

# Opportunities and challenges in collecting and studying national webs

Daniel Gomes

# Collecting the national Web of Portugal

Broad crawls: .PT + user suggestions

Daily crawls: 361 selected websites

Special crawls: events such as elections

High-quality crawls: on-demand

An attempt to archive the .EU  
domain

# Tried to perform a **collaborative** about R&D projects

**WE**ather **hAZARDs** for aeronautics

- Home
- The Project
- The Consortium
- The Advisory Board
- Deliverables
- Documentation
- Publications
- Events
- Related projects
- Contact

WEZARD Home15 Feb 2013

### Overview

The European WEZARD project (acronym standing for Weather Hazards for Aeronautics) aims at preparing the future research community in the area of air transport system robustness when it is faced with weather hazards. Its precise objectives are to provide:

(i) an interdisciplinary and cross-sector network comprising relevant experts; (ii) a state-of-the-art review of the on-going research actions; (iii) an analysis which will identify the shortcomings, areas for improvements and the type of activity needed to limit the effects of disruptive events; (iv) a set of recommendations and a roadmap validated by the main stakeholders of the aeronautics community.

The WEZARD consortium consists of 3 airframers, 2 engine manufacturers, 1 system supplier, 1 network of meteorological offices, 4 research centers, 1 provider of test facilities and 1 civil aviation authority over 2 years. An Advisory Board gathering a panel of international experts in relevant domains has been set up to provide advice on the vision, priorities and directions proposed by the project.

The project runs for 24 months from July 2011 until June 2013.



Funded by:



Collaborative list of Research and Development project websites				
File Edit View Insert Format Data Tools Add-ons Help Accessibility Last edit was made on 14 April 2016 by Dani...				
Collaborative list of sites of Research and Development projects				
	A	B	C	D
1	<b>Collaborative list of sites of Research and Development projects</b>			
	Websites of R&D projects provide valuable information but quickly disappear.			
	Arquivo.pt - the Portuguese Web Archive is making an experiment to preserve websites of R&D projects but <b>we need your help to start identifying them</b> .			
2	We need sites and lists of R&D projects in all scientific areas. It is very important that you also provide us the Acronym and Title of the project. The objective is to use this information to make the process of identifying R&D websites more automatic in the future.			
	Could you help by contributing to the list below?			
	Thank you. /Daniel Gomes			
3	<b>Project URL</b>	<b>Project Acronym</b>	<b>Project Title</b>	
4	<a href="http://4cproject.eu/">http://4cproject.eu/</a>			
5	<a href="http://alexandria-project.eu/">http://alexandria-project.eu/</a>			
6	<a href="http://amires.eu">http://amires.eu</a>			
7	<a href="http://aparsen.eu">http://aparsen.eu</a>			
8	<a href="http://atlas.fcsh.unl.pt/">http://atlas.fcsh.unl.pt/</a>			
9	<a href="http://axleproject.eu/">http://axleproject.eu/</a>			
10	<a href="http://base-adaptation.eu">http://base-adaptation.eu</a>			
11	<a href="http://biobankcloud.eu/">http://biobankcloud.eu/</a>			
12	<a href="http://blogforever.eu/">http://blogforever.eu/</a>			
13	<a href="http://bridge-project.eu/">http://bridge-project.eu/</a>			
14	<a href="http://byte-project.eu">http://byte-project.eu</a>			
15	<a href="http://clarin.eu/">http://clarin.eu/</a>			
16	<a href="http://coroado-project.eu/dissemination/">http://coroado-project.eu/dissemination/</a>			
17	<a href="http://dream2020.eu/">http://dream2020.eu/</a>			

*Google Sheet* to gather websites about R&D and development projects.

Over 25 000 project on the EU database.

# Archive the .EU domain: a “brute-force” attempt to preserve R&D websites

## 3 crawls of .EU domain

Time of crawl	# of Files collected	Data volume (TB)
21/11/2014 to 2014/12/16	129 793 987	5.8
2016/01/07 to 2016/01/26	61 863 684	3.1
02/06/2017 to 10/07/2017	105 823 552	11

Table 1. Comprehensive .eu crawls 2014-2017

Main problem: **web spam**

Searchable and accessible at:

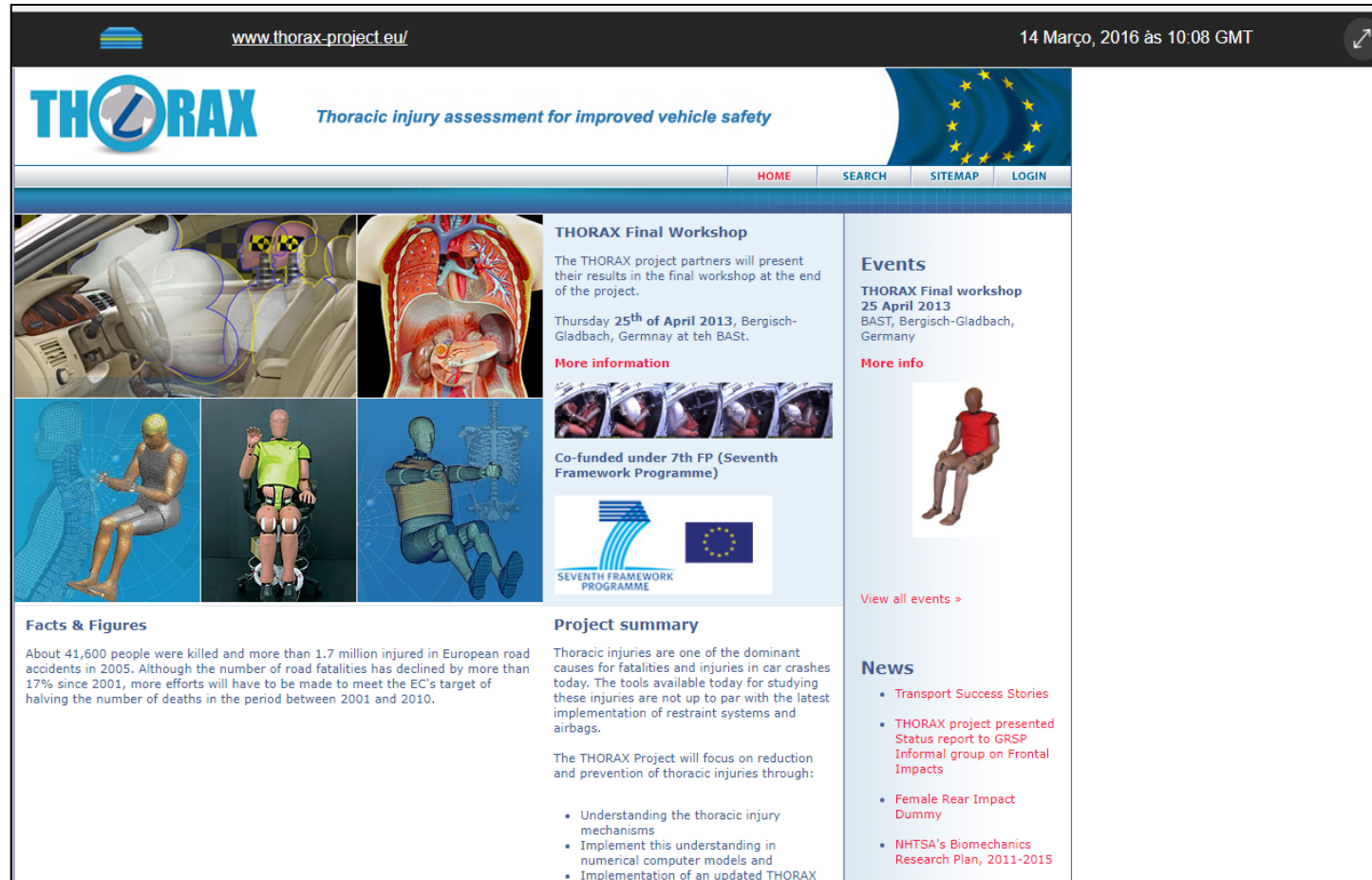
<https://arquivo.pt/resawdev>

Main problem: **web spam**



The screenshot shows the RESEARCH.EU search interface. At the top, the word "RESEARCH" is in a large, green, serif font, with ".EU" in a smaller, green, sans-serif font below it. Below this is a search bar with a light blue border and a blue "Search" button to its right. A small "x" icon is visible in the search bar. Below the search bar, the text "Search pages from the past" is displayed in a black, sans-serif font. A link "Meet the service" is visible below the text. A link "Advanced search" is visible to the right of the search bar.

# Automatic selection and preservation of websites related to R&D projects



[thorax-project.eu](http://thorax-project.eu), 2014

# Studying past webs

# Training courses on web preservation and research

## **New ways of searching the past**

Any Internet user

## **Publishing preservable information on the web**

Web authors

## **Automatic processing of information preserved from the Web**

Developers



[arquivo.pt/training](http://arquivo.pt/training)



# Investiga XXI (Research XXI)

## Communication Studies

Transformations of the Websites of Portuguese Newspapers



Short link to this page: [arquivo.pt/newspapers](https://arquivo.pt/newspapers)

## Information Science

FCSH on the Web: virtual exhibition



Short link to this page: [arquivo.pt/fcshontheweb](https://arquivo.pt/fcshontheweb)

## Social Sciences

Straight-Edge in the Lisbon metropolitan area



Short link to this page: [arquivo.pt/straightedgen](https://arquivo.pt/straightedgen)

All videos, presentations, reports at: [arquivo.pt/research](https://arquivo.pt/research)



Any subject

Arquivo.pt as main source  
of information

Submissions in Portuguese

[arquivo.pt/prizes](https://arquivo.pt/prizes)

1st place: 10 000 €

2nd place: 3 000 €

3rd place: 2 000 €