

NOVA  
UNIVERSIDADE NOVA  
DE LISBOA

Palacký University  
Olomouc

PÉCS TUDOMÁNTUDYTEM  
UNIVERSITY OF PÉCS

Humanities Going Digital (HUGOD)  
[2020-1-CZ01-KA226-HE-094363]

Compilação de Corpora -  
OCR

Sessão 04

12/10/2022

Chiara Barbero  
Sílvia Barbosa

HUGOD

Erasmus+

Na sequência das atividades propostas em sala de aula e do exercício individual pedido, disponibilizamos este breve guião como material de formação contínua ao longo do curso e para a vossa autocorreção em vista da avaliação final.

O Reconhecimento Ótico de Caracteres (mais comumente conhecido como OCR, usando o acrônimo inglês) é um procedimento, que faz uso de ferramentas digitais, para reconhecer caracteres a partir de um arquivo de imagens não editáveis. Resumindo, os programas de OCR desbloqueiam o acesso ao texto dos ficheiros que precisamos usar, de forma a torná-los manipuláveis, arquiváveis e compatíveis para outros programas.

🔍 Como escolher a ferramenta mais adequada para transformar um ficheiro não editável (JPEG, PNG ou PDF) para um ficheiro editável (.doc, .txt..)?

A escolha da ferramenta mais adequada depende de uma série de fatores:

- do tipo de ficheiros de partida que temos (manuscrito ou impresso, nem todas as ferramentas aceitam qualquer tipo de ficheiro de *input*);
- do tipo de ficheiro final como *output* que precisamos;
- da acessibilidade das ferramentas (licenças pagas), possibilidade de usar versões demo (ex. Abby Fine Reader, Adobe, etc.);

- da quantidade de ficheiros;
- da língua dos ficheiros de partida.

Devemos testar diversas ferramentas e avaliar os resultados obtidos para podermos escolher o melhor *output* que facilite a revisão manual em curto espaço de tempo.

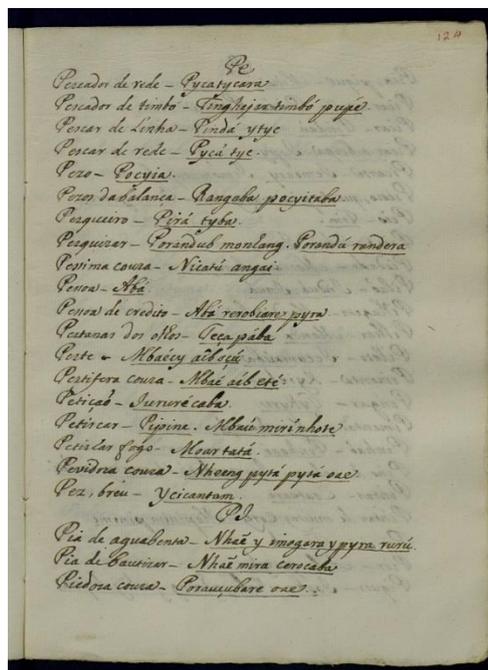
Apesar de cada ferramenta ter uma interface própria, de um modo geral todas apresentam:

- um espaço para fazer o *upload* do ficheiro que queremos transformar: é prática corrente existir um espaço para fazer o *upload* do ficheiro, através da função *arrastar e soltar* (drag and drop) numa caixa da página do software ou, em alternativa, de um botão para seleccionar a pasta do seu computador e carregar o ficheiro em questão.
- uma série de configurações internas que podem (e devem) ser utilizadas para adaptar e calibrar o processo de OCR ao ficheiro de partida, por exemplo: a língua a ser usada (é importante seleccionar a língua pois melhora substancialmente o desempenho do *software*); a possibilidade de escolher o formato de ficheiro final de *output*;
- um comando para iniciar a conversão do ficheiro;
- um comando para o *download* do ficheiro final como *output* (quando não é realizado em automático).

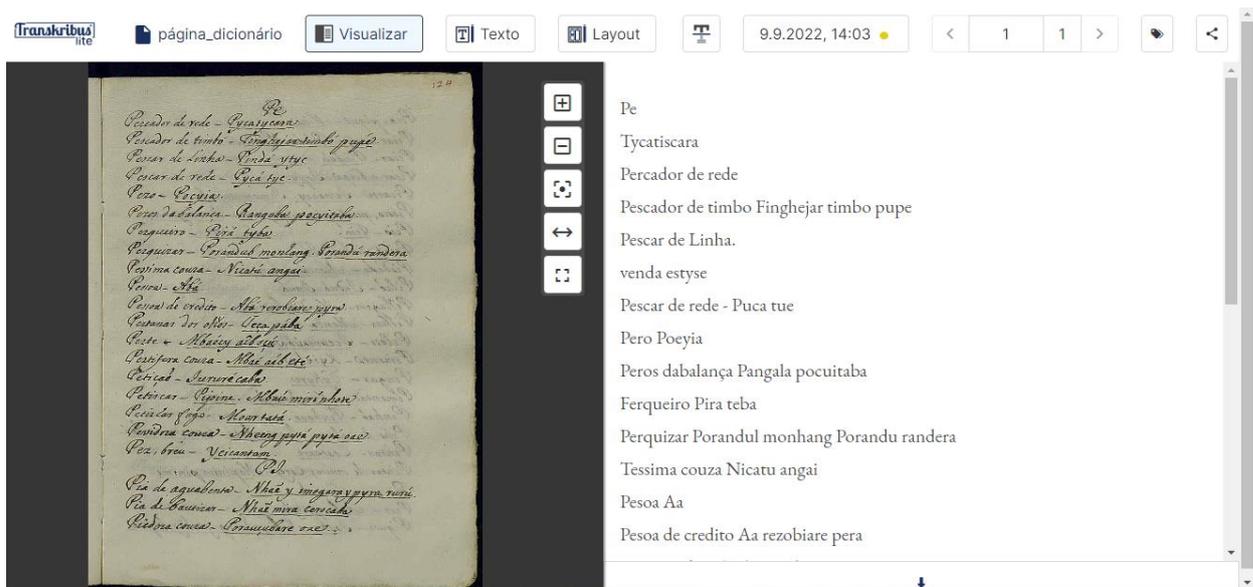
De seguida, apresentamos alguns exemplos de ficheiros de imagens com a respetiva metodologia a aplicar consoante a ferramenta seleccionada para o efeito, que pode ser replicada para outros ficheiros.

### Exemplo 1 (de JPEG para DOCX)

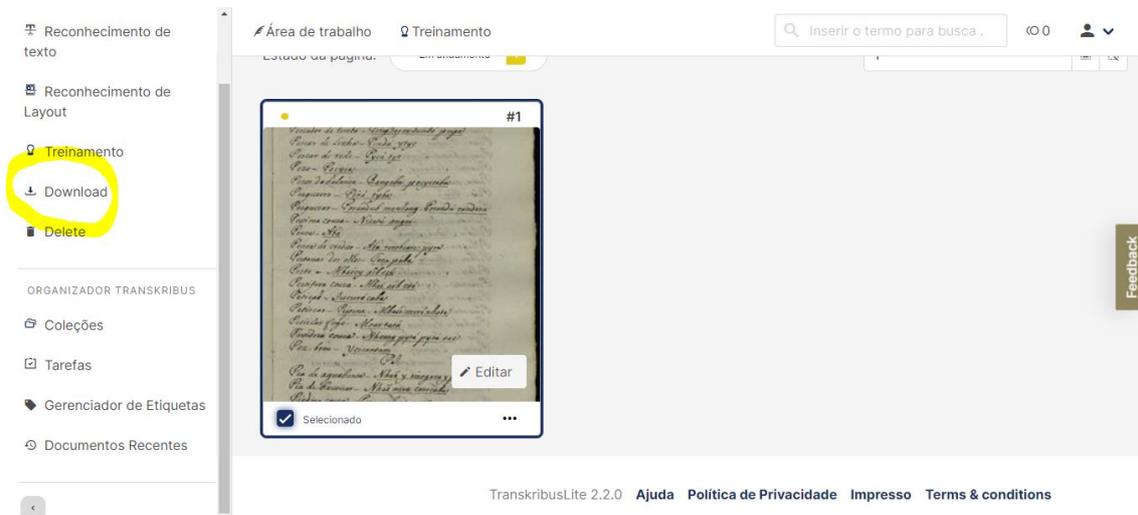
- Descrição: Ficheiro de uma imagem de uma página com texto manuscrito retirada de um dicionário “Diccionario da lingua geral do Brasil”.
- Esta imagem encontra-se disponível em [https://digitalis-dsp.uc.pt/bg3/UCBG-Ms-81/UCBG-Ms-81\\_item1/P275.html](https://digitalis-dsp.uc.pt/bg3/UCBG-Ms-81/UCBG-Ms-81_item1/P275.html) e está no formato JPEG.



- Ferramenta utilizada: Transkribus (disponível em <https://transkribus.eu/lite>). Apesar de ser uma ferramenta gratuita é necessário registar-se e fazer *login*.
- Esta ferramenta aceita como ficheiros de partida: JPEG, PNG ou PDF.
- Esta ferramenta disponibiliza como ficheiros finais: XML, Alto, PDF, TEI, DOCX.



- Esta ferramenta possibilita a extração do texto através de (i) a função de copiar e colar manualmente, ou (ii) de fazer o *download*.



- Relativamente à imagem submetida obtivemos os seguintes resultados:
  - Ao nível da disposição do texto na página, na interface *online* a visualização é alinhada entre texto original e transcrição, no *download* temos desformatação quase total do ficheiro de partida;
  - Ao nível do reconhecimento de caracteres obtivemos alguns erros (i.e. transcrição de letras erradas)

- De modo geral as palavras em português foram bem reconhecidas (ferramenta usada e calibrada para o português), os equivalentes em língua indígena (não é especificada qual) não o foram.
- Para saber mais como utilizar esta ferramenta consulte o guia de utilização (<https://readcoop.eu/transkribus/howto/getting-started-with-transkribus-lite/>).

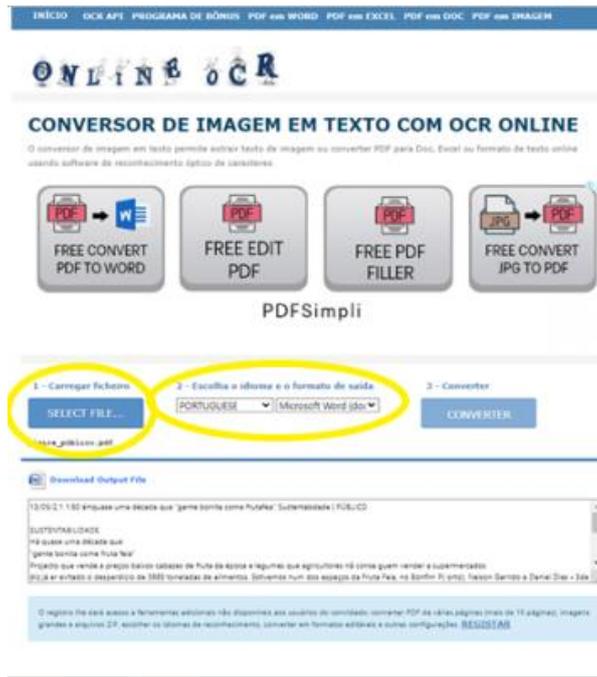
## Exemplo 2 (de PDF para DOCX)

- Descrição: Ficheiro de uma página do jornal *Público*, em formato digital.
- Esta imagem encontra-se disponível em

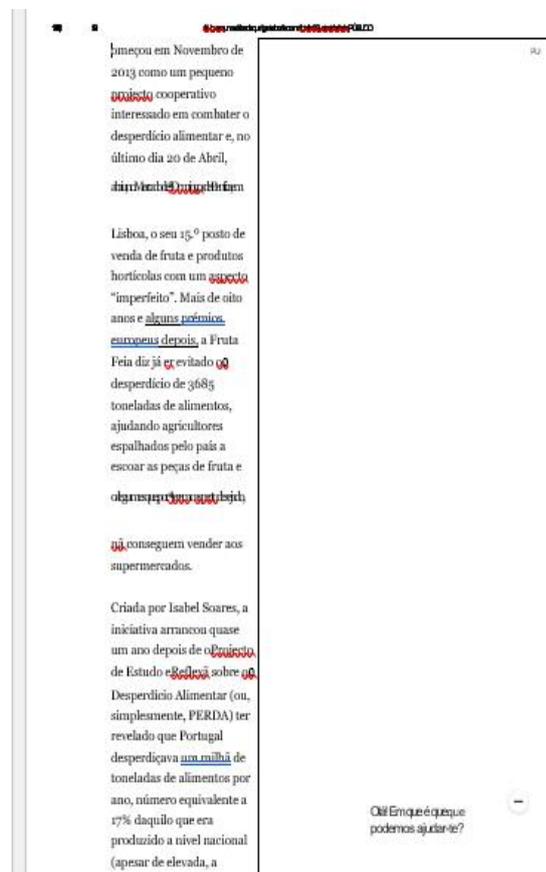


<https://www.publico.pt/2022/05/03/azul/fotogaleria/fruta-feia-sao-ja-15-postos-venda-cooperativa-combate-desperdicio-alimentar-407910> e está no formato PDF.

- Ferramenta utilizada: ONLINE OCT (<https://www.onlineocr.net/pt/>). Apesar de ser uma ferramenta gratuita há restrições quanto ao tamanho do ficheiro de partida.



- Esta ferramenta aceita como ficheiros de partida: JPEG, PNG ou PDF.
- Esta ferramenta disponibiliza como ficheiros finais: WORD, EXCEL, TXT.



- Relativamente à imagem submetida obtivemos os seguintes resultados:
  - Ao nível da disposição do texto na página, temos desformatação parcial, nomeadamente no tipo de letra e na disposição das linhas que não estão totalmente emparelhadas com o ficheiro de partida.
  - Mantém algumas marcas de formatação do texto de partida (e.g. imagens, *links*, títulos etc.)
  - Ao nível do reconhecimento de caracteres obtivemos alguns problemas:
    - caracteres diferentes do original (acrescenta ou perde caracteres que não existem no original)
    - não reconhece caracteres especiais (e.g. parêntesis)
    - Erros de reconhecimento ortográfico (i.e. diacríticos, espaços a mais ou a menos)
    - Exemplos dos erros listados:
      - nã conse guem → não conseguem
      - P( orto) → (Porto)
      - ~~abriu no Mercado de São Domingos de Benfica~~ → abriu, no Mercado de São Domingos de Benfica, em

### **Exemplo 3 (de PDF para DOCX)**

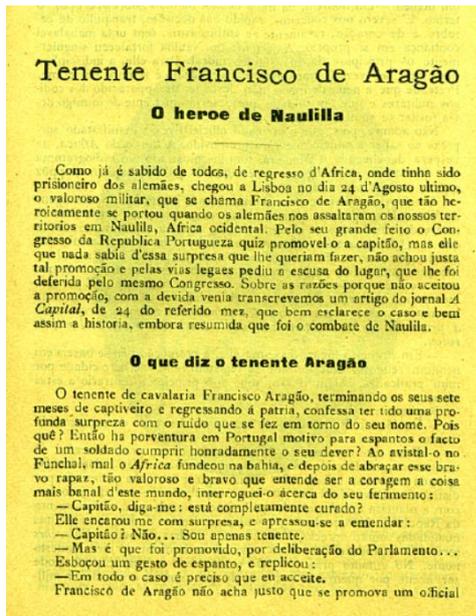
---

- Descrição: Ficheiro de uma página do jornal *O Século Ilustrado*, em formato papel.
- Esta imagem encontra-se disponível em <https://purl.pt/39776/2/> e está no formato PDF.

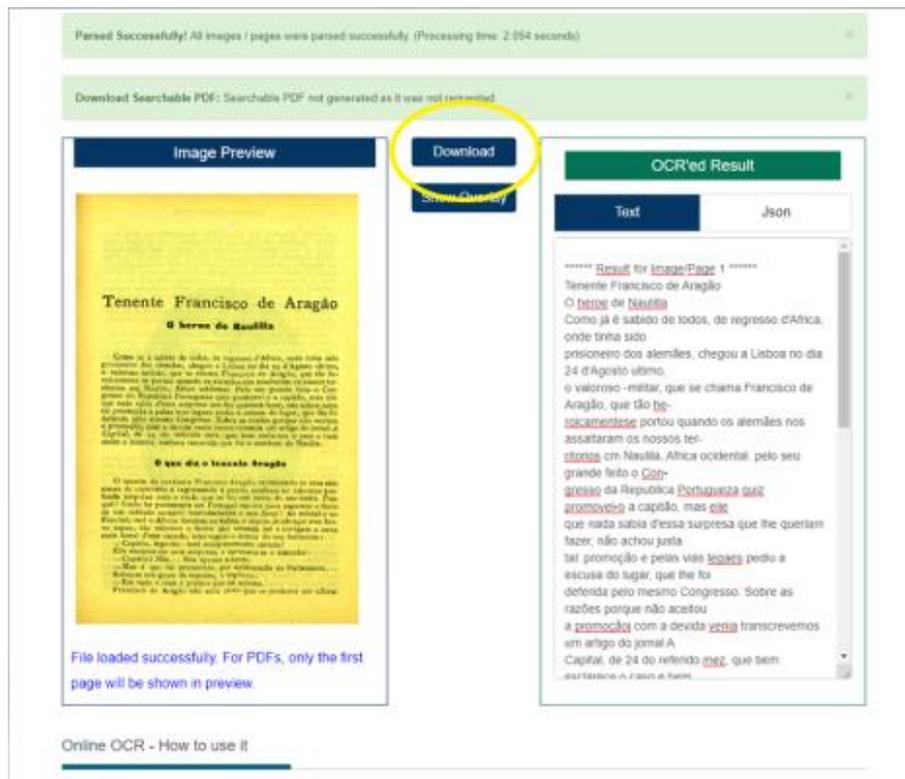


- Ferramenta utilizada: ONLINE OCR (<https://www.onlineocr.net/pt/>). Apesar de ser uma ferramenta gratuita tem restrições quanto ao tamanho do ficheiro de partida.
- Esta ferramenta aceita como ficheiros de partida: JPEG, PNG ou PDF.
- Esta ferramenta disponibiliza como ficheiros finais: WORD, EXCEL, TXT.





- Ferramenta utilizada: OCR.SPACE (<https://ocr.space/>). Apesar de ser uma ferramenta gratuita há restrições quanto ao tamanho do ficheiro de partida.



- Esta ferramenta aceita como ficheiros de partida: JPEG, PNG ou PDF.
- Esta ferramenta disponibiliza como ficheiros finais: TXT para download (apesar de permitir a visualização *pre-view* numa caixa de texto na página *online*).

```

1
2 ***** Result for Image/Page 1 *****
3 Tenente Francisco de Aragão
4 O heroe de Naulilla
5 Como já é sabido de todos, de regresso d'Africa, onde tinha sido
6 prisioneiro dos alemães, chegou a Lisboa no dia 24 d'Agosto ultimo,
7 o valoroso -militar, que se chama Francisco de Aragão, que tão he-
8 roicamente portou quando os alemães nos assaltaram os nossos ter-
9 ritorios em Naulila, Africa ocidental. pelo seu grande feito o Con-
10 gresso da Republica Portugueza quiz promovel-o a capitão, mas elle
11 que nada sabia d'essa surpresa que lhe queriam fazer, não achou justa
12 tal: promoção e pelas vias legais pediu a escusa do lugar, que lhe foi
13 deferida pelo mesmo Congresso. Sobre as razões porque não aceitou
14 a promoçãoj com a devida venia transcrevemos um artigo do jornal A
15 Capital, de 24 do referido mez, que bem esclarece o caso e bem
16 assim historia, embora resumida que foi o combate de Naulila.
17 O que diz o tenente Aragão
18 O tenente de cavalaria Francisco Aragão, terminando os seus sete
19 meses de cativo e regressando á patriá, confessa ter tido uma pro-
20 funda -surpresa com o ruido que se fez em torno do seu nome. Pois
21 uê ? Então ha porventura em Portugal motivo para espantos o facto
22 e um soldado cumprir honradamente o seu dever? Ao avistal-o no
23 Funchals mal o Africa fundeu na bahia, e depois de abraçar esse bras
24 vo ra aza, tão valoroso e bravo que entende ser a coragem a coisa
25 mais anal d'este mundo, T interroguei-o ácerca do seu ferimento :
26 -Capitão, diga-me : está completamente curado?
27 Elle encaroune com surpresa, e apressou-se a emendar :
28 Não.
29 Sou apenas tenente.
30 .-Mase é que foi promovido, por deliberação do Parlamento.
31 Esboçou um gesto de espanto, e replicou :
32 '-Em todo o caso é preciso que eu aceite.
33 Francisco de Aragão não acha nustO que se promova um oficial
34

```

- Relativamente à imagem submetida obtivemos os seguintes resultados:
  - Ao nível da disposição do texto na página, a ferramenta acrescenta linhas de metadados que não existem no texto de partida – e.g. “Result for image/Page 1). Mantém parcialmente a formatação do texto de partida (e.g. marca de fim de linha) mas perde-se a identificação dos títulos.
  - Apresenta erros de reconhecimento ortográfico (e.g. hifenização, diacríticos, acrescenta/elimina espaços que não existem no texto original)
  - Não corrige a grafia antiga (e.g. d'Agosto, mez).
  - Exemplos do erros listados:
    - he – roicamente → heroicamente
    - Com gresso → (acrescenta espaços no meio das palavras e hiper-corrige a preposição “com”) congresso
    - patriá → pátria

Através da apresentação dos diferentes exemplos e da utilização das diferentes ferramentas esperamos que seja mais fácil em trabalhos futuros saber escolher qual a ferramenta que melhor se adequa ao seu trabalho e qual apresenta melhores resultados para que possa executar esta tarefa da melhor forma.

*Equipa Hugod*



<b>EN</b>	The European Commission's support for the production of this publication does not constitute an endorsement of the contents, which reflect the views only of the authors, and the Commission cannot be held responsible for any use which may be made of the information contained therein.
<b>PT</b>	O apoio da Comissão Europeia à produção desta publicação não constitui um aval do seu conteúdo, que reflete unicamente o ponto de vista dos autores, e a Comissão não pode ser considerada responsável por eventuais utilizações que possam ser feitas com as informações nela contidas.