

Reproducible Web Corpora: Interactive Archiving with Automatic Quality Assessment

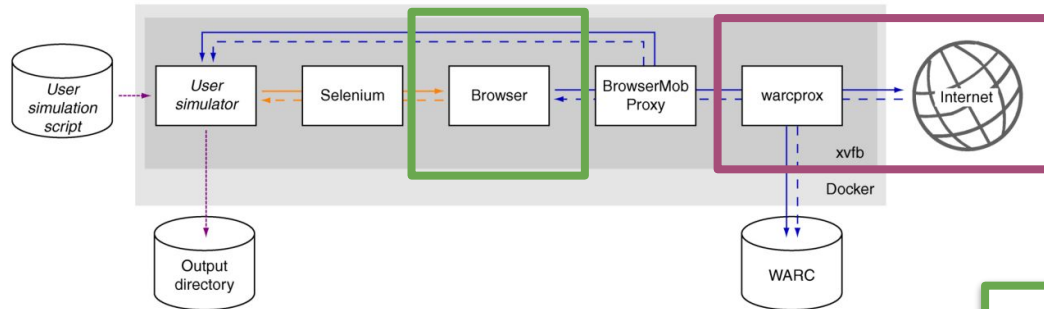
<https://dl.acm.org/doi/10.1145/3239574>

JOHANNES KIESEL
FLORIAN KNEIST
MILAD ALSHOMARY
BENNO STEIN
MATTHIAS HAGEN
MARTIN POTTHAST

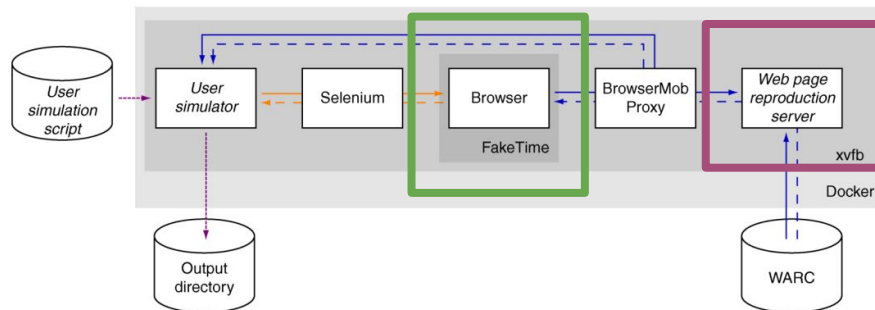
Bauhaus-Universität Weimar
Ulm University
Paderborn University
Bauhaus-Universität Weimar
Martin-Luther-Universität Halle-Wittenberg
Leipzig University

Presented by David Calano
CS 895 - Web Archiving Forensics
Old Dominion University
November 7th, 2022

Webis Web Archiver



(a)



(b)

Key:
 - HTTP request (solid blue arrow)
 - HTTP response (dashed blue arrow)
 - Browser/DOM control (solid orange arrow)
 - Browser/DOM status (dashed orange arrow)
 - File read/write (dashed purple arrow)

Mostly identical crawl and reproduction process

<https://github.com/webis-de/webis-web-archiver>

<https://github.com/webis-de/scriptor>

Figure 1

Breakdown of Dataset

Category	Sites	Pages	Most frequent site sub-categories	Most frequent page sub-categories
Adult	40	129	Computers (25), Image_Galleries (8), Society (2)	Computers (94), Image_Society (3)
Arts	187	471	Television (59), Movies (19), Music (18)	Television (152), Movie (50)
Business	229	489	Financial_Services (35), News_and_Media (24), Arts_and_Entertainment (22)	Financial_Services (116), News_and_Media (84), Arts_and_Entertainment (821), Software (143), Companies (60)
Computers	444	1502	Internet (184), Software (143), Companies (60)	Internet (821), Software (261)
Games	58	150	Video_Games (48), Gambling (6), Board_Games (1)	Video_Games (139), Gambling (2)
Health	28	36	Medicine (15), Conditions_and_Diseases (7), Nursing (5)	Medicine (21), Nursing (5)
Home	42	88	Consumer_Information (19), Cooking (9), Personal_Finance (7)	Consumer_Information (17), Personal_Finance (15)
Kids_and_Teens	55	99	School_Time (21), Games (13), Computers (7)	School_Time (43), Games (20), Computers (10)
News	149	303	Newspapers (77), Breaking_News (13), Magazines_and_E-zines (8)	Newspapers (110), Headlines (24)
Recreation	70	106	Travel (32), Autos (8), Humor (5)	Travel (63), Autos (10)
Reference	164	275	Education (105), Dictionaries (25), Libraries (21)	Education (128), Ask_a_Question (47)
Regional	605	1182	North_America (395), Europe (112), Asia (78)	North_America (829), Europe (144)
Science	74	116	Technology (19), Social_Sciences (16), Biology (12)	Technology (34), Social_Sciences (27), Biology (20)
Shopping	182	413	General_Merchandise (30), Clothing (26), Entertainment (16)	Entertainment (108), General_Merchandise (105), Auctions (51)
Society	100	123	Issues (20), Religion_and_Spirituality (15), Government (12)	Issues (27), Religion_and_Spirituality (18), Government (16)
Sports	73	123	Resources (18), Soccer (9), Baseball (6)	Resources (48), Soccer (16), Cricket (10)
World	1453	3437	Chinese_Simplified_CN (305), Russian (245), Deutsch (213)	Chinese_Simplified_CN (866), Russian (731), Deutsch (593)

Table 1

<https://alexa.com/topsites>
<https://webis.de/data/webis-web-archive-17.html>

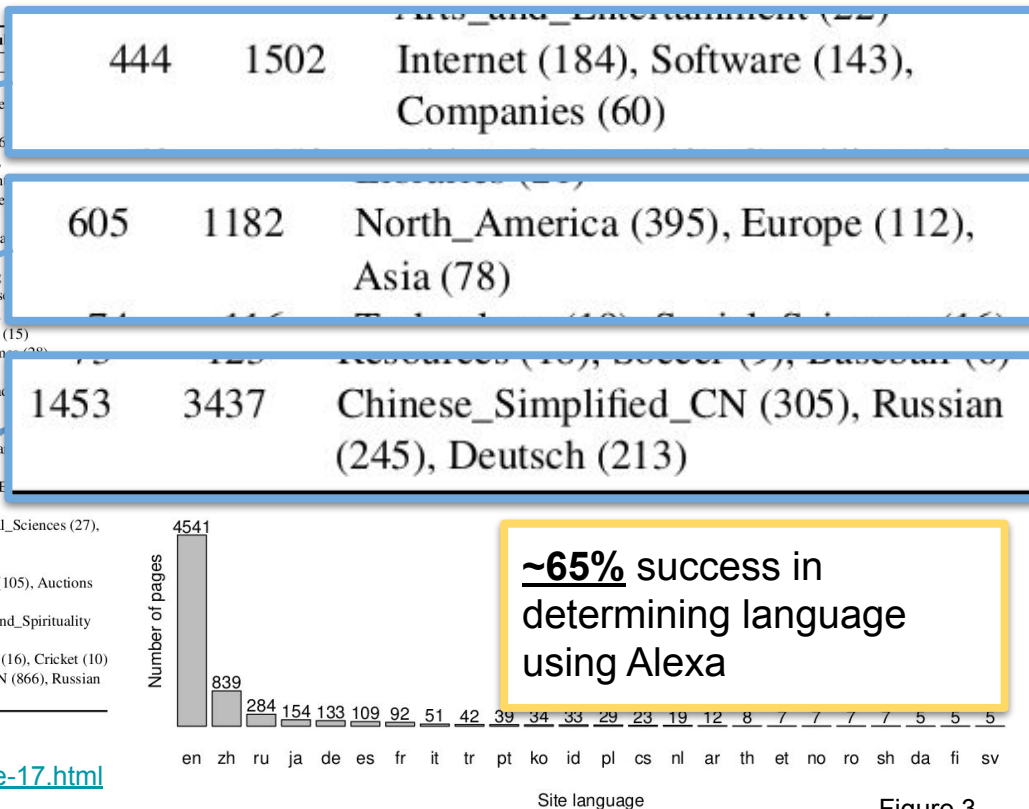


Figure 3

Web Capture vs Playback

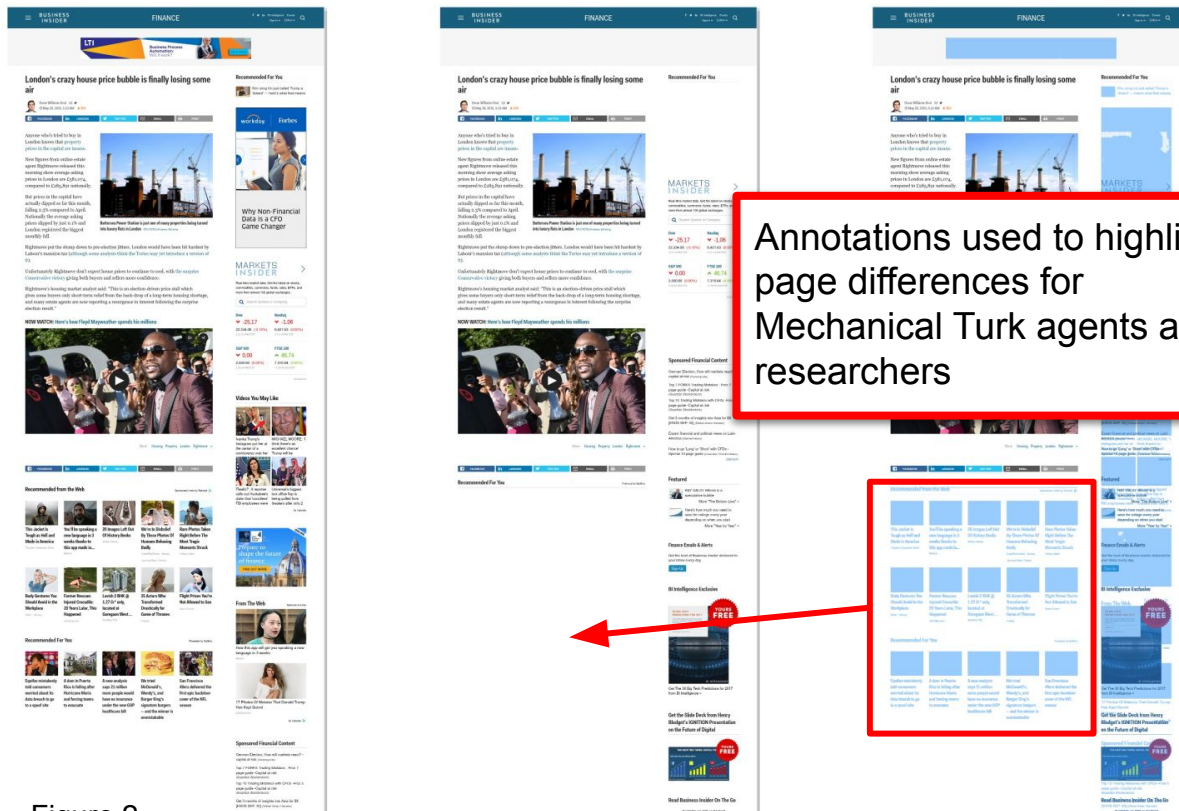


Figure 2

Mechanical Turk Used to Test Perceived Differences

Task 3

Left page

Right page

☒ Highlight differences (ignore the blue highlight for answering the question)

Imagine someone wants to visit the left page, but gets the right page. How much would this difference affect the visitor?
Examples for each score:

- Score 1 (not affecting): Parts of the page are just moved up or down a bit.
- Score 2 (small effect on a few visitors): Social media buttons, ads, or unimportant images or text are missing.
- Score 3 (small effect on many or all visitors): Comments on the main content are missing.
- Score 4 (affecting, but page can still be used): Striking difference in colors, background, or layout.
- Score 5 (unusable page, important/main content is missing and/or visitors can't use the right page due to the differences).

Difference will not affect visitor ☐ 1 ☐ 2 ☐ 3 ☐ 4 ☒ 5 Difference makes page useless for visitor

Comments for this task (optional):

1 2 3 4 5 Next task

Complete all tasks

1 for "best"

5 for "worst"

Figure 4

<https://www.mturk.com/>

1: Page Differences Have Negligible Effects

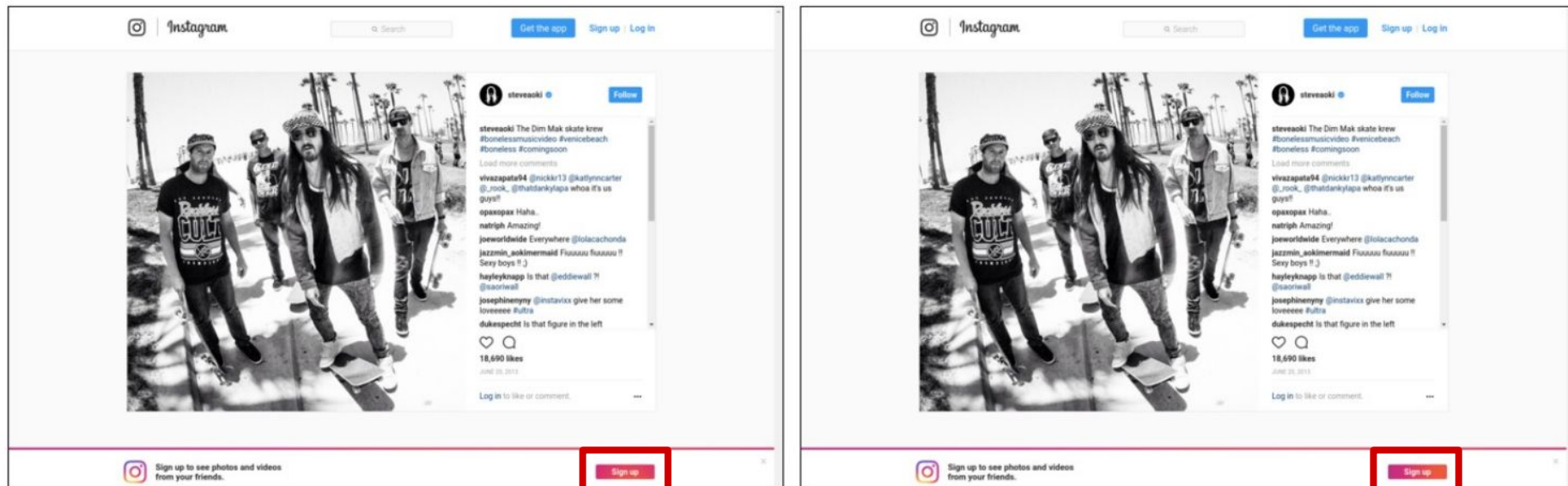


Figure 5

Slight difference in colors
of animated button

2: Small Issues Affecting Small Numbers of People

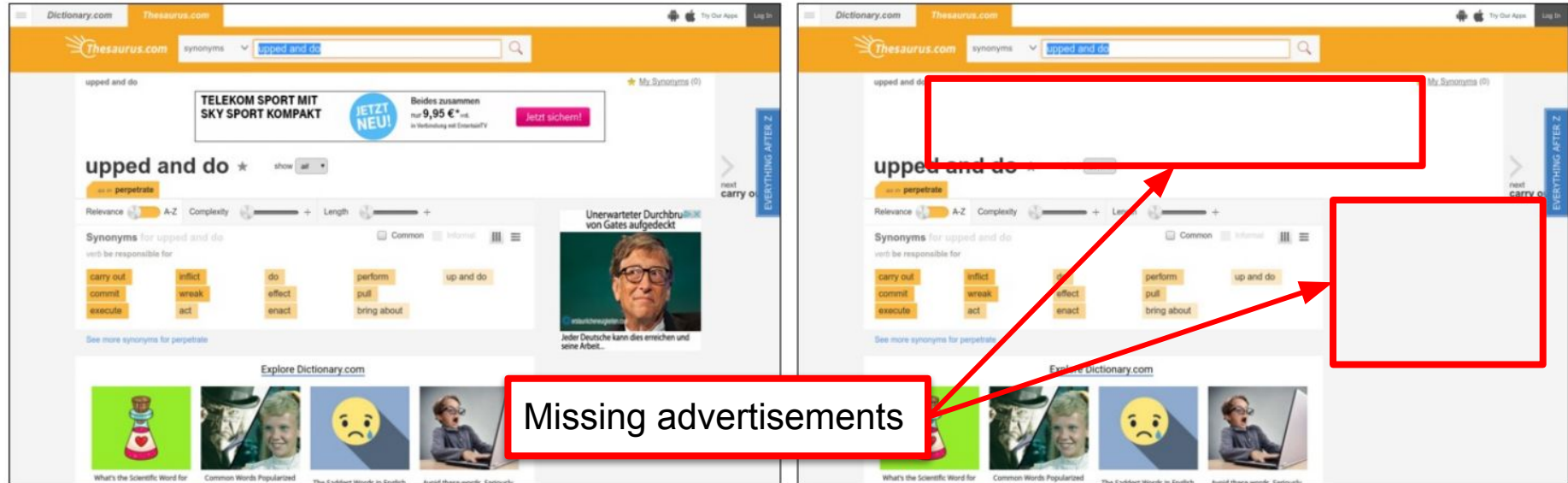


Figure 5

3: Page Has Small But Pervasive Differences

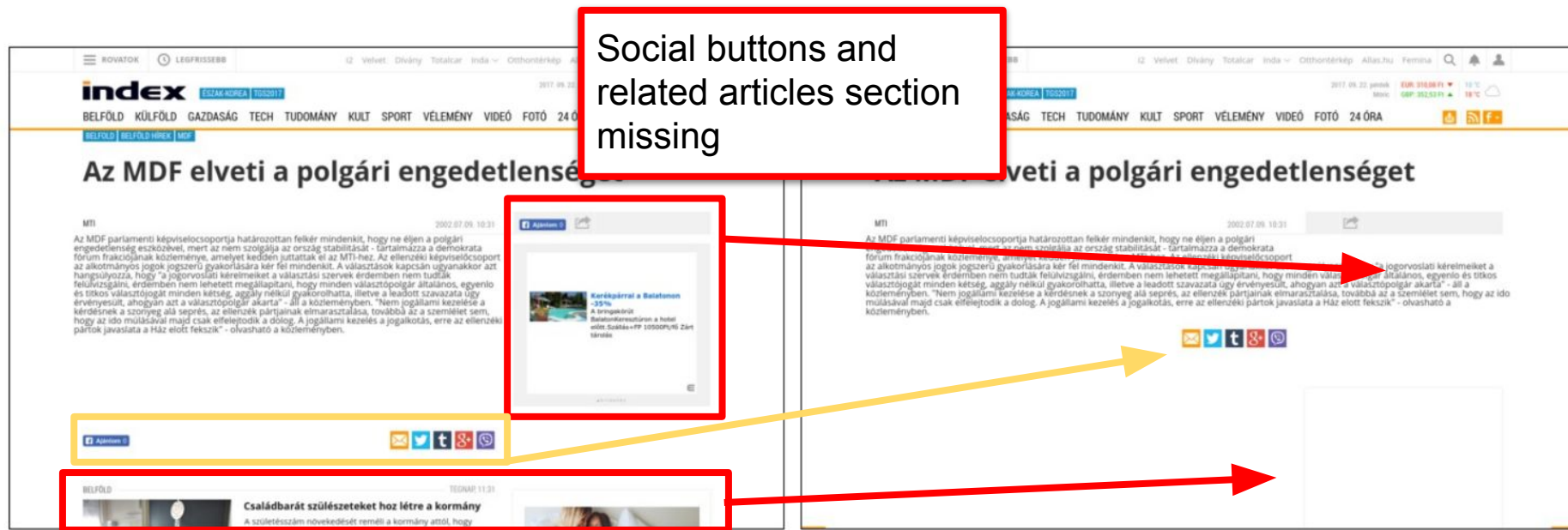


Figure 5

4: Page Is Usable Despite Significant Differences

Missing content not loaded
at page initialization

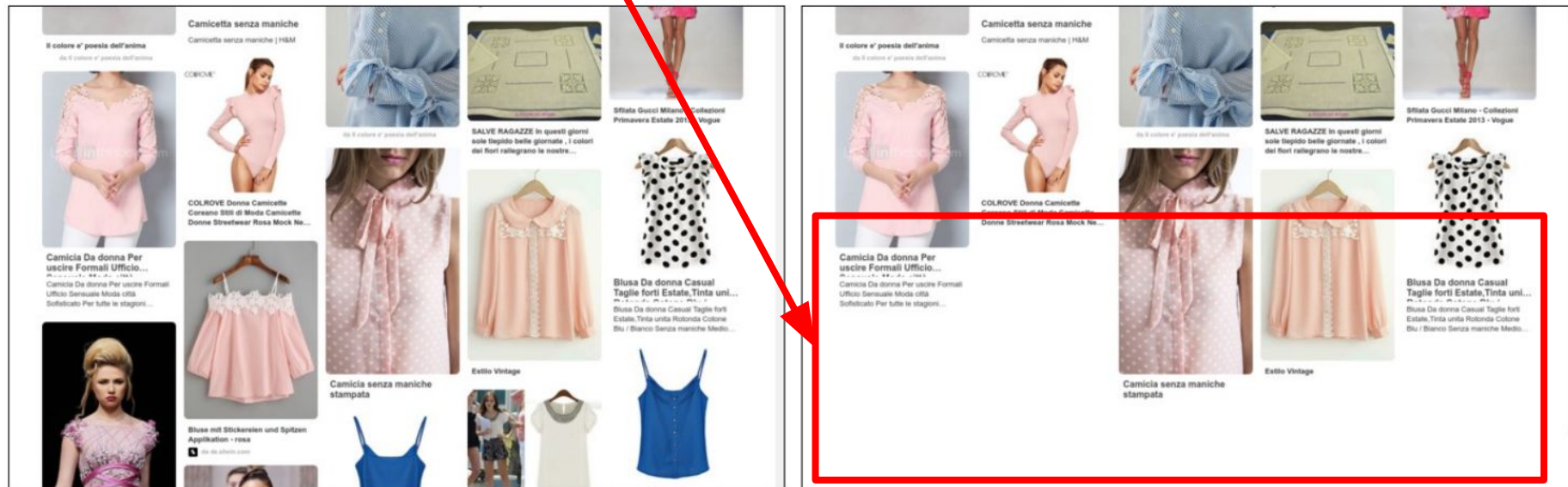


Figure 5

5: Page Is Effectively Unusable

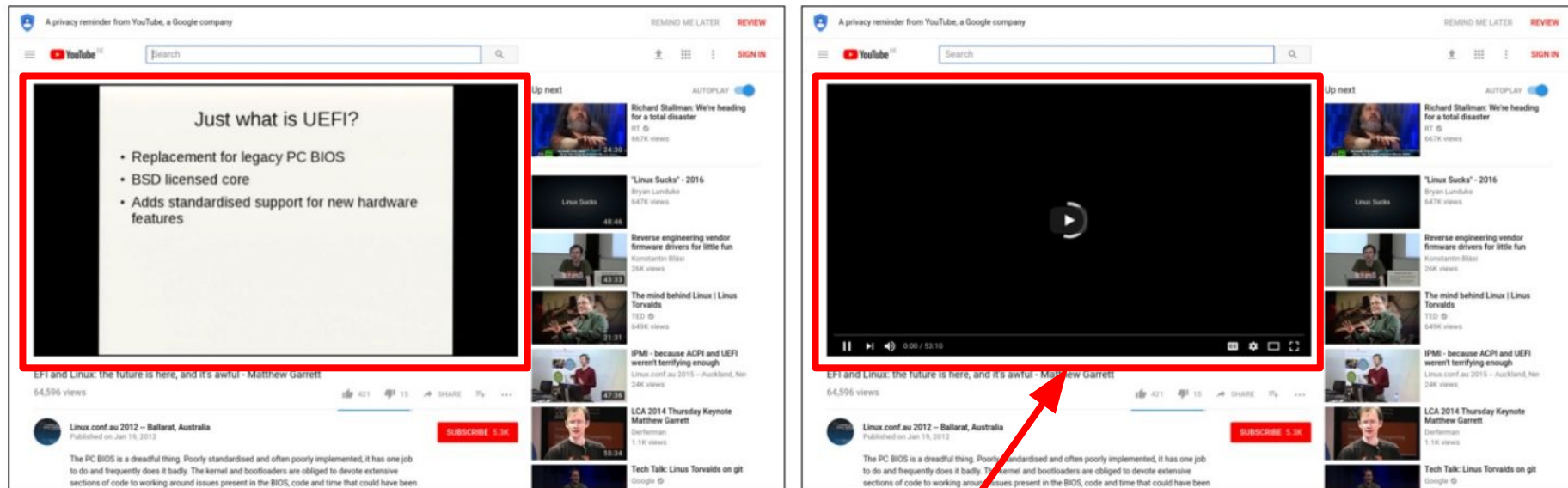


Figure 5

Primary content of page is missing or non-functional

Crowdsourced Score Breakdown

Reproductions	Score					Σ
	1	2	3	4	5	
Annotated	1942 (30.6%)	3307 (52.1%)	422 (6.7%)	318 (5.0%)	359 (5.7%)	6348
All	5594 (55.9%)	3307 (33.1%)	422 (4.2%)	318 (3.2%)	359 (3.6%)	10000

Table 2

Unannotated pages were reconfirmed by Mechanical Turk agents

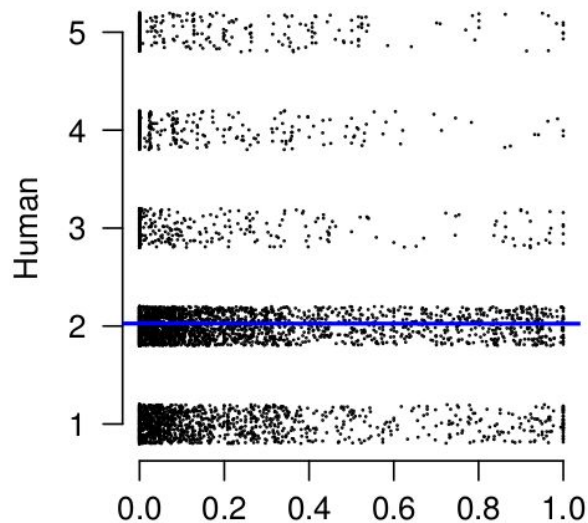
Algorithm Comparison

https://en.wikipedia.org/wiki/Pearson_correlation_coefficient

https://en.wikipedia.org/wiki/Root-mean-square_deviation

<http://justinfbunelle.com/pubs/jcdl-2014-brunelle-damage.pdf>

<https://ws-dl.blogspot.com/2018/09/2018-09-03-lets-compare-memento-damage.html>



Brunelle15

$r = 0.0$

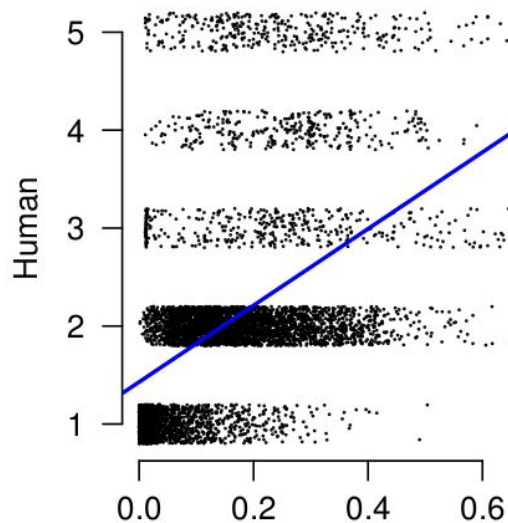
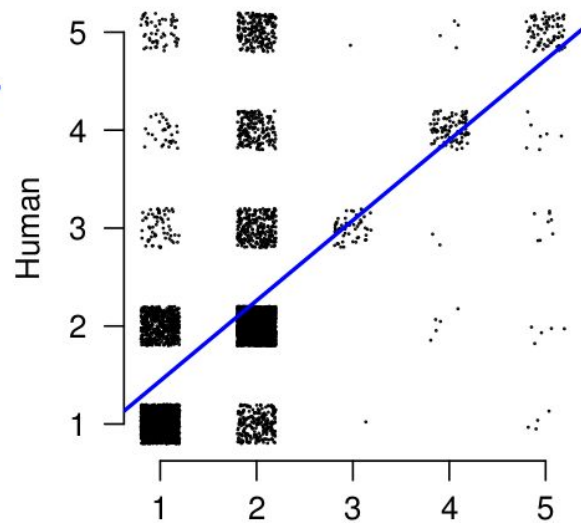


Figure 6

RMSE

$r = 0.48$



Neural network

$r = 0.57$

Pearson Correlation Coefficient measures linear correlation between algorithmic results

Neural Network Accuracy

Min. Quality	Acc.	Prec.	Rec.	F ₁
1	90.8%	94.1%	84.4%	89.0%
2	91.4%	94.4%	22.9%	36.9%
3	94.8%	88.1%	27.3%	41.7%
4	97.0%	76.0%	22.0%	34.1%

Table 3

Useful for automatically determining pages if were not faithfully reproduced

Most important metric for debugging

Calculations exclude trivial cases with no difference to screenshots

Comparison of Faithful Reproducibility for Replay Engines

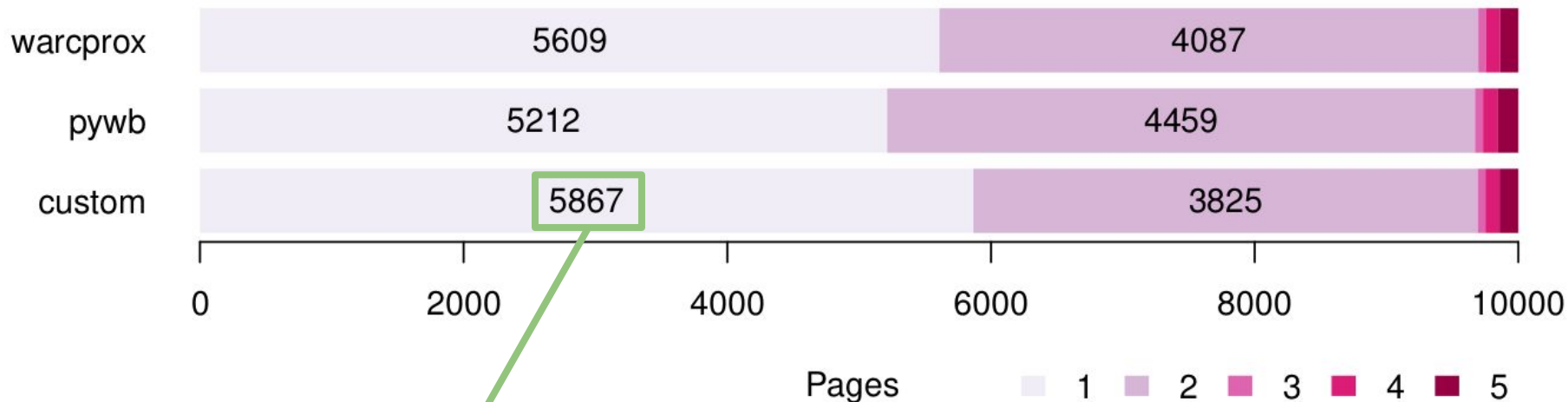


Figure 7

2.58% improvement over warcprox
6.55% improvement over pywb

<https://github.com/internetarchive/warcprox>

<https://github.com/webrecorder/pywb>

Conclusions

- Sample data set could be larger to improve accuracy and edge cases
- Using screenshot comparisons, neural networks can be trained to rapidly assess reproduced web page quality
- Move to automated calculation approaches over heuristics
- Check out Justin Brunelle's WS-DL blog post for more details:
<https://ws-dl.blogspot.com/2018/09/2018-09-03-lets-compare-memento-damage.html>