



NICOLAUS COPERNICUS  
UNIVERSITY  
IN TORUŃ  
Faculty of History

# **Studying the past Web in Poland – current state and perspectives**

Bartłomiej Konopa

NCU Toruń, State Archive in Bydgoszcz

06.06.2019



## Table of contents

- Need of Web archiving
- Web archiving projects
- State of research
- Researches on the past Web
- Survey results
- Perspectives

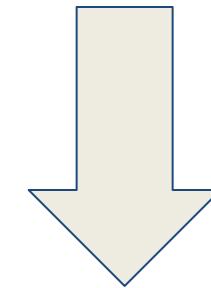
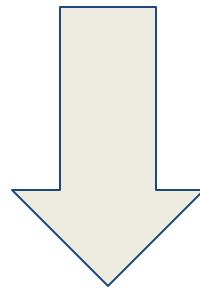


***Program for digitizing cultural heritage and collecting,  
storing and accessing digital objects in Poland 2009-2020***

Archiving polish Web as one of four task of  
expanding Polish digital resources



# Institutions that could be responsible for archiving the Polish Web



National Library  
of Poland

National Digital  
Archives



## Internet Archive by National Digital Archives

- launched in april 2009
- opened in march 2010
- last crawl in april 2011
- closed in end of 2015
- scope: 41 sites (governmental and State Archives websites)



## archiwuminternetu beta

Wprowadź adres www:  2009 ▾ Szukaj  
Wyszukiwanie zaawansowane



### Witamy w Archiwum Internetu!

Narodowe Archiwum Cyfrowe, aby zabezpieczyć informacje o historycznym znaczeniu dla państwa polskiego archiwizuje strony internetowe najważniejszych urzędów państwowych. W wersji beta Archiwum Internetu dostępne są strony poniższych instytucji.

- Naczelną Dyrekcję Archiwów Państwowych
- Archiwum Główne Akt Dawnych
- Archiwum Akt Nowych
- Narodowe Archiwum Cyfrowe
- Archiwum Państwowe w Białymostku
- Archiwum Państwowe w Bydgoszczy
- Archiwum Państwowe w Częstochowie
- Archiwum Państwowe w Elblągu z siedzibą w Malborku
- Archiwum Państwowe w Gdańsku
- Archiwum Państwowe w Gorzowie Wielkopolskim
- Archiwum Państwowe w Kaliszu
- Archiwum Państwowe w Katowicach
- Archiwum Państwowe w Kielcach
- Archiwum Państwowe w Koszalinie
- Archiwum Państwowe w Krakowie
- Archiwum Państwowe w Lesznie
- Archiwum Państwowe w Lublinie
- Archiwum Państwowe w Łodzi



Projekt | Ludzie | Prawa autorskie



## National Library of Poland

- took part in '09 European Election Web Harvesting Project
- only 16 websites captured
- unpublished



## Anonymous Internet Archive

- at the turn of 2012 and 2013
- websites from polish ccTLD found in Alexa rankings
- 14 thousand websites (only homepages)
- crawl every hour
- collected data is not available



### Archiwum Internetu

- [Dane](#)
- [O Archiwum](#)
- [Kontakt](#)

### Dane

- Obierają wszystkie adresy z końcówką .pl znajdujące się w rankingu miliona najczęściej odwiedzanych stron w Internecie wg [alexa.com](#). W tej chwili jest to nieco ponad 14 tys. adresów.
- Strony główne tych serwisów kopiowane są co godzinę i pakowane do pliku.
- W tej chwili dostępne są tylko surowe dane, tj. **co godzinę do poniższej listy dopisywana jest aktualna paczka**.

Oto one:

[2012-11-17 20:00:01](#)  
[2012-11-17 19:00:02](#)  
[2012-11-17 17:00:01](#)  
[2012-11-17 16:00:02](#)  
[2012-11-17 14:00:02](#)  
[2012-11-17 13:00:01](#)  
[2012-11-17 12:00:01](#)  
[2012-11-17 11:00:01](#)  
[2012-11-17 10:00:02](#)  
[2012-11-17 09:00:01](#)  
[2012-11-17 08:00:02](#)  
[2012-11-17 07:00:01](#)  
[2012-11-17 06:00:01](#)  
[2012-11-17 05:00:02](#)  
[2012-11-17 03:00:02](#)  
[2012-11-17 00:00:02](#)  
[2012-11-16 22:00:01](#)  
[2012-11-16 21:00:01](#)  
[2012-11-16 20:00:02](#)  
[2012-11-16 18:00:02](#)  
[2012-11-16 17:00:02](#)  
[2012-11-16 16:00:02](#)  
[2012-11-16 15:00:01](#)  
[2012-11-16 14:00:02](#)  
[2012-11-16 13:00:01](#)



## My Country by ePanstwo

- political related content gathered from Facebook and Twitter



MOJE PAŃSTWO SEARCH

Analizowany okres: Ostatnia doba Tydzien Miesiąc Rok

Maj 2019 Więcej ▾

3 czerwca 2018 r. 23:10 – dzisiaj 23:10

### Najbardziej angażujące tweety

Tweety, które uzyskały najwyższą liczbę retweetów, polubień i komentarzy.

Stanisław Janecki Komentator 28 września 2018 r. 10:18:53

Ciekawe, czy Wojciech Smarzowski zrobiłby film "Aktorzy". O alkoholikach, seksoalkoholikach, narkomanach, oszustach fin... <https://t.co/oxecRB57dN>

0 0 95

Sławek Neumann Polityk 28 września 2018 r. 10:18:43

Partia Panów jak mówi prezes. Ich prawo nie obowiązuje.

0 0 17

Konrad\_Piasecki Komentator 28 września 2018 r. 18:30:30

I tak środkiem trawnika? I akurat Statua Wolności w tle? I zero potu na czole?

0 0 0

Pozwól monitorować swoje tweety!

Dodaj konto

### Najbardziej angażujące profile

Profile, których tweety uzyskały największe liczby retweetów, polubień i komentarzy.

	Stanisław Janecki	95
	Prawo i Sprawiedliwość	38
	Sławek Neumann	17
	Konrad_Piasecki	15
	Kancelaria Prezydenta	14

11

My Country web service



01 października 2018

## Nowy serwis internetowy MNiSW

Od 1. października serwis internetowy MNiSW przechodzi do centralnego serwisu rządowego GOV.PL. Aktualne strony będą dostępne jako archiwum.

→

!! ▲ Ukryj rotator



# Polish scientific literature about Web archiving projects

- mostly about Internet Archive foundation
- about using Wayback Machine
- also introducing other national or grassroot initiatives



# Polish scientific literature about Web archiving

- Web archiving as a process
- methods and tools
- legal issues
- grassroot archiving
- cooperation with users
- information management



# Polish scientific literature about archived Web as sources for researches

- Web archeology as a way to study websites
- important source for contemporary history and humanistics
- going to replace old sources
- requires new skills



## webArch

Pracownia archiwizacji Webu LaCH UW

O pracowni Blog About @DHLab\_UW

Szukaj ...

Szukaj

O PRACOWNI

Zespół

ABOUT

BAZA WIEDZY

Archiwizacja Webu

Periodyzacja archiwizacji Webu

Krótką historią

Bibliografia (pol)

Do pobrania

Seminarium (2017)

Tworzenie przypisów

Jak zadbać o stabilne linki?

Sprawdź swoją stronę



### webArch LaCH UW

Celem pracowni jest gromadzenie wiedzy i rozwijanie kompetencji w zakresie profesjonalnej archiwizacji Webu i zaawansowanej pracy z archiwami WWW.

### Archiwizacja Webu

Celem pracowni jest gromadzenie wiedzy i rozwijanie kompetencji w zakresie profesjonalnej archiwizacji Webu i zaawansowanej pracy z archiwami WWW.

### Warsztaty

Zachęcamy do udziału lub współorganizacji warsztatów z zagadnień archiwistyki Webu.



M. Roszkowski, B. Włodarczyk

*Web Citations in Polish Library and Information Science Journals: the Analysis of URLs' Validity (2016)*

- 4 scientific journals analized
- 4593 unical URLs from 670 articles - approx. 30% inactive
- 53% of inactive was found in the Internet Archive resources



M. Król, D. Zdonek and J. Gorzelny

*Information value of Internet domain (2017)*

- searching for methods to evaluate reliability
- comparison of WHOIS search engine and Wayback Machine



## B. Michalska-Bednarek

### *History of websites of NCU Library (2016)*

- based on official reports, interviews and archived Web

## A. Surdyk

### *The activities of academic associations on the internet based on the example of PTBG. A summary of seven years of the existence of the association on the net (2017)*



## M. Król

### *From the archives of the Internet: changes in the mode of presentation of agrotouristic offer (2017)*

- analysis of websites of agritourism centers
- author observed changes in technical issues and designing style



## *The use of archived Web in scientific research - survey results*

- survey was addressed to scholars and PhD students
- division into users and nonusers
- 52 responses were collected



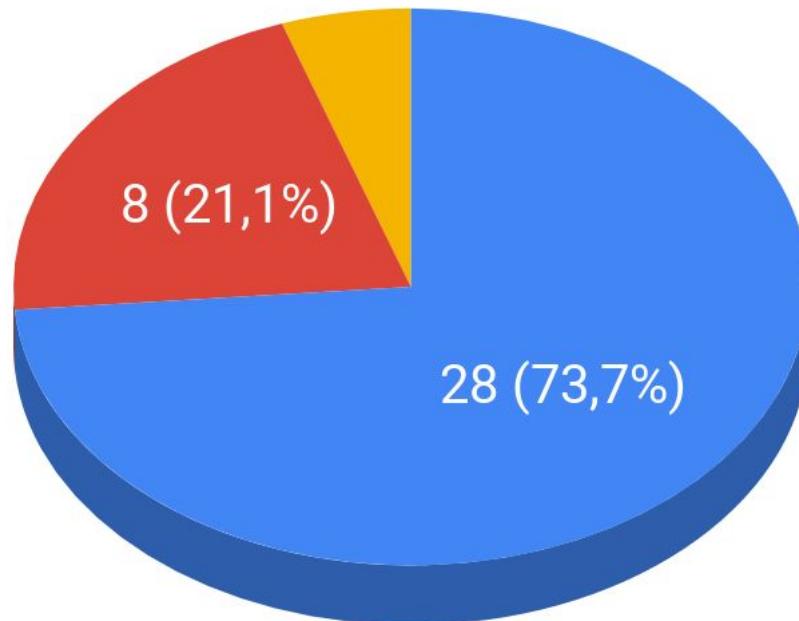
## ***The use of archived Web in scientific research - users***

- 11 respondents were using archived Web among others from social communication sciences, IT, economics, archeology and fine arts
- mostly using Internet Archive (project Minerva and Google Cache also mentioned)
- commonly used to obtain inaccessible materials
- average rating 3.9 (1-5 scale)



## *The use of archived Web in scientific research - nonusers*

Why you did not use archived Web in your researches

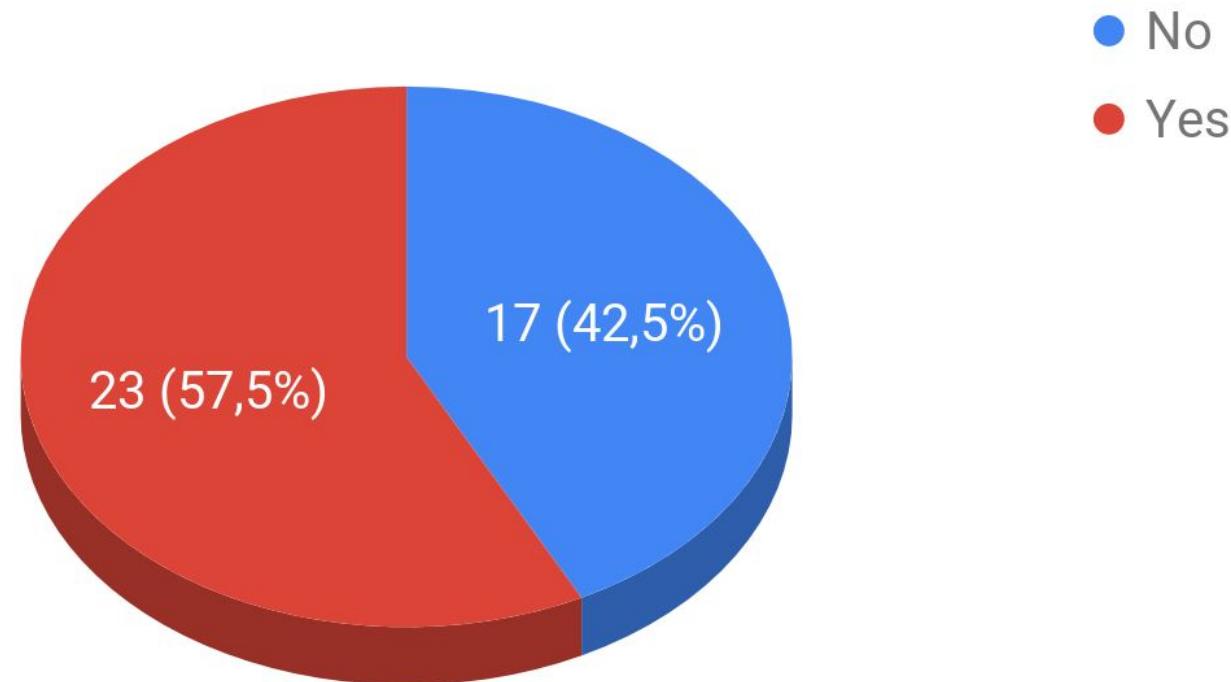


- I was not aware of them
- There was no usage for them
- I did not know how to use them



## *The use of archived Web in scientific research - nonusers*

Do you see the need for archiving Web for researches





## ***The use of archived Web in scientific research - user and nonusers***

Which scientific disciplines could use the archived Web?

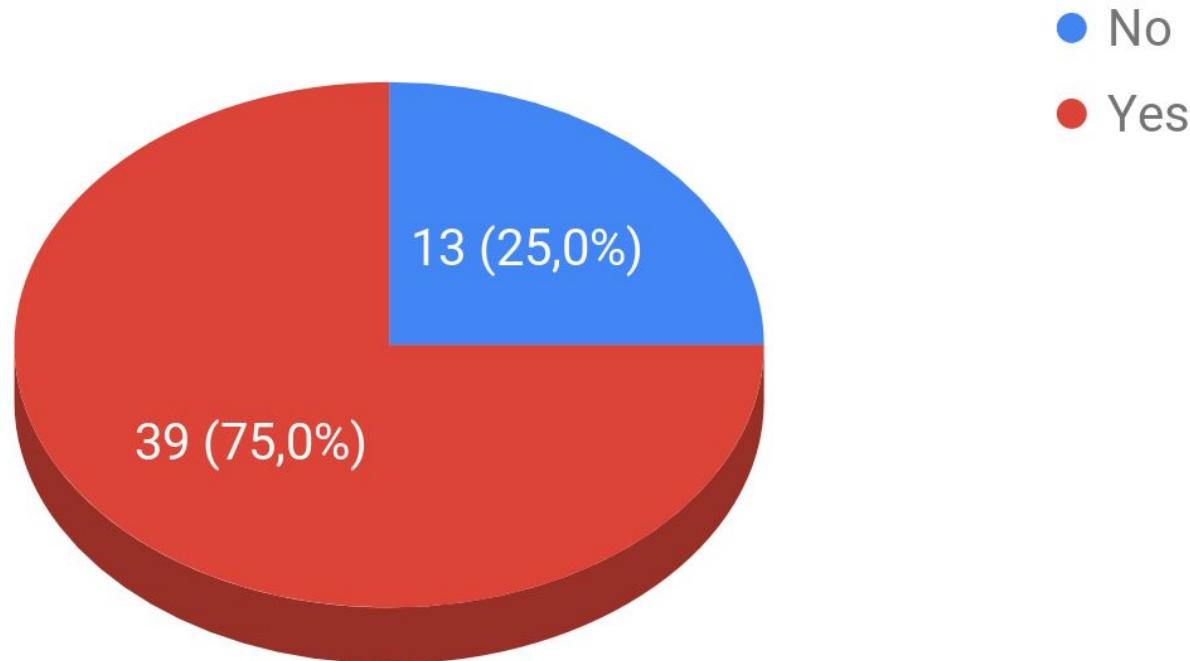
(most commonly responses):

- history (22)
- social communication and media sciences (13)
- sociological sciences (9)
- cultural and religious sciences (8)
- legal sciences (6)
- economics (6)
- political sciences (5)
- linguistics (5)
- computer sciences (5)



## *The use of archived Web in scientific research - user and nonusers*

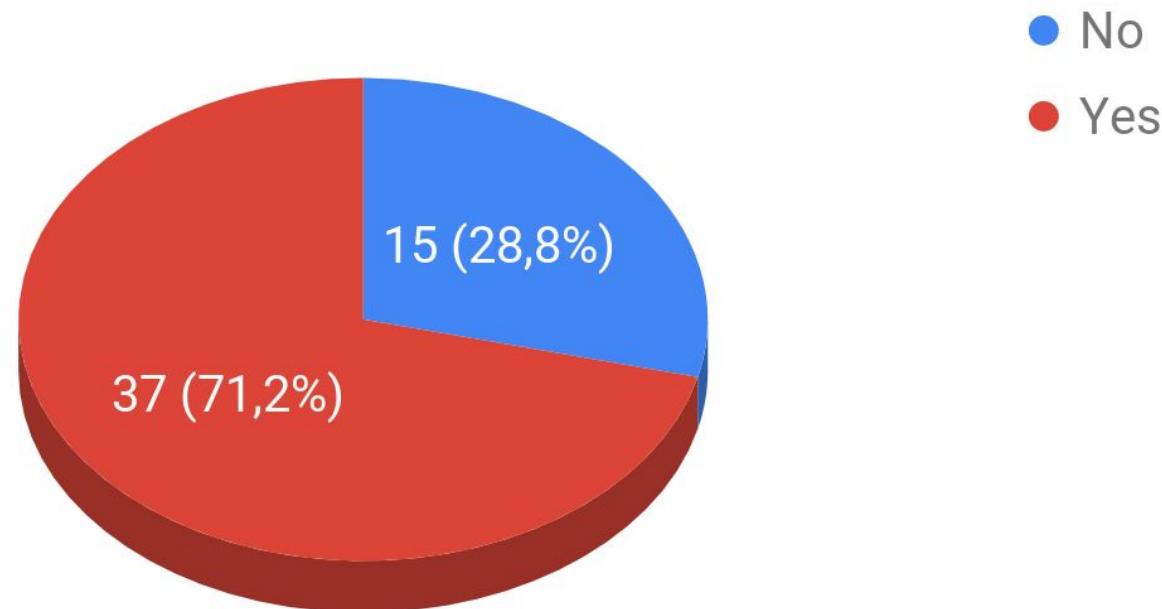
Should there be a Polish Web archive?





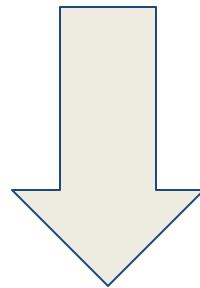
## ***The use of archived Web in scientific research - user and nonusers***

Does the issue of archived Web require popularization (eg in the form of trainings, workshops, etc.)?

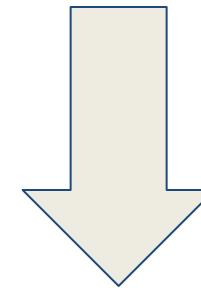




# Main aims for the future of studying past Web in Poland



establishing  
Polish Web  
archive



popularization  
and educating



# Thank you

Bartłomiej Konopa  
[bartlomiejkonopa@gmail.com](mailto:bartlomiejkonopa@gmail.com)