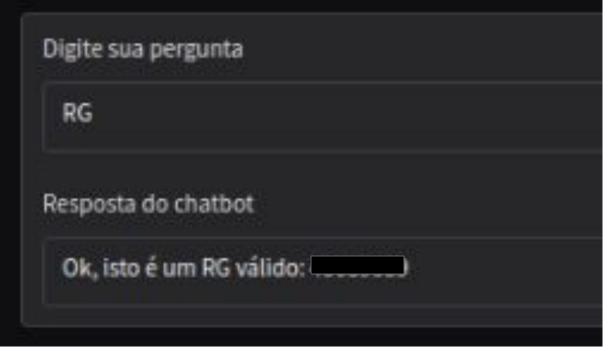


Texto reconhecido:



O QUE DEU CERTO







O QUE DEU CERTO





Implementação

Diploma

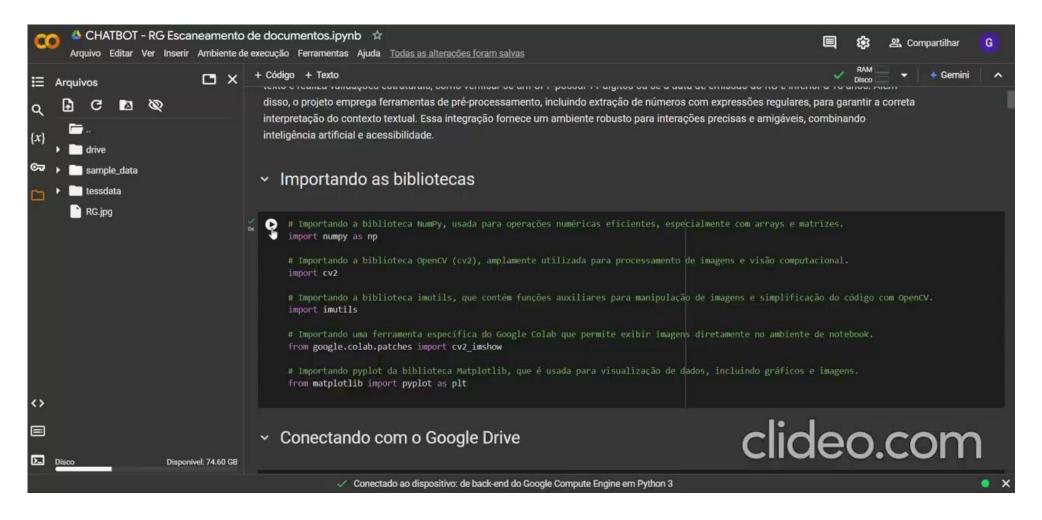
O código usa OCR e ferramentas de processamento de imagens para detectar e destacar termos em documentos, validando e extraindo informações automaticamente.

RG

O projeto usa NLP e **DistilBERT** para criar um chatbot que valida informações de documentos, como CPF e RG, com uma interface gráfica feita com **Gradio**.

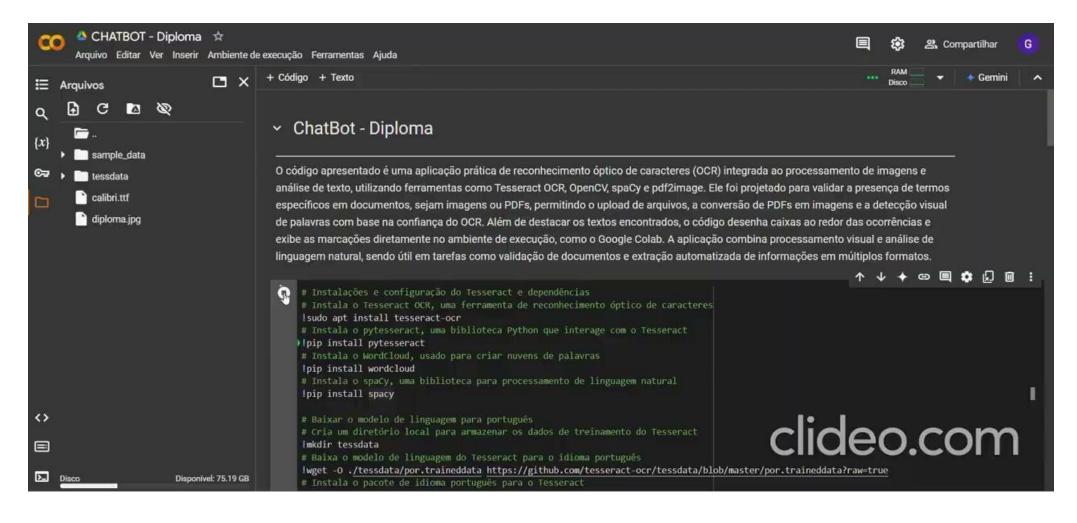


RG e CPF





Diploma





Pontos positivos

A implementação do sistema na UNESP proporcionaria:

- Maior eficiência em processos administrativos.
- Redução de erros na validação de documentos.
- Atendimento ágil e acessível para a comunidade acadêmica.
- Redução da carga de trabalho dos funcionários, permitindo maior foco em atividades estratégicas.



PERSPECTIVAS FUTURAS

- Implantar o CHATBOT de RG e CPF para leitura de pdf.
- De posse de um banco de dados estruturado seria possível realizar o treinamento de uma rede yolo para detectar se a imagem inserida se trata de um RG ou não, e em seguida extrair as informações de acordo com o CHATBOT implementado.
- Com o banco de dados seria possível realizar o treinamento personalizado de um extrator de texto, específico para extrair informações de grande relevância nestes documentos.
- Melhor a detecção de contornos para que seja possível usar o CHATBOT de RG e CPF para o DIPLOMA.



Referências

- GOODFELLOW, Ian; BENGIO, Yoshua; COURVILLE, Aaron. Deep Learning. MIT Press, 2016.
- SMITH, Ray. An Overview of the Tesseract OCR Engine. *Proceedings of the Ninth International Conference on Document Analysis and Recognition*, IEEE, 2007.
- BROWN, Tom et al. Language models are few-shot learners. Advances in Neural Information Processing Systems, 2020.



VALIDAÇÃO DE DOCUMENTOS RG, CPF E DIPLOMAS, COM AUXÍLIO DE OCR-TESSERACT

Gabriel Mancini

Pedro Marcelo Prado