

MOTIVATION

- Memento TimeMaps [2] list identifiers for archival web captures (URI-Ms).
- Some URI-Ms redirect when dereferenced.
- Ergo, obtaining an accurate count of non-forwarding URI-Ms is not possible by solely using the contents of a TimeMap without dereferencing.
- Some archives have various endpoints beyond Memento.
- Endpoints' counts may differ depending on which is consulted.
- Therefore a standard-based mechanism must be used to obtain an accurate count.

CANONICALIZATION

URI canonicalization associates differently formatted URIs.
For example, <http://example.com> might be associated with:

<http://www.example.com> <https://www.example.com>
<http://example.com/> <http://example.com/index.html>
<http://example.com/#articles> <//example.com>

DIRECT-TO-INDIRECT METRIC

$$M = \text{URI-M} \subset \text{TimeMap}, \text{ if } \text{"memento"} \subset \text{valuesOf}(\text{rel}) \quad (1)$$

$$|TM|_{rel} = \sum_{m=1}^{len(M)} 1 \quad (2)$$

$$|TM|_D = \sum_{m=1}^{len(M)} \begin{cases} 0 & 300 \geq \text{httpStatus}(m) < 400, \\ 1 & \text{otherwise.} \end{cases} \quad (3)$$

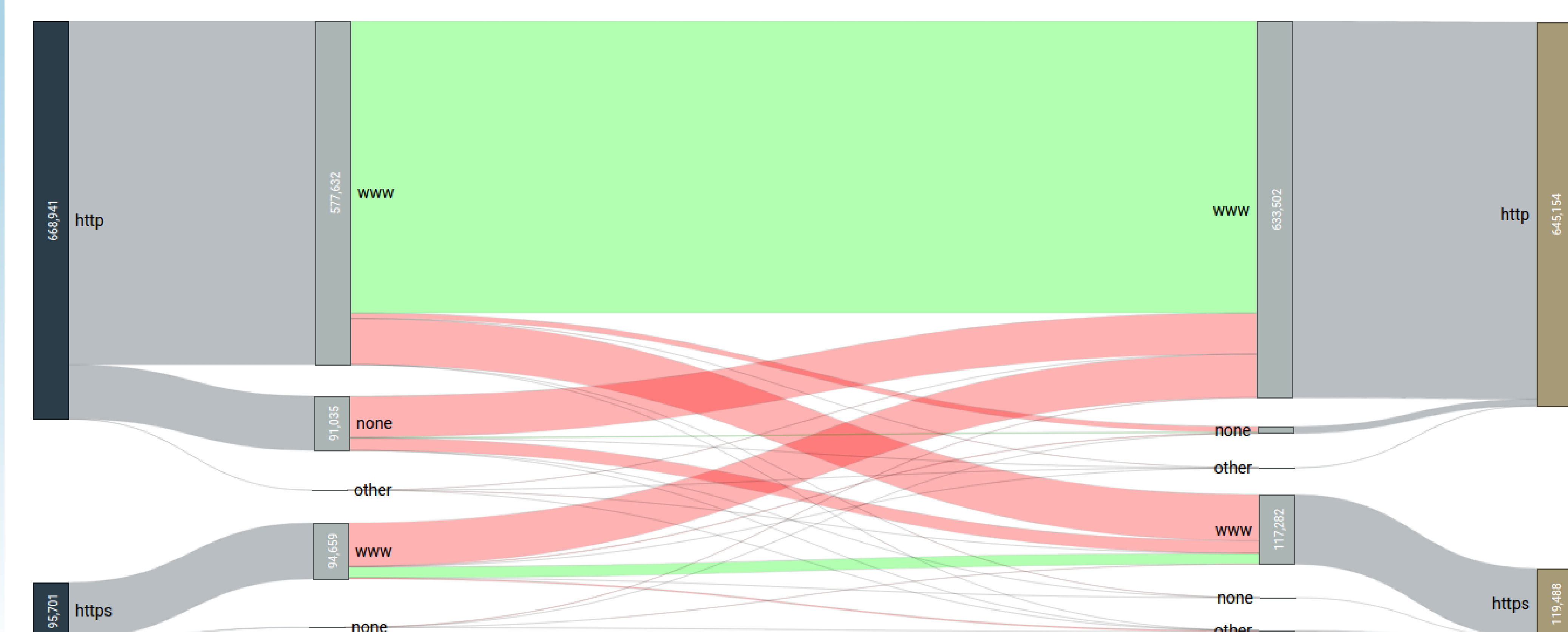
$$|TM|_I = |TM|_{rel} - |TM|_D \quad (4)$$

$$DI = \begin{cases} \frac{|TM|_D}{|TM|_I} & |TM|_I > 0, \\ \infty & \text{otherwise.} \end{cases} \quad (5)$$

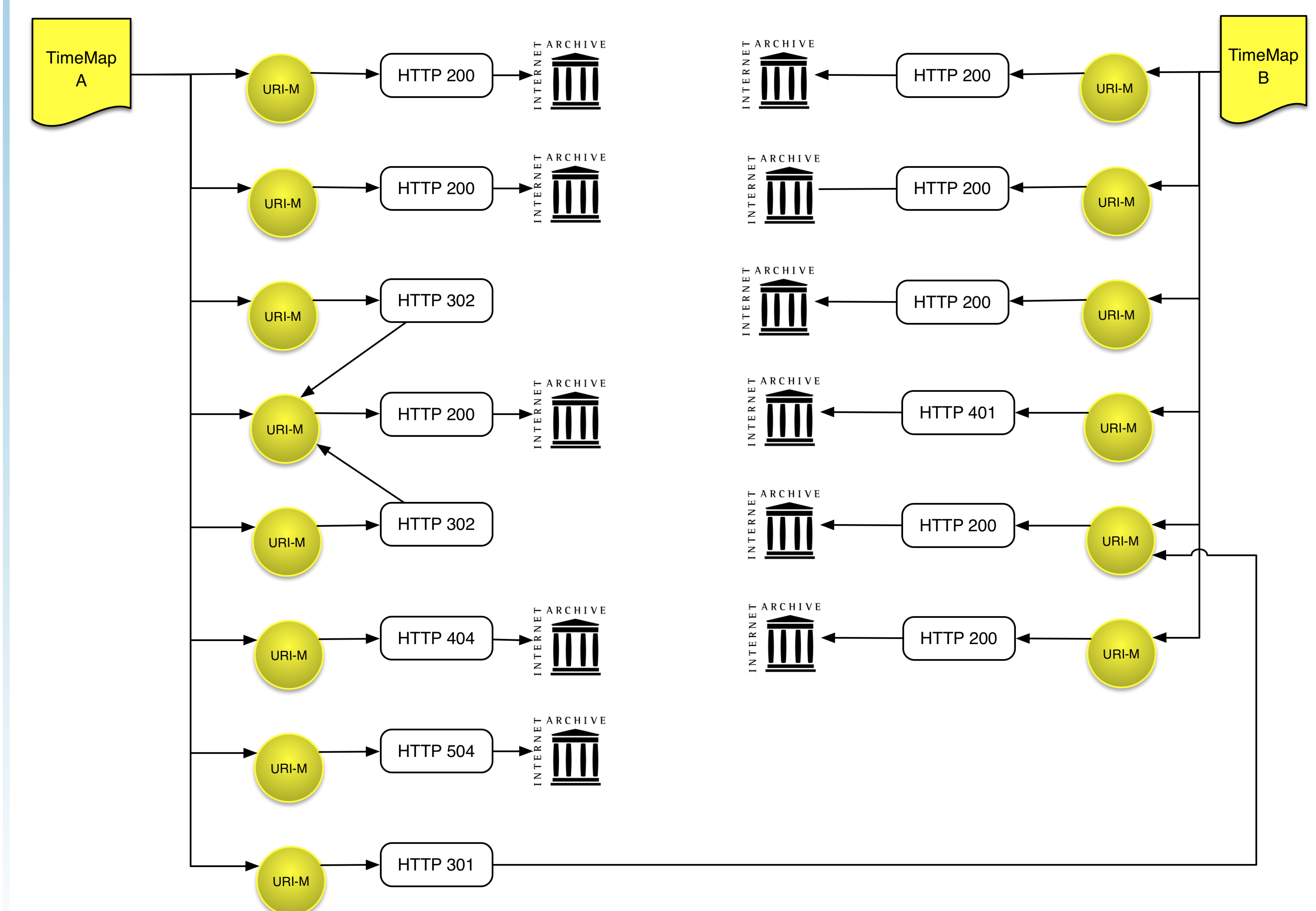
REDIRECTION FOR ARCHIVED URI-Rs

host	% 3XX	% 200	M _{TM}	DI	year	M _{TM}	M _{RC}	DI
google	84.89	15.11	695,525	0.178	1998	4	4	∞
yahoo	88.16	11.83	418,896	0.134	1999	19	19	∞
sourceforge	73.34	26.63	31,408	0.363	2000	132	87	1.933
instagram	67.32	32.65	55,228	0.485	2001	1,185	579	0.955
vimeo	57.04	42.94	199,262	0.752	2002	176	137	3.513
cnn	49.97	50.01	87,148	1.001	2003	75	55	2.750
wikipedia	44.62	55.19	25,973	1.240	2004	197	143	2.648
whitehouse	44.57	55.24	26,006	1.243	2005	1,236	414	0.504
stanford	62.14	37.84	19,309	0.609	2006	735	483	1.917
princeton	60.10	39.88	9,355	0.663	2007	1,055	842	3.953
columbia	48.01	51.88	9,882	1.082	2008	1,376	894	1.855
harvard	33.91	65.96	7,699	1.948	2009	6,074	4,335	2.493
caltech	33.13	66.86	5,474	2.017	2010	9,326	6,530	2.335
mit	26.57	73.24	6,379	2.763	2011	20,634	9,279	0.817
gatech	26.03	73.94	3,907	2.841	2012	102,533	16,240	0.188
ufl	24.76	75.23	4,927	3.038	2013	228,405	25,203	0.124
vt	23.07	76.92	4,061	3.334	2014	164,865	22,738	0.160
lsu	15.06	84.93	2,974	5.638	2015	17,978	11,286	1.686
nsu	13.82	86.00	1,208	6.233	2016*	139,520	5,805	0.043
odu	9.727	90.27	1,727	9.279	<i>google.com redirection over time</i>			
tcc	5.429	94.57	884	17.41	<i>*Data collected May 2016</i>			

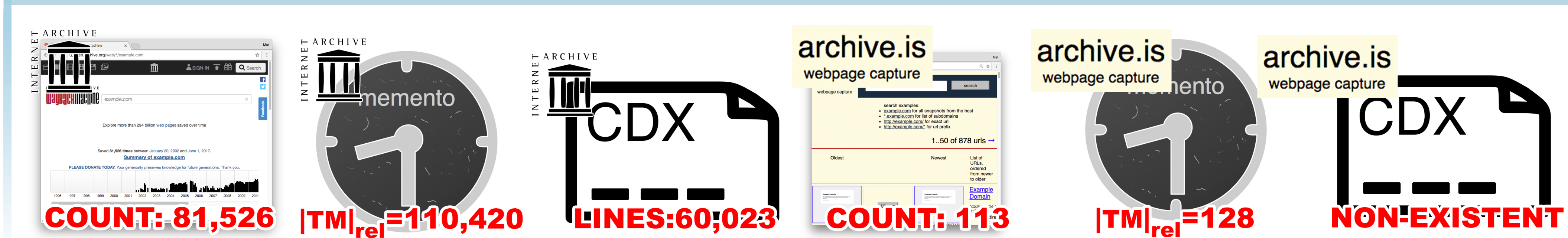
INTERSCHEME/INTER-SUBDOMAIN REDIRECTS (GOOGLE.COM)



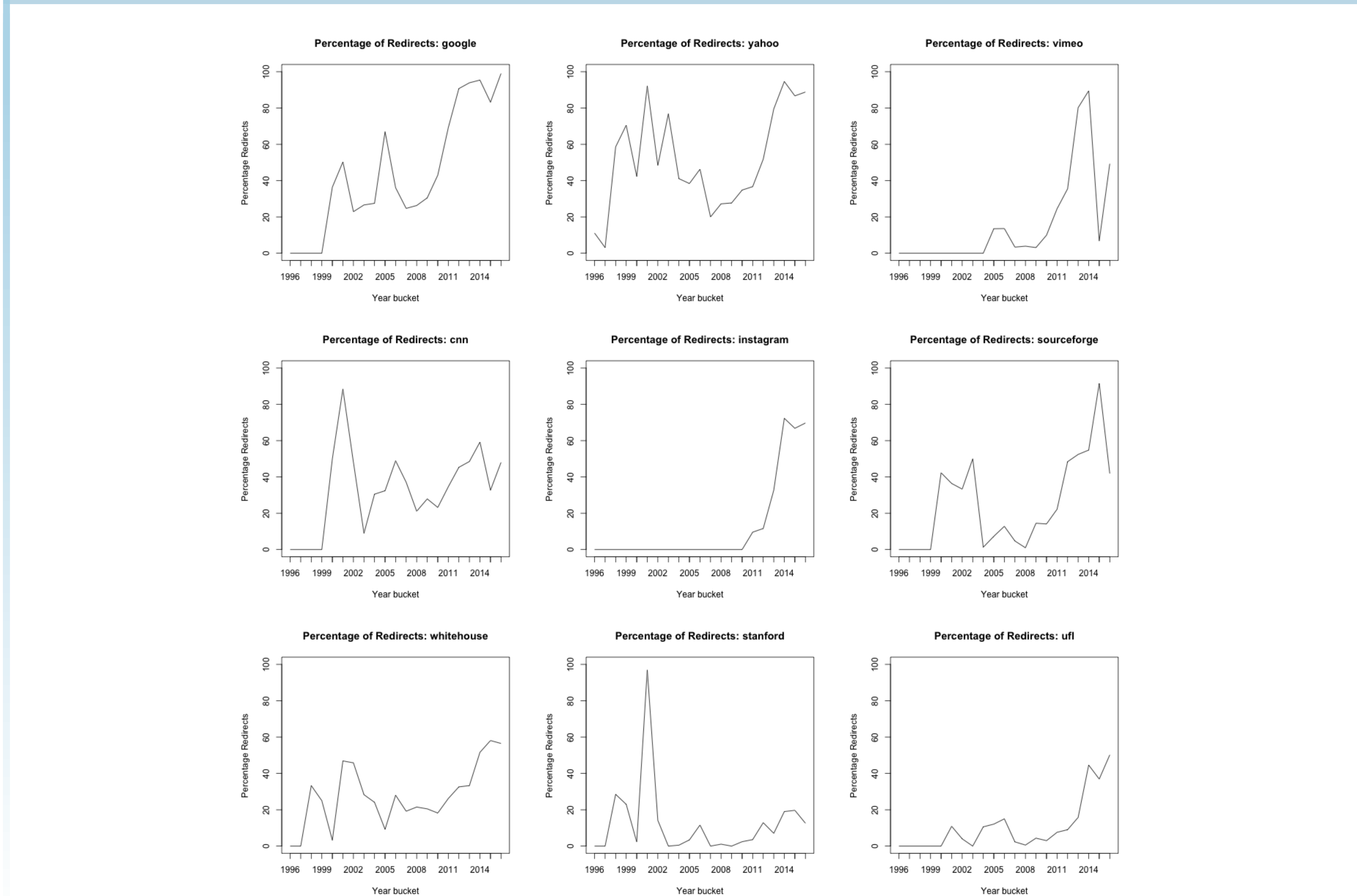
DEREFERENCING URI-MS FROM A TIMEMAP



INTRA-ARCHIVE VARIANCE IN "COUNT"



REDIRECTION OVER TIME FOR NINE URI-Rs



CONCLUSIONS

- Counting captures is impossible using only a TimeMap.
- Non-standard capture endpoint listings differ in count interpretation intra-archive.
- TimeMaps need more info about mementos (e.g., status code) for accurate counting.
- Alluding to non-standard endpoints (e.g., CDX) may be N/A when aggregating with archives without these endpoints.

REFERENCES

- [1] Mat Kelly, Lulwah M. Alkwai, Michael L. Nelson, Michele C. Weigle, and Herbert Van de Sompel. 2017. Impact of URI Canonicalization on Memento Count. (March 2017). arXiv:1703.03302
- [2] Herbert Van de Sompel, Michael Nelson, and Robert Sanderson. 2013. HTTP Framework for Time-Based Access to Resource States – Memento. IETF RFC 7089. (December 2013)