

Criação e gestão de sites preserváveis

Recomendações do Arquivo.pt



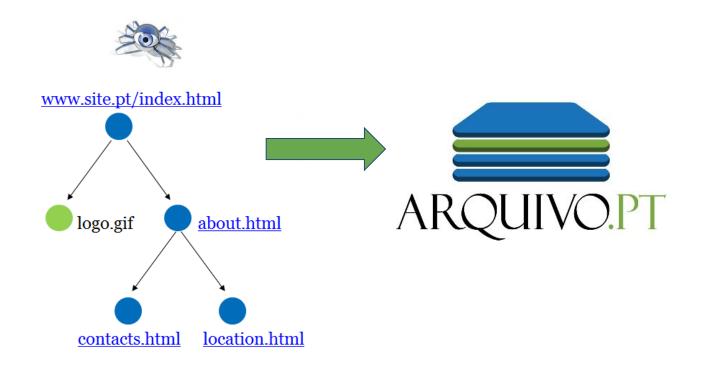
Introdução

Quando se fala em preservar a Web

- Recolher, armazenar e disponibilizar
- Fragmentação inevitável
- Conservação de recursos únicos para memória futura



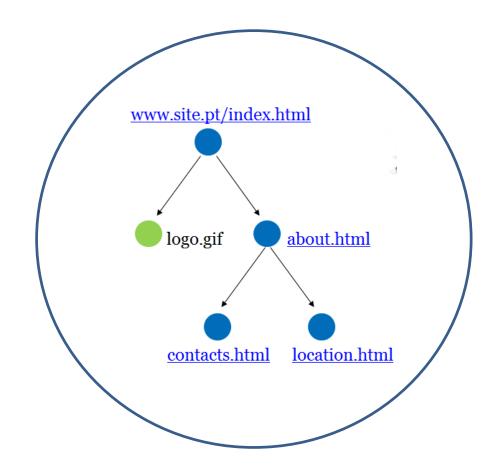
Recolha é feita de forma automática por "robots"



Recolha é feita dentro de certos limites:

Relacionados com o âmbito de recolha do Arquivo.pt

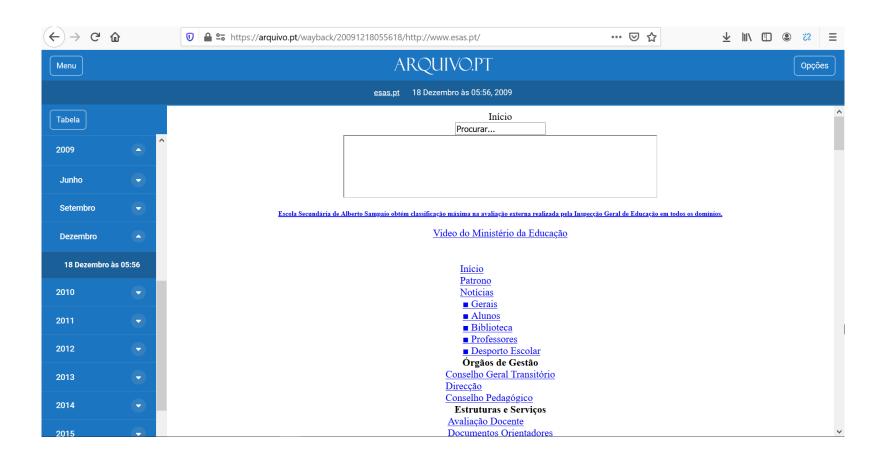
Relacionados com os websites a recolher



Reprodução de páginas preservadas



Nem sempre corre tudo bem!



Recomendações para **Publicar** Informação **Preservável**

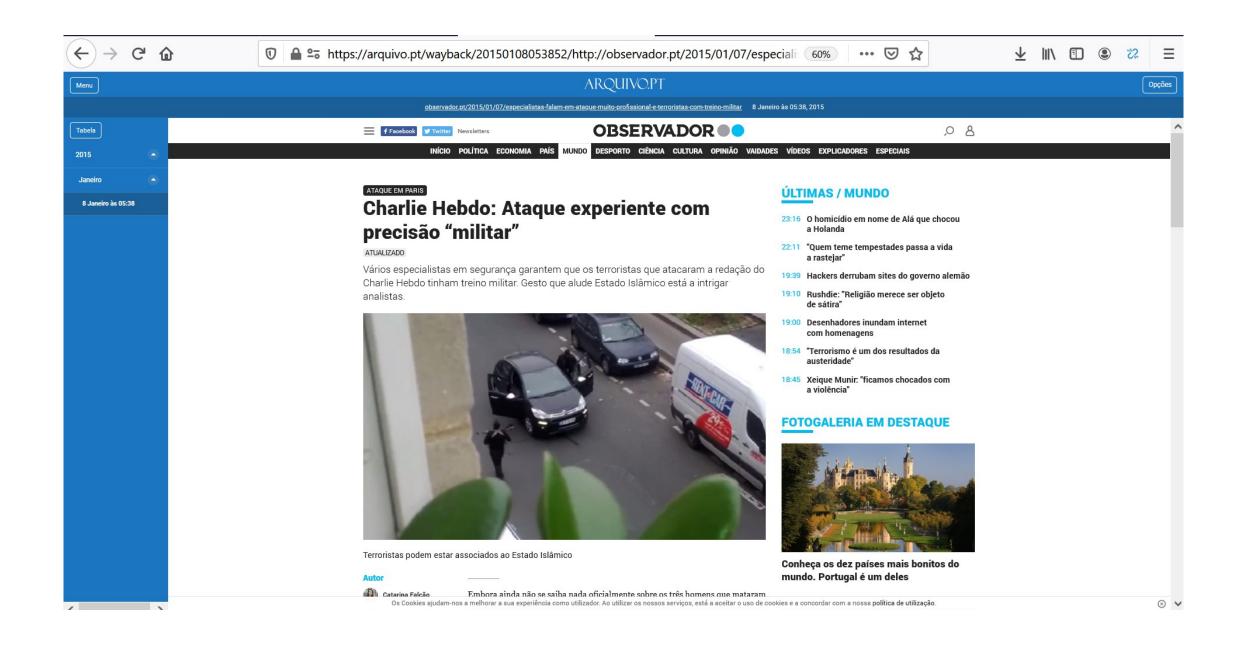
1

Identifique corretamente a data de publicação

Identifique corretamente a data de publicação (#1)







ATAQUE EM PARIS

Charlie Hebdo: Ataque experiente com precisão "militar"

ATUALIZADO

Vários especialistas em segurança garantem que os terroristas que atacaram a redaç Charlie Hebdo tinham treino militar. Gesto que alude Estado Islâmico está a intrigar analistas.

Identifique corretamente a data de publicação (#1)

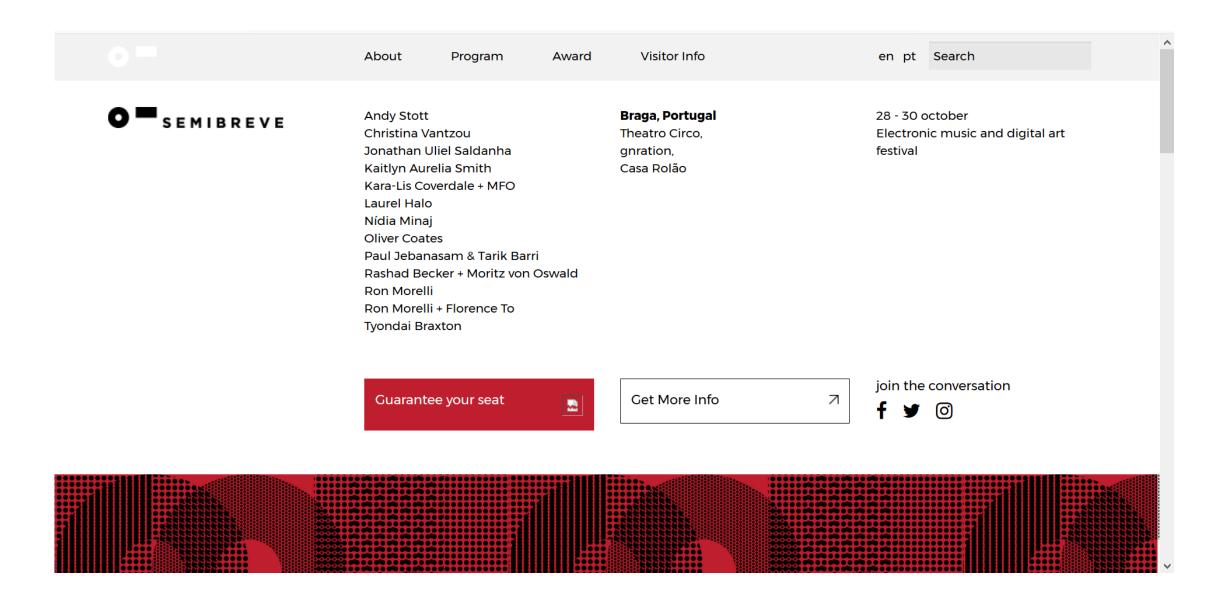




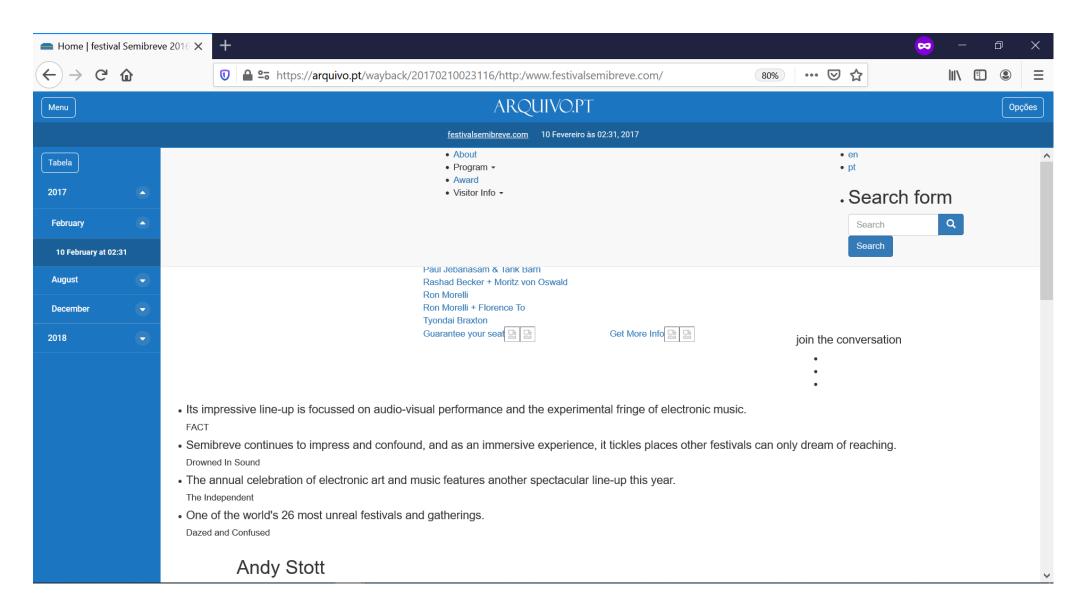


2

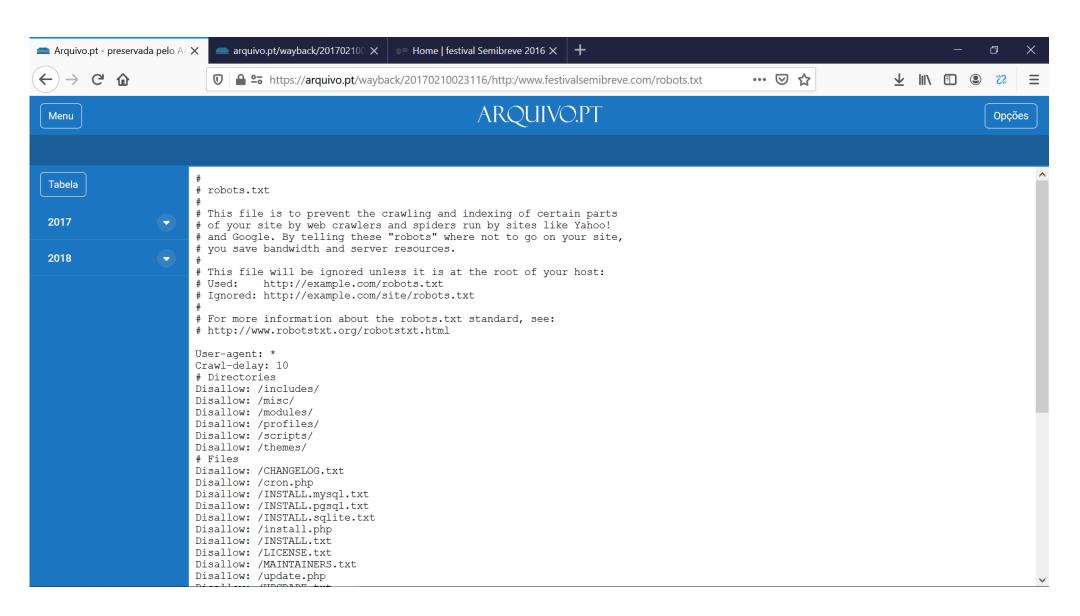
Site do festival Semi-Breve: como era em 2017



Site do festival Semi-Breve: como foi preservado



Robots Exclusion Protocol é a origem do problema de preservação



```
# robots.txt
# This file is to prevent the crawling and indexing of certain parts
# of your site by web crawlers and spiders run by sites like Yahoo!
# and Google. By telling these "robots" where not to go on your site,
# you save bandwidth and server resources.
# This file will be ignored unless it is at the root of your host:
        http://example.com/robots.txt
# Ignored: http://example.com/site/robots.txt
# For more information about the robots.txt standard, see:
# http://www.robotstxt.org/robotstxt.html
User-agent: *
Crawl-delay: 10
# Directories
Disallow: /includes/
Disallow: /misc/
Disallow: /modules/
Disallow: /profiles/
Disallow: /scripts/
Disallow: /themes/
```



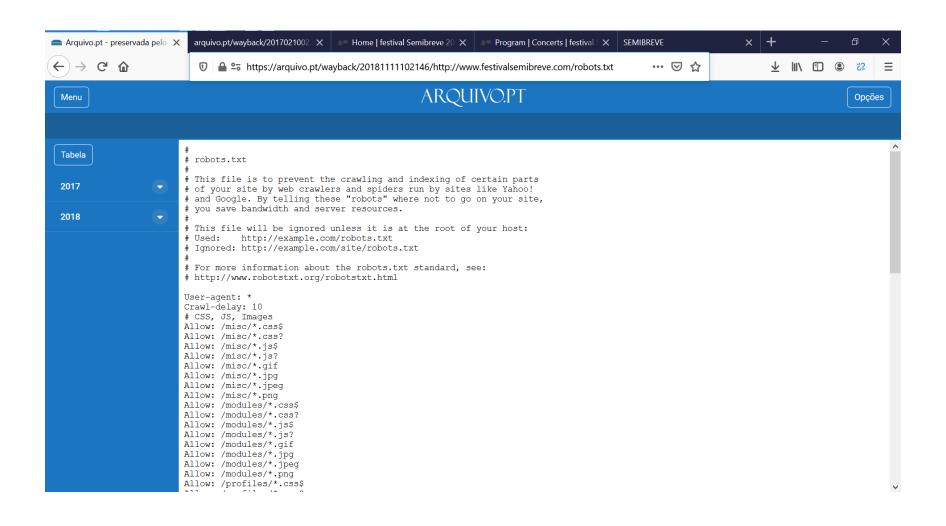
<u> 으</u>

https://arquivo.pt/wayback/20170210023116/http:/www.festivalsemibreve.com/robots.txt

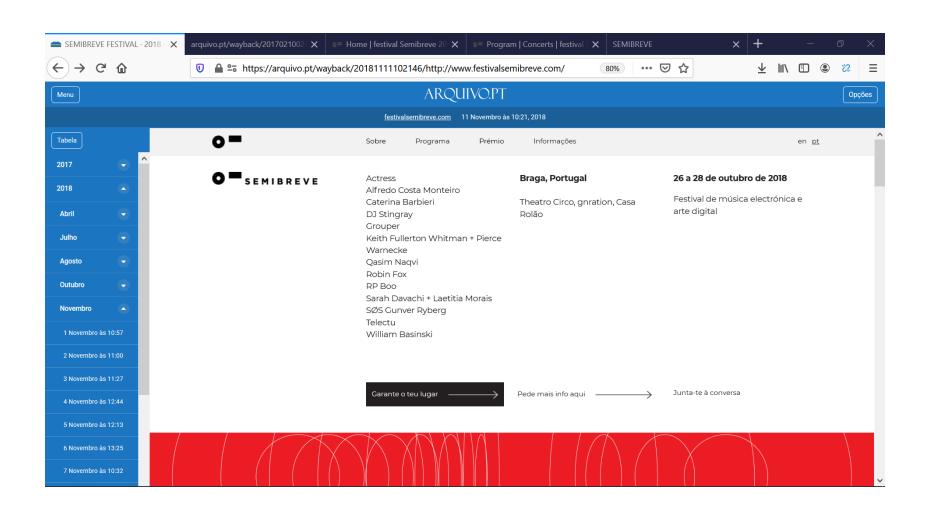
https://arquivo.pt/wayback/20170210023116/http://www.festivalsemibreve.com/robots.txt

User-agent: *

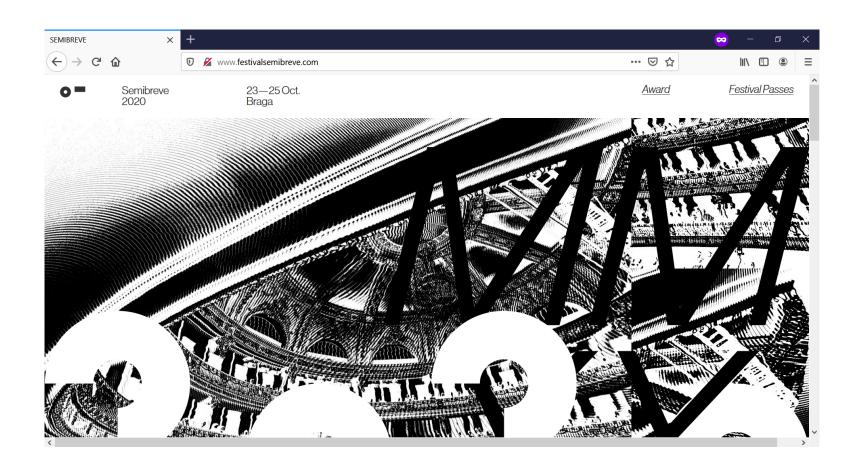
Disallow:



```
# For more information about the robots.txt standard, see:
# http://www.robotstxt.org/robotstxt.html
User-agent: *
Crawl-delay: 10
# CSS, JS, Images
Allow: /misc/*.css$
Allow: /misc/*.css?
Allow: /misc/*.js$
Allow: /misc/*.js?
Allow: /misc/*.gif
Allow: /misc/*.jpg
```



Em alternativa: utilize de um mapa do site



Em alternativa: utilize de um mapa do site

```
-<urlset>
  -<url>
     <loc>http://festivalsemibreve.com/index.html</loc>
     <lastmod>2020-06-04</lastmod>
     <changefreq>weekly</changefreq>
     <priority>0.5</priority>
    -<image:image>
      -<image:loc>
          http://festivalsemibreve.com/images/pasted-svg-84244x36.svg
       </image:loc>
     </image:image>
    -<image:image>
      -<image:loc>
          http://festivalsemibreve.com/images/moshed-2019-11-7-18-56-211024x683.jpg
       </image:loc>
     </image:image>
    -<image:image>
      -<image:loc>
         http://festivalsemibreve.com/images/moshed-2019-11-7-18-56-21 2x.jpg
       </image:loc>
```

http://festivalsemibreve.com/sitemap.xml

Porque continua a haver bloqueios no Robots.txt?

"Disallowing crawling of Javascript or CSS files in your site's robots.txt directly harms how well our algorithms render and index your content and can result in suboptimal rankings."

https://webmasters.googleblog.com/2014/10/updating-our-technical-webmaster.html?m=1



Exclusões pré-definidas pelos Sistemas de Gestão de Conteúdos (CMS) causaram problemas



User-agent: *

Disallow: /administrator/

Disallow: /cache/

Disallow: /components/

Disallow: /editor/

Disallow: /help/

Disallow: /images/

Disallow: /includes/

Disallow: /language/

Disallow: /mambots/

Disallow: /media/

Disallow: /modules/

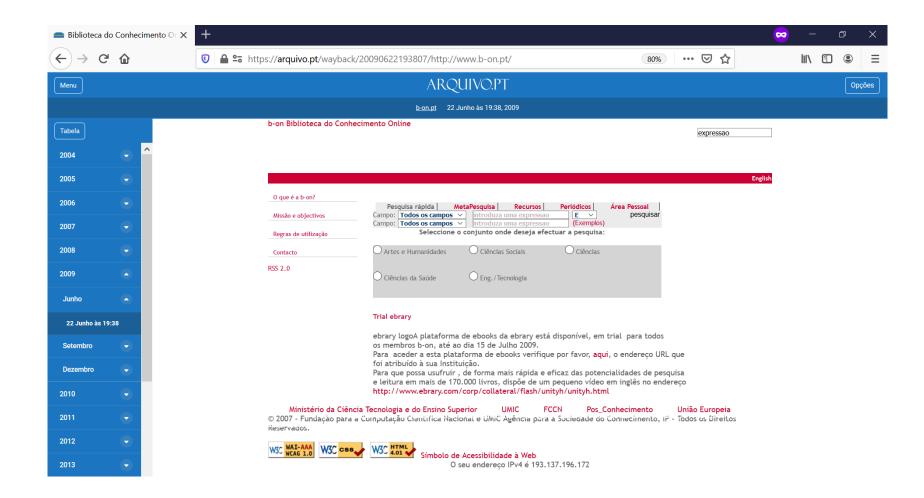
Disallow: /templates/

Disallow: /installation/

Disallow: /dmdocuments/

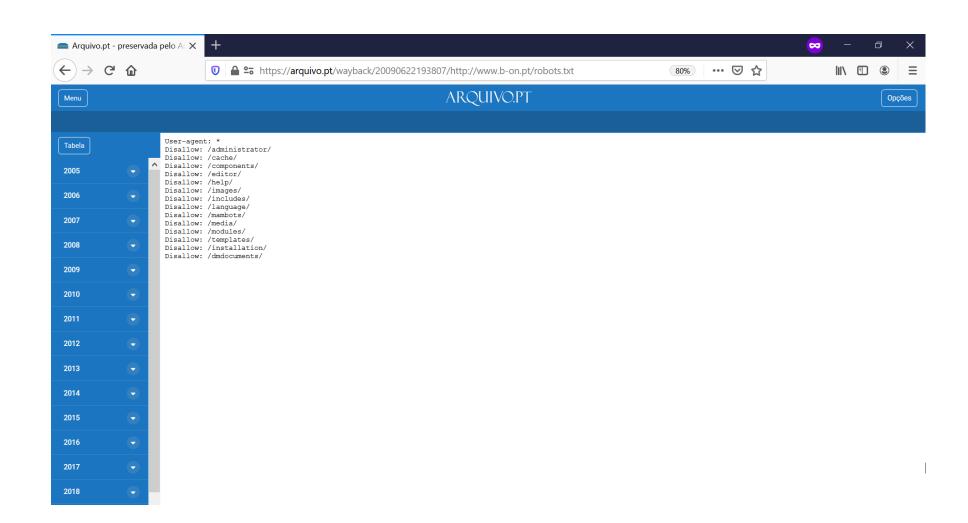
...

B-on.pt: como foi preservado



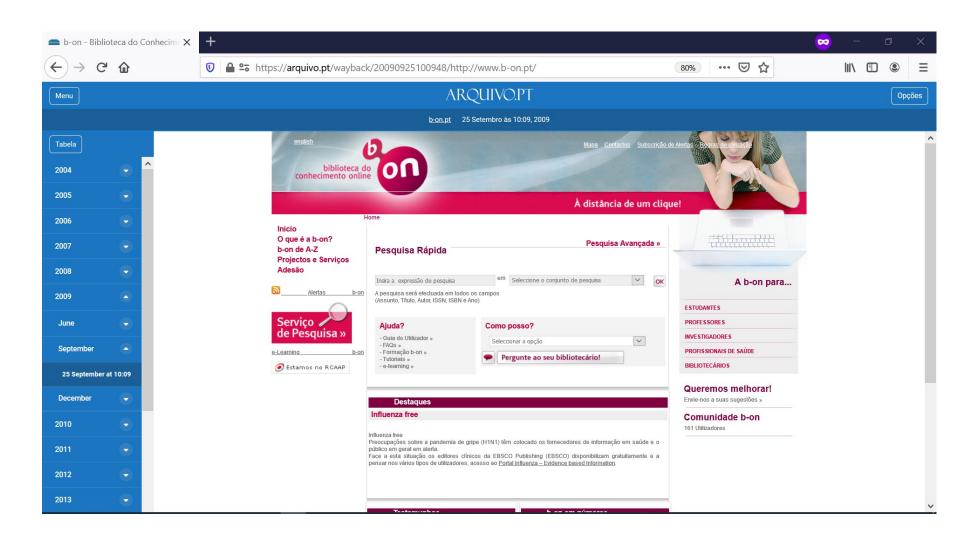
...

B-on.pt: como foi preservado



...

B-on.pt: como foi preservado



O Sistema de Recolha do Arquivo.pt está devidamente identificado.

User-agent: Arquivo-web-crawler

Disallow:



3

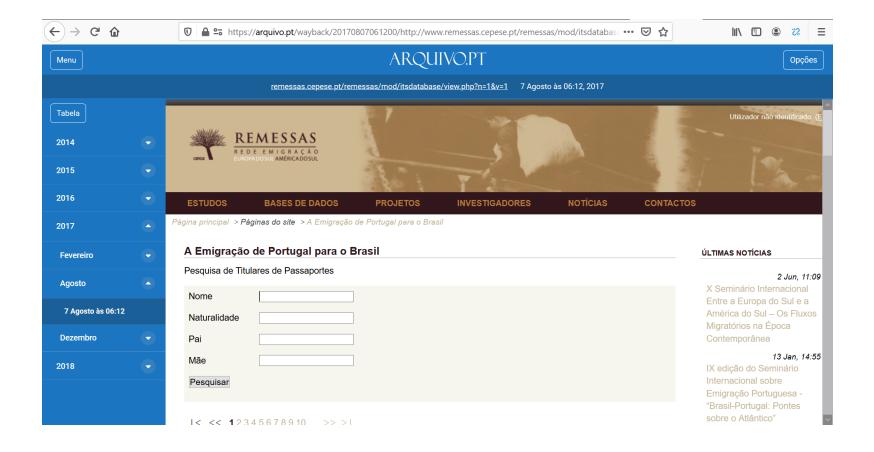
Utilize um endereço para cada conteúdo

Utilize um endereço para cada conteúdo

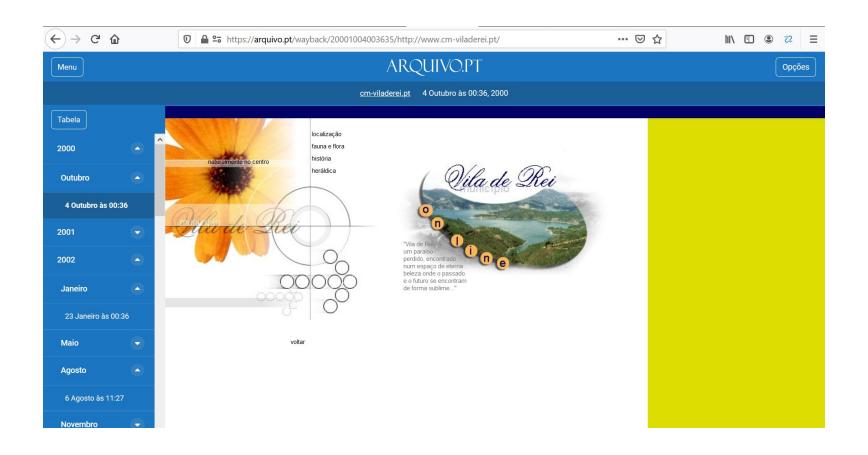
Conteúdos escondidos atrás de formulários escapam às recolhas



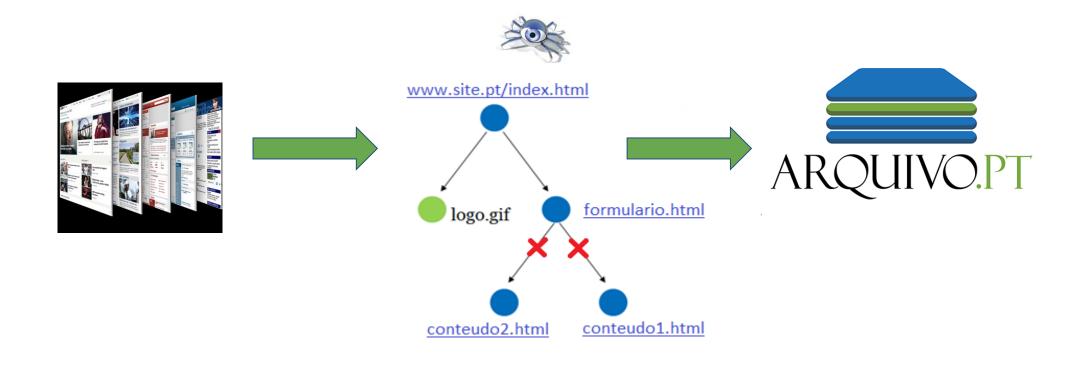
Conteúdos escondidos atrás de formulários escapam às recolhas



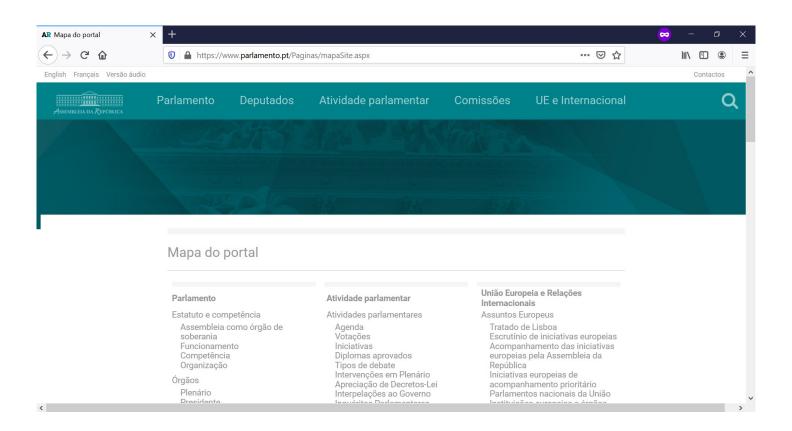
Conteúdos não identificados por um endereço único escapam às recolhas



Conteúdos escondidos atrás de formulários escapam às recolhas



Mapa do Site facilita acesso a pessoas (usabilidade) e máquinas (SEO)



Conteúdos escondidos atrás de formulários escapam às recolhas

Alternativa para recuperar conteúdo perdido



Pesquisar noutros arquivos

4

Mantenha o mesmo endereço ao longo do tempo

Problema: quebra de histórico devido a mudança de endereço do site



Mantenha o histórico redirecionando os endereços antigos para os novos

	Tabela de versões 460 versões de iscte.pt												
2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	
<u>22 Jan</u>	28 Jan	3 Fev	7 Jan	<u>13 Jan</u>	4 Jan	<u>6 Jan</u>	20 Mai	<u>31 Mai</u>	21 Jan	<u>21 Jan</u>	5 Nov	26 Set	
<u>25 Mar</u>	2 Fev	<u>21 Mar</u>	<u>12 Jan</u>	2 Fev	9 Jan	<u>18 Jan</u>	20 Mai	<u>6 Jun</u>	20 Mai	<u>22 Jan</u>	7 Nov	30 Set	
25 Mai	<u>18 Fev</u>	21 Abr	<u>13 Jan</u>	<u>6 Fev</u>	<u>13 Jan</u>	<u>26 Jan</u>	<u>21 Mai</u>	5 Aug	22 Mai				
27 Mai	<u>24 Mar</u>	<u>5 Jun</u>	1 Fev	7 Fev	<u>19 Jan</u>	28 Jan	23 Jun						
<u>5 Jun</u>	<u>27 Mar</u>	<u>11 Jun</u>	4 Fev	2 Mar	<u>21 Jan</u>	<u>6 Fev</u>	24 Set						
<u>24 Set</u>	9 Abr	<u>11 Jun</u>	4 Fev	2 Abr	<u>24 Jan</u>	<u>15 Fev</u>	<u>26 Set</u>						
<u>27 Set</u>	<u>10 Abr</u>	<u>12 Jun</u>	8 Fev	<u>25 Abr</u>	<u>27 Jan</u>	<u>11 Mar</u>	<u>17 Dez</u>						
<u>13 Nov</u>	<u>10 Abr</u>	<u>14 Jun</u>	8 Fev	<u>27 Abr</u>	<u>5 Fev</u>	<u>14 Mar</u>	<u>18 Dez</u>						
<u>26 Nov</u>	<u>31 Mai</u>	<u>23 Jun</u>	9 Fev	2 Jun	<u>20 Fev</u>	21 Out							

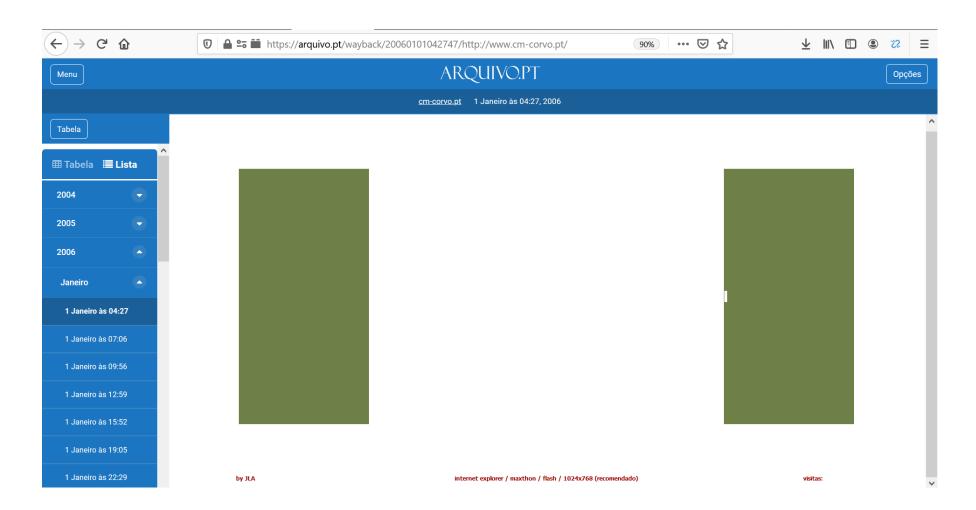
http://iscte.pt

http://iscte-iul.pt

5

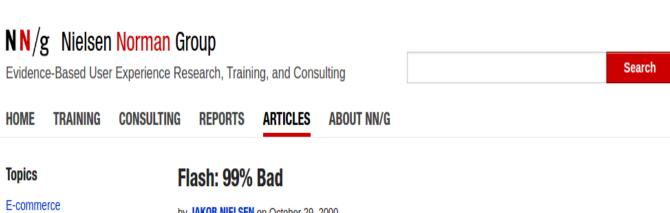
Utilize formatos adequados para preservação

Problema: Utilização de Flash



https://arquivo.pt/wayback/20060101042747/http://www.cm-corvo.pt/

"It breaks with the Web's fundamental interaction principles"



Intranets

Mobile & Tablet

User Testing

Web Usability

See all topics...

Author

Jakob Nielsen

Don Norman

Bruce "Tog" Tognazzini

See all authors...

Recent Articles

by JAKOB NIELSEN on October 29, 2000

Topics: Web Usability

Summary: Although multimedia has its role on the Web, current Flash technology tends to discourage usability for three reasons: it makes bad design more likely, it breaks with the Web's fundamental interaction style, and it consumes resources that would be better spent enhancing a site's core value.

About 99% of the time, the presence of Flash on a website constitutes a usability disease. Although there are rare occurrences of good Flash design (it even adds value on occasion), the use of Flash typically lowers usability. In most cases, we would be better off if these multimedia objects were removed.

Flash tends to degrade websites for three reasons: it encourages design abuse, it breaks with the Web's fundamental interaction principles, and it distracts attention from the site's core value.



http://www.occupyflash.org/

Escolha formatos adequados:

Condições de licenciamento que permitam a sua utilização.

Normas emitidas por um organismo oficial (W3C).

Documentados abertamente através de uma especificação pública.

Lidos e escritos por múltiplas plataformas de software, incluindo código-aberto.

Amplamente usados.

Escolha formatos adequados:

Texto

HTML, XHTML ou XML
Open Document Text (.odt)
PDF/A-1 segundo a norma ISO 19005-1 (.pdf)

Imagem

PNG (.png)
JPEG2000

Video

AVI sem compressão (.avi)

Escolha formatos adequados

Evite formatos não adequados para preservação

Texto:

Microsoft Word (.doc)

Imagem:

Macromedia Flash (*.swf)
PhotoShop (.psd)

Vídeo:

Windows Media Video (.wmv)

Novas formas de inserir os conteúdos nas páginas



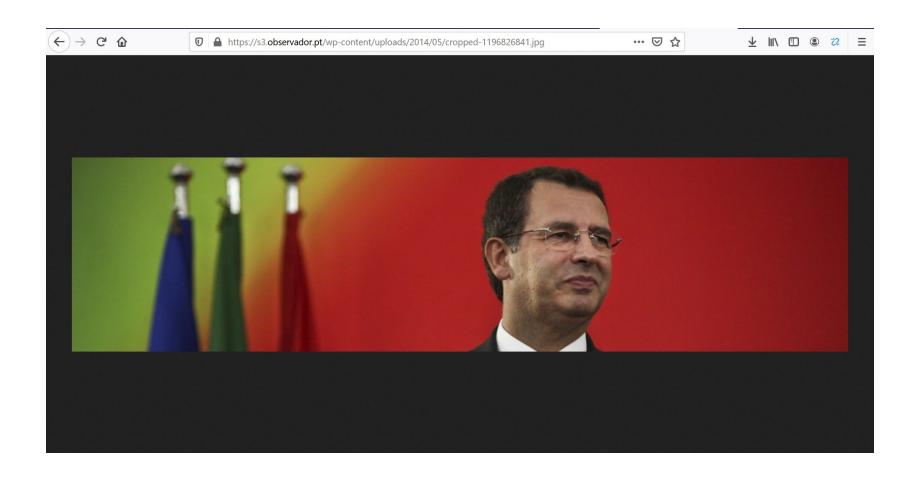
imagem

Novas formas de inserir os conteúdos nas páginas



http://cdn.observador.pt/wp-content/uploads/2014/05/cropped-1196826841.jpg

Novas formas de inserir os conteúdos nas páginas

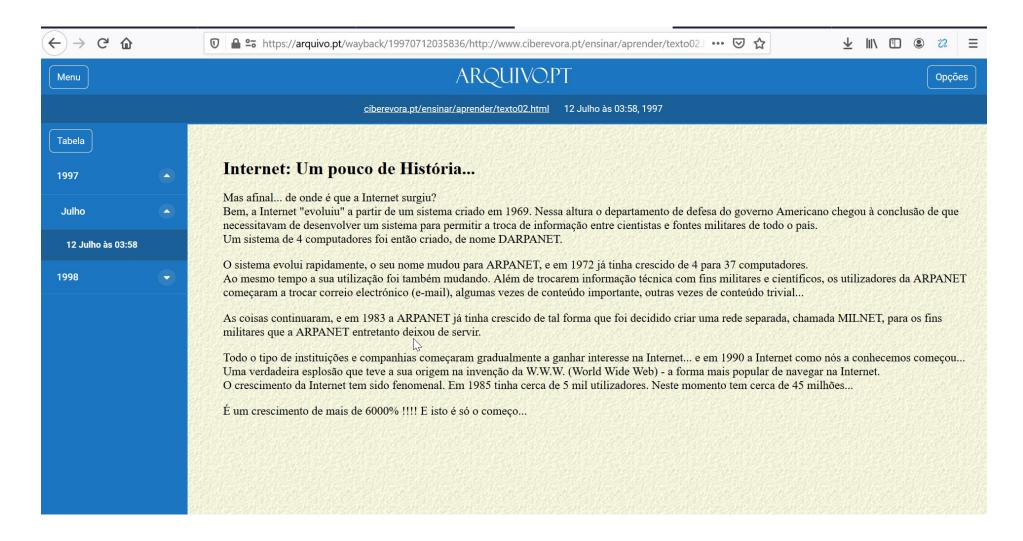


http://s3.observador.pt/wp-content/uploads/2014/05/cropped-1196826841.jpg

6

Utilize metadados para descrever os conteúdos

Quem é o autor da página?



Utilize metadados (Dublin Core)

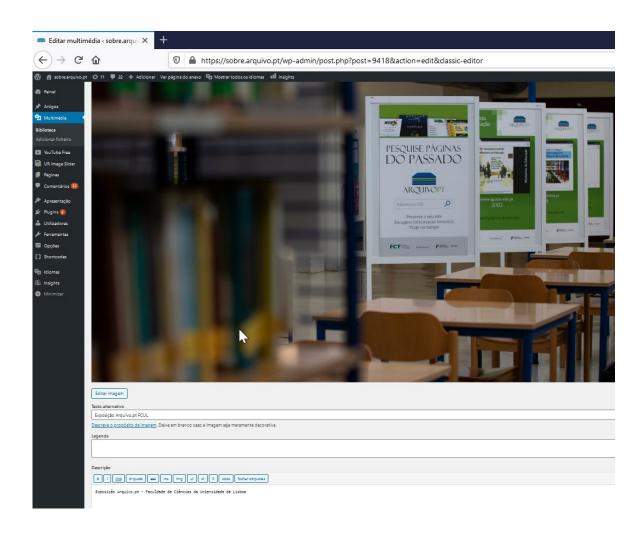
Exemplo de campos de descrição do Wordpress image title

```
<meta name="DC.Type" content="Text" />
<meta name="DC.Creator" content="Daniel Gomes" />
<meta name="DC.Date.Created" content="2009-08-21" />
<meta name="DC.Date.Modified" content="2009-11-10" />
```

Esta informação resume, **enriquece** ou complementa os conteúdos, produzindo assim um potencial incremento de informação.

Computadores conseguem utilizar esta informação.

Utilize metadados



7

Torne-se curador dos seus websites

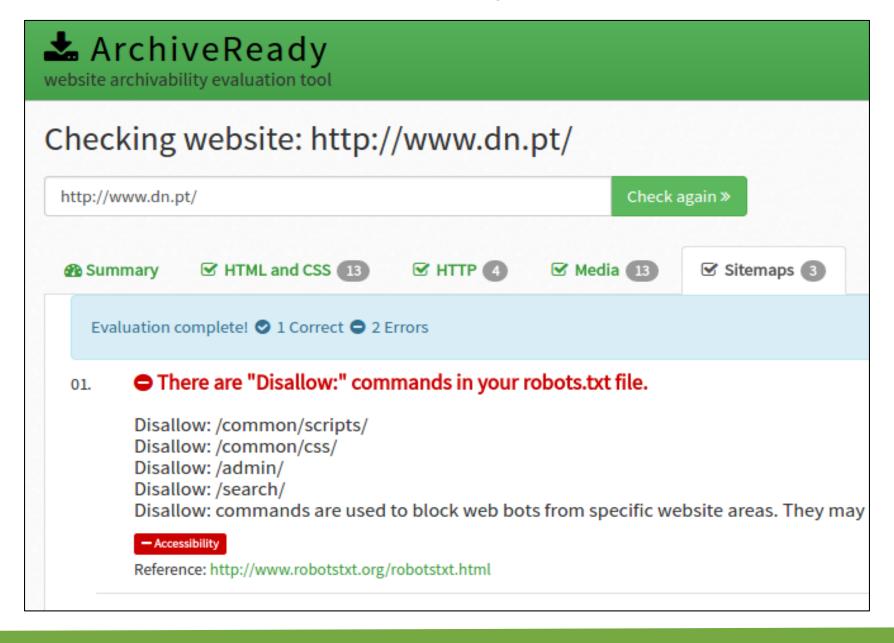
Use ferramentas para avaliar se uma página é preservável

Archive Ready

http://archiveready.com



http://archiveready.com



Recupere conteúdos perdidos

Soft404

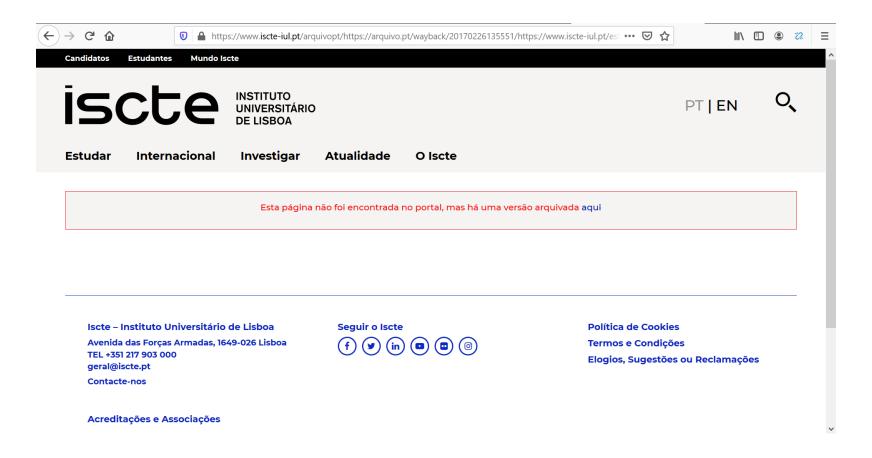


Páginas Web com ligações quebradas

Utilizadores seguem a ligação para uma página preservada no Arquivo.pt

Soft404

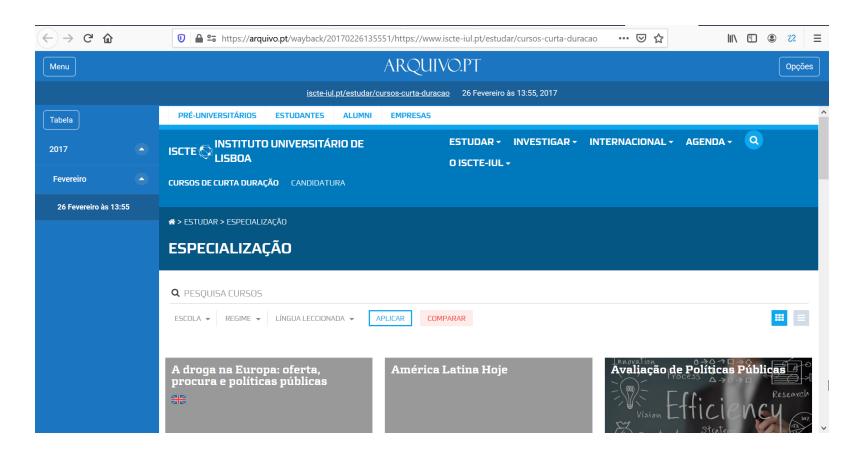
https://www.iscte-iul.pt/estudar/cursos-curta-duracao



https://github.com/arquivo/example-cdx-api

Soft404

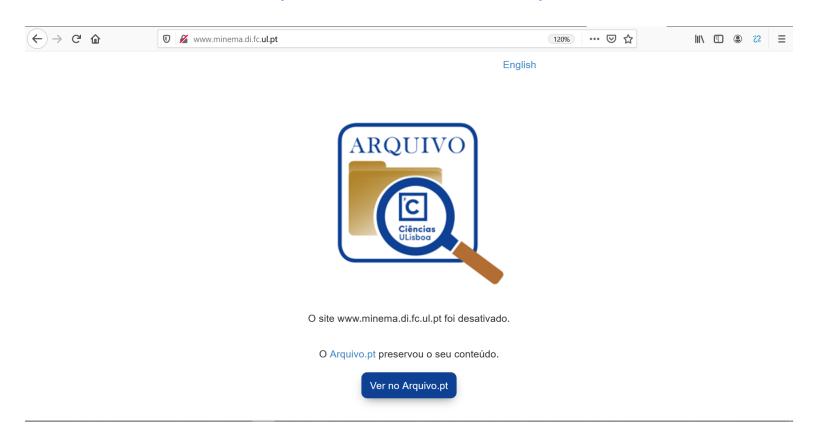
https://www.iscte-iul.pt/estudar/cursos-curta-duracao



https://github.com/arquivo/example-cdx-api

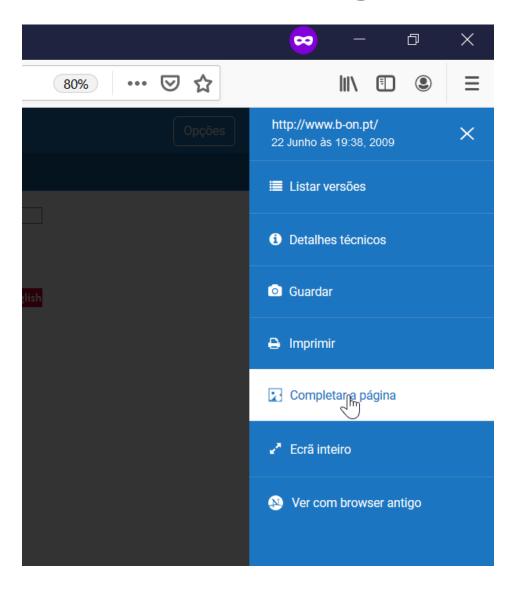
Memorial do Arquivo.pt

Exemplo, site do projeto Minema da Faculdade de Ciências da Universidade de Lisboa http://www.minema.di.fc.ul.pt/



https://arquivo.pt/memorial

Completar Página



WARC (Web ARChive)
ISO 28500:2017

Webrecorder

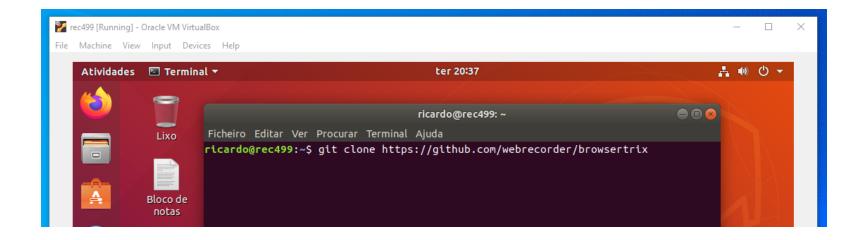
Selecionar

Capturar



Browsertrix

Instalar no próprio computador - Guia de instalação



Short link: https://tinyurl.com/instalar-browsertrix

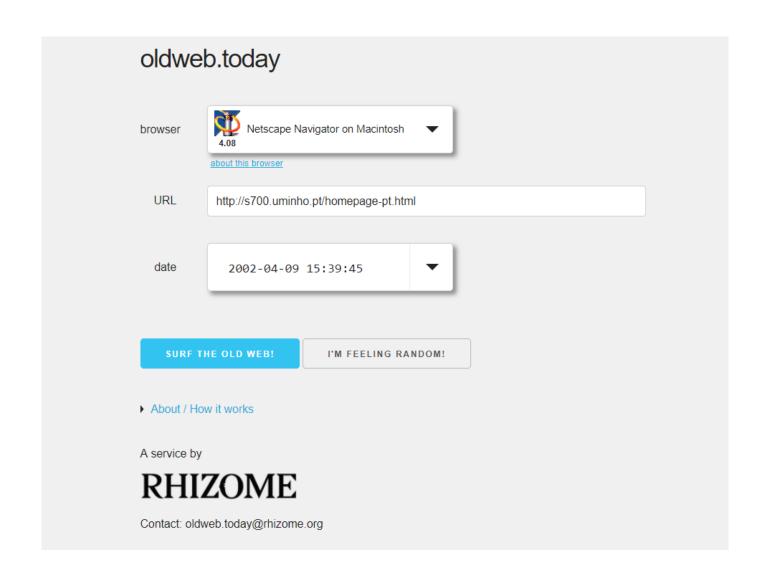
Torne-se curador dos seus websites

Recolhas locais feitas pela instituição ou pessoa proprietária do site:

- Recolhe páginas escolhidas
- Determina a periodicidade que acha adequada
- É objeto de uma verificação humana
- Utiliza ferramentas de recolha de alta qualidade ex. Webrecorder e Browsertrix
- Gera compromisso com o desenvolvedor do site
- Envolve outros intervenientes comunicação e imagem, arquivo e a própria gestão

Oldweb.today

Look and Feel do passado



http://oldweb.today



Time Left 08:27



about this browser

Current Page Archived On:

2002-02-04 12:23:05

Requested Date/Time:

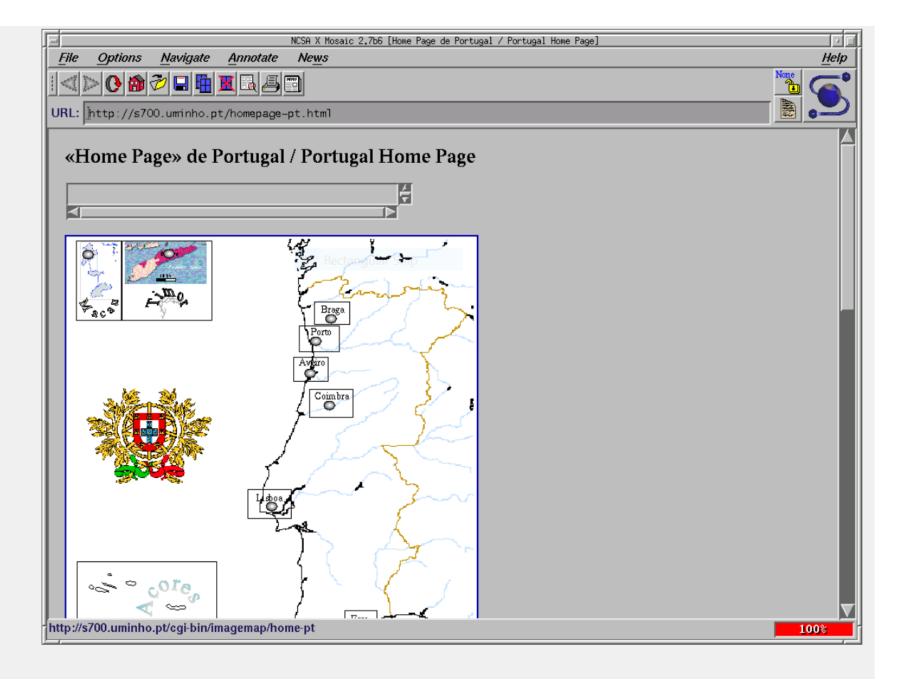
2002-04-09 15:32:31

Loaded 11 resources, spanning

2001-08-05 to 2013-08-02 10:27:41 16:22:08

from public web archives:

- Internet Archive
- Portuguese Web Archive



Conclusão

Recomendações criar e gerir websites preserváveis

- 1. Identifique corretamente a data de publicação
- 2. Use corretamente o protocolo de exclusão de robots
- 3. Use um endereço para cada conteúdo
- 4. Mantenha endereços ao longo do tempo
- 5. Utilize **formatos adequados** para preservação
- 6. Publique **metadados** para enriquecer os conteúdos
- 7. Torne-se curador dos seus websites

Recomendações:

arquivo.pt/recomenda

