

Humanities Going Digital (HUGOD)

[2020-1-CZ01-KA226-HE-094363]

Compilação de Corpora -
OCR

Sessão 04

12/10/2022

Chiara Barbero
Sílvia Barbosa



AULAS

Introdução



01

02

Linguística de
Corpus



Ferramentas
computacionais



03

Compilação de
Corpora



Digitalização



04

05

OCR



06



OCR (*Optical Character Recognition* | Reconhecimento Óptico de Caracteres)

O que é o OCR?

- O OCR permite a conversão de documentos em diversos formatos (GIF, JPG, PNG, TIFF e PDF) em dados em formato de texto que podem ser pesquisados e editados.

Como funciona?

- A ferramenta analisa o documento, compara os caracteres com fontes armazenadas e/ou reconhece características típicas de determinado caractere.



OCR (*Optical Character Recognition* | Reconhecimento Óptico de Caracteres)

Como funciona o reconhecimento do OCR?

- compara cada caractere identificado previamente com uma base de símbolos para definir padrões e encontrar semelhanças.
- depois da identificação e definição dos caracteres, o OCR compara as informações extraídas com uma base de palavras da língua em questão, de forma a confirmar ou não os dados extraídos.

Quais as vantagens?

- permitir o acesso completo ao textos para pesquisa, edição e armazenamento
- reduzir ao mínimo a intervenção manual
- possibilitar a conversão para diferentes tipos de formatos



HWR / HTR

(Handwritten Recognition / Handwritten Text / Recognition |
Reconhecimento Óptico de Caracteres)

O que é o HWR/HTR?

- Ferramenta para reconhecimento de textos manuscritos
- Utilizada em projetos que lidem com documentação histórica (cartas, diários, ...)



- <https://readcoop.eu/transkribus/>

The screenshot displays the Transkribus web interface. On the left, a scanned image of a handwritten document is shown with a toolbar for editing. The text on the document is in German cursive. On the right, the digital transcription of the document is displayed, with the text converted into a structured format. The transcription includes the title 'Germteig.' and the main text: 'In laue Milch, Germ hinein, und etwas Mehl, von den 50 dkg Mehl versprudeln, u. am Herd lau machen u. aufgehen lassen. 50 dkg Mehl 1-2 dkg Germ 1/8 l laue Milch, salzen Verfeinerung 2 Eier 4 dkg Butter ode. Fett Panille, Zitronengeschmack. Für Milchbrot, Kipferl, Gugelhupf, Strudel.'





OCR e Corpora

Quando recorrer ao uso de OCR para a compilação de corpora?

- Quando for necessário recorrer à digitalização, para a transformação das “imagens” até se tornarem “textos”
- Na consulta e uso de materiais presentes em repositórios em formato eletrónico mas não editável (muito frequente!!)



- Porque a maioria das ferramentas para exploração de corpora exige textos em formato editável, maioritariamente em formato de texto simples (.txt)



Captura de Imagem



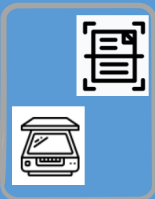
Como são as imagens?

- formato
- tamanho
- cores
- ...



Que ferramentas?

- digitalizador/*scanner*
- máquina fotográfica
- telemóvel/*smartphone*



Captura de Imagem

Formatos


- imagem com texto a cores (TIFF/JPEG2000, 300 ppi resolução, 24bit, sem compressão)
- imagem com texto p/b (TIFF/JPEG2000, 300 ppi resolução, 8bit, sem compressão)
- JPEG/PNG também são possíveis

Otimização da imagem obtida (alterações para melhor OCR)

- corte das áreas não necessárias
- orientação vertical/horizontal (para melhorar a leitura)
- formatação (das cores, da tonalidade, das curvaturas)
- limpeza de marcas (borrões, marcas, outros problemas)

Gravação da nova imagem

- descrição dos procedimentos usados para replicar nas seguintes no formato desejado



Como são as imagens?

- formato
- tamanho
- cores
- ...



Digitalizador / Scanner
Máquina fotográfica
Telemóvel / smartphone

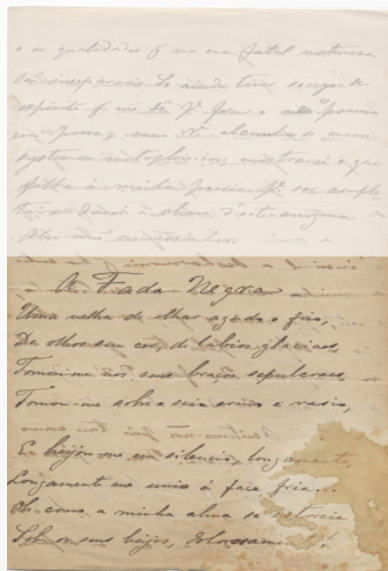


Captura de Imagem: exemplos

Existe um conjunto diversificado de material em diferentes formatos que pode ser alvo de digitalização

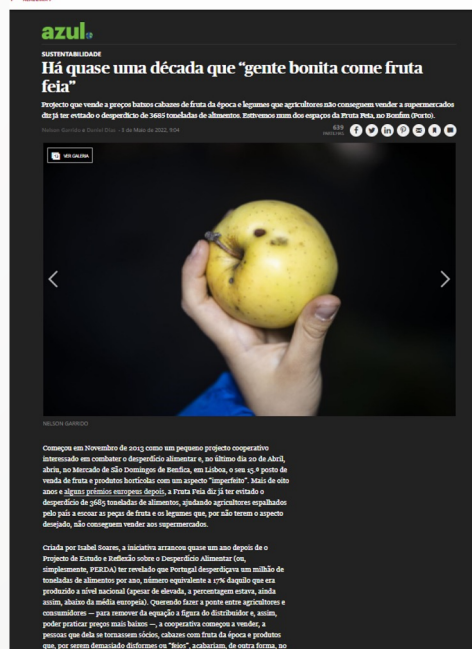


A



Tipo image/ibex Tamarho 712 KB

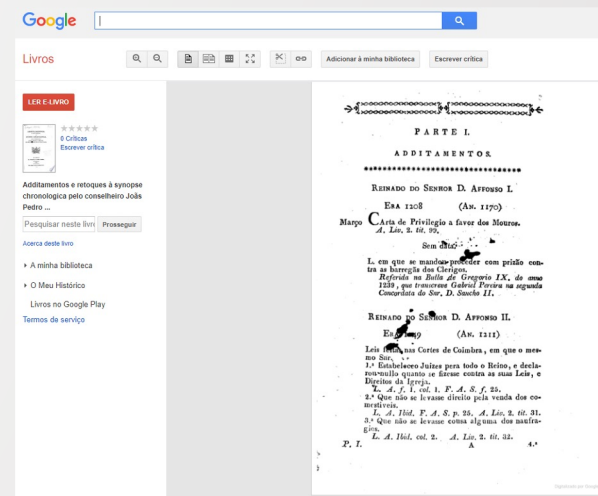
B



C



D



E



F

Que ferramenta usar?

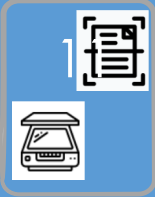
Como escolher?

- De acordo com o objeto - documento manuscrito (HTR) ou ficheiro eletrónico (OCR)
- De acordo com o formato de ficheiro eletrónico (imagem ou PDF)
- De acordo com a língua de trabalho
- De acordo com a formatação do texto (texto corrido, diferentes colunas..)
- De acordo com a variedade linguística em uso (contemporânea ou grafias arcaicas) e o alfabeto (cirílico, hebraico, etc.) ou sistema de caracteres (japonês, coreano, etc.)
- De acordo com as licenças à nossa disposição (softwares pagos)



Testar sempre uma amostra e avaliar os resultados para escolher a ferramenta que melhor se adapta à nossa investigação





Tratamento da imagem

Programas/software de HTR

- Transkribus:
(<https://readcoop.eu/transkribus/>)
- Amazon Textract
(<https://aws.amazon.com/pt/textract/>)

Programas / software de OCR

Pagos

- Adobe Acrobat Pro DC
- PDF Reader
- OmniPage Ultimate
- ABBYY FineReader PDF
- SimpleOCR
- Tesseract

Gratuitos

- Online OCR (<https://www.onlineocr.net/pt/>)
- Tools PDF24 (<https://tools.pdf24.org/pt/ocr-pdf>)
- OCR Space (<https://ocr.space/>)
- OCR 2 Edit (<https://www.ocr2edit.com/pt/converter-para-txt>)
- New OCR (<https://www.newocr.com/>)



Avaliação dos resultados/problemas mais frequentes

Problemas de ordem gráfica:

- falta de sinais diacríticos (ou sinais diacríticos manipulados, acentuação errada)
- falta de caracteres
- caracteres a mais que não existem no original
- falta de espaços brancos entre duas palavras
- espaços brancos a mais no meio das palavras
- hifenização alterada

Problemas de formatação:

- desformatação geral do texto
- manipulação de metadados (desaparecem ou acrescentar marcas do *software*)
- reorganização do texto (no caso de textos com mais de uma coluna) que não respeita a ordem original

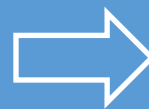


Trabalho em aula – apresentar no dia 17 os resultados

- (1) Selecionar a imagem (1 imagem por grupo) / (2) selecionar um conjunto de imagens
- (1) Testar a imagem com *softwares* diferentes / (2) testar todas as imagens do conjunto selecionado num dos software propostos
- Observar os resultados (pontos forte e os pontos fracos)
- Escolher o *software* que melhor se adapta
- Anotar o procedimento metodológico

IMAGEM

? Formato, ? tamanho, ?
Cores, ? Caracteres, ...



TEXTO EDITÁVEL

(resultado é um texto num ficheiro com extensão .TXT)

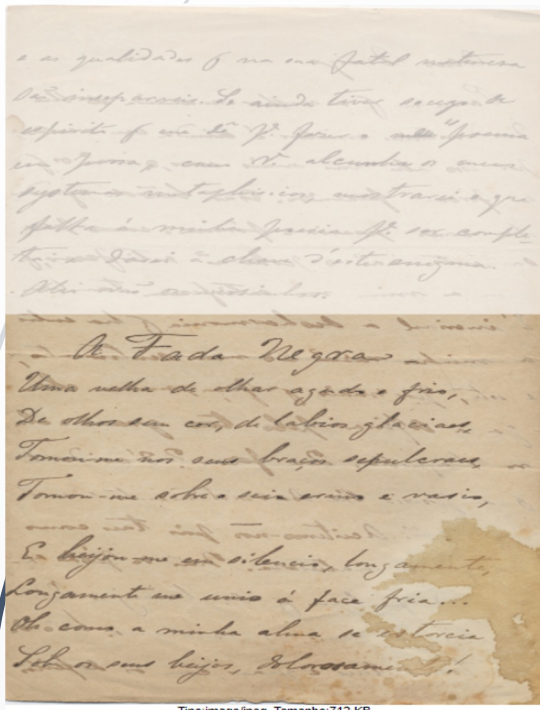
“O choco de Setúbal não é o único a estar frito.”



EXEMPLOS DE IMAGENS

- Existe um conjunto diversificado de material em diferentes formatos que pode ser alvo de digitalização

G3

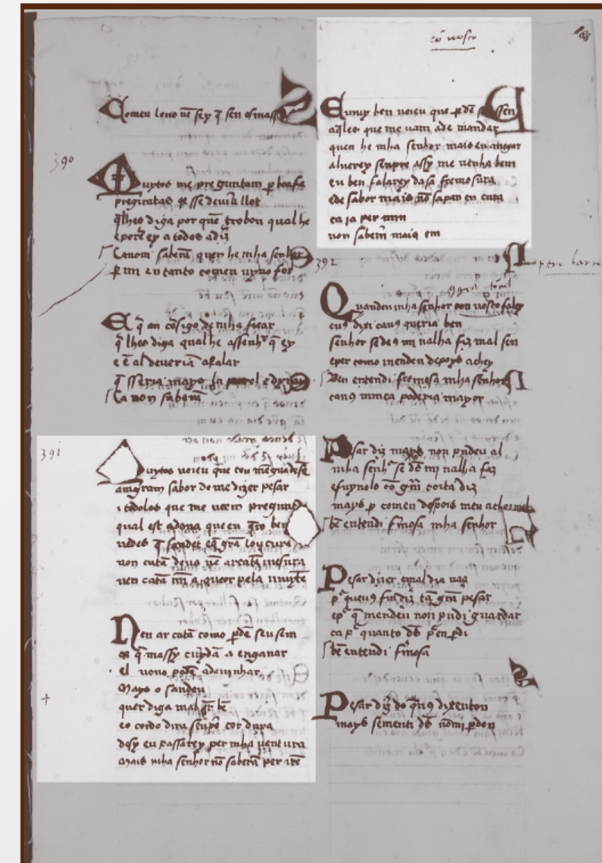


Tipo: image/jpeg. Tamanho: 712 KB

G1



G2





EXEMPLOS DE IMAGENS

- Existe um conjunto diversificado de material em diferentes formatos que pode ser alvo de digitalização

G4



* tejo

ADEGA COOPERATIVA DE ALCANHOES, C.R.L.
Alcanhões, Tel: 243 429 151
Fax: 243 429 78
E-mail: adega.coop.alcanhoes@sapopt
Website: www.adg.alcanhoes.pt
Adega Cooperativa de Alcanhões

ADEGA COOPERATIVA DE ALMEIRIM
Almeirim, Tel: 243 57 0560
Fax: 243 570 561
E-mail: geral@adegaalmeirim.pt
Website: www.adg.almeirim.pt
Adega Almeirim

14,50 € 1,75
TERRAS DO PAÇO
Reg. Tejo Reserva tinto 2013
Castelão, Trincadeira e Aragonez. Aroma bem composto, com matos secos, fruta vermelha fresca, notas mineiras como barro húmido. Macio e feito, com uma nota de madeira verde a dar uma rugosidade vegetal. (13,5%)

14,50 € 1,59
AA BRANCO LEVE
Reg. Tejo branco 2014
Muito leve e floral, com citrinos puros e algumas notas de pedras. Ligeiro gás na boca, doçura perceptível integral, com acidez alta. Leve e intenso, muito fácil, franco e directo, com boa presença. (10%)

14,50 € 1,68
PLANÍCIE
Reg. Tejo branco 2014
Fervor Pires e treme. Frutos amarelos, resinas doces, tostados discretos, citrinos suaves. Ligeiro, algo diluído, com sabor discreto e algum equilíbrio, termina floral e ligeiramente adocicado. (12%)

15,50 € 2,99
VARANDAS
Do Tejo branco 2014
Notas de fruta amarela com boa expressividade, cheio na boca, maduro e de acidez presente mas discreta. Um branco volumoso e de boa presença, para pratos com mais peso. (13%)

15,50 € 5
VARANDAS
Do Tejo Grande Escolha branco 2014
Chardonnay e Arioto parcialmente fermentados em barrica. Frutos amarelos, fumados, citrinos discretos. Ligeiro mas com boa definição e equilíbrio, ligeira rugosidade, final contido, muito versátil. (13%)

14,50 € 3,89
VARANDAS
Do Tejo tinto 2012
Aroma discreto a fruta escondida, notas herbáceas como hortelã ou poejo, muito suave. Tanninos marcantes, boa acidez, austero e sólido, final seco e áspero, a pedir pratos de substância. (13,5%)

G5

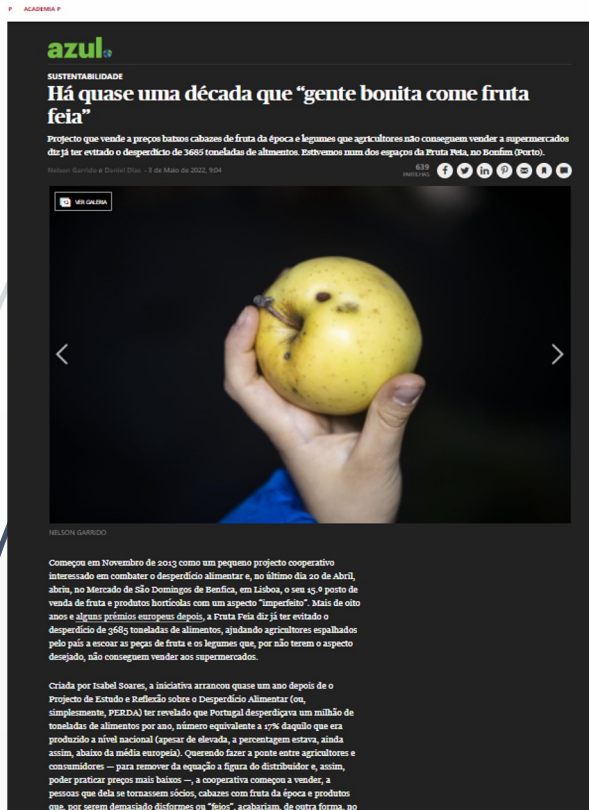
G6





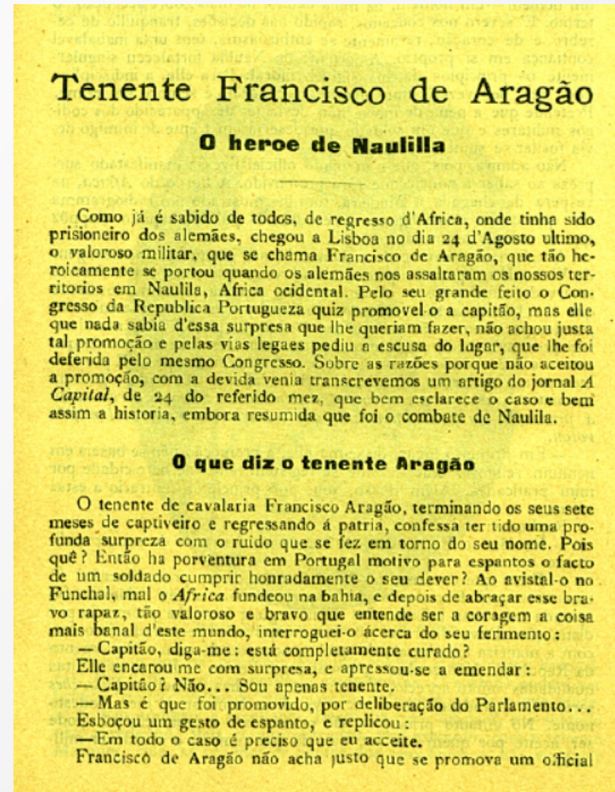
EXEMPLOS DE IMAGENS

- Existe um conjunto diversificado de material em diferentes formatos que pode ser alvo de digitalização



G7

G8



G9



EXEMPLOS DE IMAGENS

- Existe um conjunto diversificado de material em diferentes formatos que pode ser alvo de digitalização

Gazela

GAZELA VINHO VERDE BRANCO

TIPO: Trankão. COR: amarelo
 DESIGNAÇÃO DE ORIGEM: DOC Vinho Verde. REGIÃO: Vinhos Verdes
 PAÍS DE ORIGEM: Portugal

O VINHO
 Gazela é um Vinho Verde de sabor suave, leve e refrescante, que sabe bem com a vida. Gazela é um vinho descomprometido, um clássico renovado, que confirma todo o potencial da região para criar vinhos leves e cativantes, ideais para o dia-a-dia.

NOTAS DE PROVA PROVA DE 2020
 Gazela tem uma cor muito leve, com um ligeiro deprimimento de gás que realça os aromas, sabores e frescura do vinho. Aromático e cativante, Gazela dá as notas de citrinos e frutas tropicais uma doce nota e estimulante. O resultado é um vinho elegante, versátil e muito atractivo.

ENÓLOGO: Diogo Sepúlveda

CASTAS: Loureiro, Ovar, Pedernã, Trajadura

VINIFICAÇÃO
 As uvas são desmontadas e esmagadas suavemente. O mosto resultante é separado das películas em pressas pneumáticas e sujeito a decantação estática durante 24 horas, devidamente protegido das condições, até atingir o grau de limpidez desejado. Segue-se a fermentação em cubas de aço inox, sob uma temperatura controlada de aproximadamente 18°C.

MATURACÃO
 Gazela é engarrafado imediatamente após a fermentação e tone, por forma a garantir toda a sua frescura inicial.

GUARDAR
 Gazela deve ser armazenado de pé, em local seco e fresco. Dado a sua frescura, Gazela é um vinho que ganha em ser consumido de imediato.

SERVI- R
 Gazela deve ser servido bem fresco a uma temperatura entre 6°C-8°C.

DESFRLUTAR
 Gazela é ideal para acompanhar copos amplos em bons momentos de diversão.

DETALHES TÉCNICOS
 Alcool: 11,5% (vol) | Açúcar Total: 5,9 g/L (doç. natural) | Açúcares Totais: 11 g/L | pH: 3,2 (doç.)
INFORMAÇÃO NUTRICIONAL (VALORES TÍPICOS PARA 100 mL)
 Alcool: 11,5 g | Açúcar: 11 g | Valor Energético: 249 kJ/59 kcal | Adequado para Vegetarianos Sim | Adequado para Veganos Não | Sem Glúten

ALERGÉNICOS:
 Contém sulfite.

DATA ENGARRAFAMENTO: 2020-09-01
CAPACIDADES DISPONÍVEIS: 1.500 mL, 1.000 mL, 750 mL, 375 mL, 187 mL
ENGARRAFADOR:
 Engarrafado por: Segrape Vinhos, S.A., Avintes, Portugal



G10

COMPARTILHE

Conceito de resiliência surgiu na física, mas é aplicado em diversas áreas do conhecimento

Se preferir, vá direto ao ponto [\[Esconder\]](#)

1. Resiliência psicológica
2. Resiliência física
3. Resiliência ambiental
4. Resiliência climática

Resiliência, do latim "resiliens (re + salio)" significa "saltar para trás", recuperar-se, voltar ao "normal". O conceito é aplicado em diversas áreas do conhecimento, incluindo a área ambiental, climática, física e psicológica. Mas foi cunhado em 1807 pelo físico Thomas Young, que o definiu como a capacidade de um objeto, material ou corpo de sofrer pressão ou impacto e, depois, voltar à forma original.

[NO ASSUNTO](#) Cidades desenvolvem estratégias de resiliência climática

Após: [Rocha](#)

Saiba onde descartar seus resíduos

O QUE PRECISA DESCARTAR?
 Seleccione o objeto

ONDE VOCE DESEJA DESCARTAR?

SEU E-MAIL
 Digite seu e-mail

Concordo em receber comunicações e ofertas?

Buscar onde descartar

G11

Google

Livros

LER E-LIVRO

★★★★★
 0 Críticas
 Escrever crítica

Aditamentos e retoques à synopse chronologica pelo conselheiro João Pedro ...

Pesquisar neste livro

Acerca deste livro

- A minha biblioteca
- O Meu Histórico

Livros no Google Play

Termos de serviço

PARTE I.
ADDITIONENTOS.

REINADO DO SENHOR D. AFFONSO I.
 ESA 1208 (AN. 1170)
 Carta de Privilégio a favor dos Mouros.
A. Lio. 2. tit. 95.
 Sem data.

L. em que se mandou prender com prisão contra as barregas dos Clerigos.
Referida na Bula de Gregorio IX, do anno 1239, que transcreve Gualtero Testa na segunda Concordata do Sr. D. Sancho II.

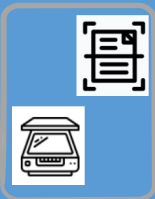
REINADO DO SENHOR D. AFFONSO II.
 ESA 1219 (AN. 1211)
 Lei dada nas Cortes de Coimbra, em que o mesmo Sr. ...

- 1.º Estabelece Juizes para todo o Reino, e declarou-nillo quanto se fizesse contra as suas Leis, e Direitos da Igreja.
- 2.º Que não se levasse direito pela venda dos comreividos.
- 3.º Que não se levasse cotna alguma dos naufragios.

L. de f. l. col. 1. F. A. S. f. 25.
L. A. lio. f. A. S. p. 25. A. Lio. 2. tit. 31.
L. A. lio. col. 2. A. Lio. 2. tit. 32.

P. I. A 4.

G12



EXEMPLOS DE IMAGENS DADAS EM AULA

Fontes das imagens recolhidas e apresentadas:

- G1 - https://purl.pt/14355/1/obras/n48/n48_item59/P1.html
- G2 - <https://tinyurl.com/ycxc7u7h>
- G3 - <https://cantigas.fcsh.unl.pt/cantiga.asp?cdcant=364&pv=sim>
- G4 - <https://amusearte.hypotheses.org/6492>
- G5 - <https://vdocuments.com.br/guia-de-compras-2016-revista-de-vinhos.html?page=1>
- G6 - <https://purl.pt/39776/2/>
- G7 - <https://tinyurl.com/2s3pc3ds>
- G8 - <https://tinyurl.com/y2baspze>
- G9 - <https://www.fgp.pt/store/esteva-tinto>
- G10 - <https://www.datocms-assets.com/33016/1664446617-gazela-vinho-verde-branco.pdf>
- G11 - <https://www.ecycle.com.br/resiliencia/>
- G12 - <https://tinyurl.com/5n6v2knn>



PENSAR PARA O TRABALHO FINAL

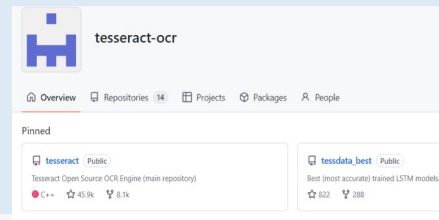
- O tópico?
- Quais documentos a digitalizar?
- Quantos documentos?
- Qual o melhor *software*?
- Esclarecer dúvidas sobre a digitalização ...

Obrigatório: Responder aos Inquéritos 1 e 2 (até dia 14/10/2022)

BIBLIOGRAFIA

Projetos

- <https://rwi.app/iurisprudencia/en/vonglueck>
- <https://rwi.app/iurisprudencia/en/vonglueck>



Palestras online





EN	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.
PT	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.