

Why We Need Private Web Archives

Hussam Hallak
WSDL Research Group
Ph.D. student

Types of web pages

- **Publicly Archivable:**
No paywall (mainly)
- **Not Publicly Archivable:**
SERPs
Maps
Social Media: Facebook, Linkedin, etc
Commerce: Amazon
Email sites: Yahoo, Gmail, Hotmail, etc
Bank statements

Dr. Steven J. Zeil



Associate Professor, Old Dominion University

Ph.D., 1981, Ohio State University

[vita](#)

Contact Information

- Office: E&CS 3208
Phone: (757) 683-4928
E-mail: zeil@cs.odu.edu
- [Office hours and Appointments](#)

Dr. Steven J. Zeil



Associate Professor, Old Dominion University

Ph.D., 1981, Ohio State University

[vita](#)

Contact Information

- Office: E&CS 3208
Phone: (757) 683-4928
E-mail: zeil@cs.odu.edu
- [Office hours and Appointments](#)

電子郵件或電話

密碼

忘記帳號？

登入

是否要加入 Facebook？

註冊

Problem

- Estimating the amount of archivable web traffic

Few websites get most of the traffic

- Top 100 websites receive %49 of the entire traffic on the web.
- Fact is consistent with a Pareto distribution.

Web Traffic Measures

- **Total visits** is the total number of non-unique visits from last month.
- **Unique visits** is the number of unique visits from last month.
- **Pages/visit** is the average number of visited pages per user's visit.

Approach

- Collect all three web traffic measures for the top 100 websites (Alexa Ranking) from 2 different web traffic analysis services (SimilarWeb, Semrush). (Feb 23rd 2018)
- Classify the top 100 websites:
Publicly Archivable or **Not Publicly Archivable**
- Calculate the percentage of publicly archivable web traffic using six different formulas that use some or all traffic measures and their combinations in different ways

65% of Web Traffic is NOT publicly archivable

- Using all measures, 65.30% of the traffic of the top 100 sites is not archivable by public web archives
- What about the remaining 1.8 billion live websites?

Best case: If all are publicly archivable,
31% of web traffic is NOT publicly archivable

Worst Case: If all are Not publicly archivable,
83% of web traffic is NOT publicly archivable

Most Likely: The bottom half is similar to the top half
65% of web traffic is NOT publicly archivable

Take Away Message

- Personalized web pages are on the rise.
Therefore, personal web archiving is crucial.
- There will always be pages that public archives cannot capture no matter how good they get.
- We need personal web archiving tools
- Available tools:
Web Recorder, Warcreate, WAIL