

Reference Rot: An Emerging Threat to Transparency in Political Science

Aaron L. Gertler, *Independent Scholar*

John G. Bullock, *University of Texas at Austin*

ABSTRACT Transparency of research is a large concern in political science, and the practice of publishing links to datasets and other online resources is one of the main methods by which political scientists promote transparency. But the method cannot work if the links don't, and very often, they don't. We show that most of the URLs ever published in the *American Political Science Review* no longer work as intended. The problem is severe in recent as well as in older articles; for example, more than one-fourth of links published in the *APSR* in 2013 were broken by the end of 2014. We conclude that "reference rot" limits the transparency and reproducibility of political science research. We also describe practices that scholars can adopt to combat the problem: when possible, they should archive data in trustworthy repositories, use links that incorporate persistent digital identifiers, and create archival versions of the webpages to which they link.

The past decade has given rise to unprecedented concern about our ability to verify the claims that appear in peer-reviewed journals (e.g., Dafoe 2014; Gerber and Malhotra 2008). In the face of this concern, political scientists have embraced new and occasionally elaborate practices that promise to bolster the transparency of published research (e.g., Lupia and Elman 2014; Monogan 2015). But one of the main tools that we use to promote transparency is one of the simplest: we publish links to websites that hold our data and to the files that we have used to analyze those data. This is a common and time-honored practice. It should do much to promote transparency. Still, it does far less than it might—because many of the links in published research are broken.

In this article, we examine all of the URLs that appeared in the *American Political Science Review* between 2000 and 2013. Most of these links are broken. The problem is not confined to articles that were published when political scientists were only beginning to use the Internet: even in articles published as late as 2012, more than 40% of all links are broken.

These findings suggest that the practice of linking to relevant information is not doing as much as it should to promote transparency in political science. It may even be creating a false sense of transparency among readers. In this article, we outline the problem and describe its manifestation in the *APSR*. We also describe three simple steps that scholars can take to combat the problem: when possible, they should archive data in trustworthy

repositories, use links that incorporate persistent digital identifiers, and create archival versions of the webpages to which they link.

BACKGROUND: STANDING ON QUICKSAND

A URL—a "uniform resource locator" or, colloquially, a "link"—is a reference to a resource that specifies the resource's location and the protocol by which it can be accessed. The protocol is usually the hypertext transfer protocol, or "http"; typical examples of links include <http://www.apsanet.org> and <http://example.com/example.html>. Including links in one's papers to data and statistical code should do much to promote reproducibility. But in practice, this strategy often fails. The reason is that the links are often broken: they do not lead to the intended resource. The phenomenon is known as "reference rot" (e.g., Van de Sompel and Treloar 2014, 197).

Published links may fail to work for many reasons, few of which are the authors' faults. The site that once contained the relevant data may no longer exist. Or its structure may have changed, such that old links to a particular part of the site no longer work. When scholars move to new universities, for example, their old universities typically stop hosting their online materials, and links to those materials typically stop working. Reorganizations of many major websites, including <http://www.un.org> and <http://whitehouse.gov>, have also broken thousands of published links (Zittrain, Albert, and Lessig 2014, 187). The reasons are not surprising, but the extent of the problem is. Zittrain, Albert, and Lessig (2014, 180) find that 70% of the URLs cited in leading law journals no longer work as intended. High rates of reference rot have also been found in communication (39%, Dimitrova and Bugeja 2007) and public health (49%, Wagner et al. 2009).

Aaron L. Gertler is an independent scholar. He can be reached at aaronlgertler@gmail.com.

John G. Bullock is assistant professor of government at the University of Texas at Austin. He can be reached at john@johnbullock.org.

As Lepore (2015, 36) has it, relying on published links in scholarly research is like “trying to stand on quicksand.”

As we show below, the problem seems to be great in political science as well. Most of the links ever published in the *American Political Science Review* are broken. The problem is noteworthy in recent volumes of the *APSR*, and in older volumes—those from the previous decade—it is overwhelming.

In the decade from 2000 through 2009, there is not one year for which the majority of links are still functioning.

DATA, RESEARCH DESIGN, AND FINDINGS

The *American Political Science Review* is the best-cited political science journal (Thomson Reuters 2015). We considered every article published in the *APSR* between 2000 and 2013—633 articles in all. Within each article, we searched for the terms **http**, **www**, **.com**, **.edu**, **.org**, and **.gov**. These searches produced a list of URLs—links to data, statistical code, software, and other resources.

Counting URLs no more than once per article, we recorded 1,135 URLs in the 56 issues of the *APSR* that were published from 2000 through 2013.¹ Of these URLs, 1,055 were unique. By examining each article, we determined that 418 of these URLs (37%) pointed to information that readers would need to reproduce the authors’ findings. In May 2016, we followed (“clicked on”) each URL and examined the result to determine whether the URL was working. URLs were classified as working if and only if they led directly to the intended resource. All others were classified as broken.²

Figure 1 reports our main results. Fully 59% of the links published in the *APSR* between 2000 and 2013 are now broken. It is to be expected that the problem increases with age, but the speed of decay is remarkable. In no year are more than two-thirds of links still functioning. In the decade from 2000 through 2009, there is not one year for which the majority of links are still

functioning. Only 30% of the links published in 2009 are still working. A slight majority of links published since then are working, but if history is any guide, most of these still-working links will soon be broken, too. The correlation between year and the working-link rate is $r = .84$.

When we restrict our focus to reproducibility URLs—the links that should lead to materials that one needs to reproduce

the authors’ results—our findings are more encouraging, but not by a lot. Fully 53% of these links are broken. Of the 14 years of articles that we examine, there is only one year—2013, the final year—for which two-thirds of these links are still working. And time is a strong predictor of decay for these links, too: for them, the correlation between year and the working-link rate is $r = .79$.

One may imagine that the high rates of reference rot in the *APSR* are due to authors linking to resources that are stored on personal sites rather than institutional sites. By this reasoning, personal sites—typically sites that individual scholars maintain to make their research available to the public—are more likely than institutional sites to change in ways that break URLs. This is a plausible explanation, but it is not correct. In every year from 2000 through 2010, personal-site URLs were less likely to be broken than institutional-site URLs. In all, 52% of the personal-site URLs in our dataset (99 of 192) are broken, against 60% of the institutional-site URLs (567 of 943).³

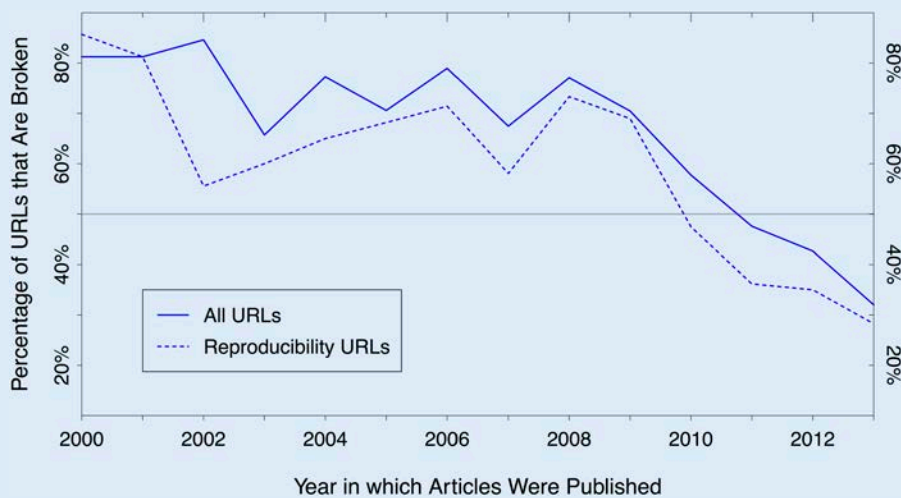
Rates of Decay

How long do links continue to work, and at what age should we expect that most links will be broken? Figure 1 offers no firm answers to these questions. For example, one might note from

the figure that most links published in 2010 are broken, and that most links published after 2010 continue to work. Given that we evaluated these links in 2016, one might therefore infer that most links published in year x will continue to work until year $x + 6$. But figure 1 does not warrant inferences of this sort. Among other problems, links published in more recent years may be either more or less robust than those published in earlier years, in which case a rule like “ $x + 6$ ” cannot be said to apply generally across years. The fundamental problem with making inferences of this sort from figure 1 is that the figure reports an examination at a single point in time—May 2016—of links from 14 different years.

To shed more light on rates of link decay, we turn to figure 2, which depicts results from two different

Figure 1
Broken links in the *American Political Science Review*, 2000–2013



Data are from a review of links undertaken in May 2016.

investigations: the May 2016 investigation described above, and an earlier investigation that we undertook in November 2014. The figure thus permits one to see how “reference rot” increased over an 18-month period among links published in any given year. Begin with the left-hand panel of the figure, which suggests that rates of link decay vary little by the age of the links: during the 18 months that we considered, reference rot increased steadily among links published in almost every APSR volume from 2000

expire very quickly. Many of them may be broken even before they are published.

Can Resources Associated with Broken Links Still Be Found Online?

Of course, just because links to critical resources are broken does not mean that the resources themselves are unavailable. They may remain online, and careful searching may reveal their location.

These results suggest that many links published in the APSR expire very quickly.

through 2013.⁴ Averaging across the 14 volumes in the study, reference rot increased by eight percentage points during the 18 months in question. The rate of decay was greatest by far among links published in 2002 and 2003: for these links, rates of reference rot increased by 27 percentage points and 14 percentage points, respectively. But in general, the left-hand panel of figure 2 suggests that reference rot is a steady process.

The right-hand panel of figure 2 shows that the process is somewhat different when we confine our study to links that one needs to reproduce the authors’ analyses. In this case, decay was greatest by far for links published from 2006 through 2011: fully 15% of the reproducibility links published in this period broke during the 18-month period under consideration. By contrast, broken-link rates increased by only 3% among links published from 2000 through 2005, and by only 8% among links published in 2012 and 2013.

Figure 2 illustrates one further critical feature of reference rot in the APSR. It shows that, by November 2014, 27% of all links published in the APSR in 2013 were broken. Reproducibility links fare somewhat better, but even 20% of the links of this type that were published in 2013 had stopped working by November 2014. These results suggest that many links published in the APSR

To examine this possibility, we randomly sampled 100 broken reproducibility links—links to resources that one needs to reproduce authors’ results—from our population of broken links. For each sampled link, a trained graduate student who had read the relevant article was given five minutes to locate the missing resource.

This exercise revealed that three of the selected links were working; we had mistakenly coded them as broken. We were able to locate the “missing” resource for 55% of the remaining links (53 of 97). But in 23% of the cases (22 of 97), we were able to locate only a related resource—for example, a later version of a paper, the cited version of which seems to have disappeared from the Internet. And in another 23% of cases (22 of 97), we were entirely unable to locate the relevant information.

Given more than five minutes, it is possible that the results would have been better. But we suspect that they would not have improved much: five minutes allows for a substantial amount of online searching. Our exercise thus suggests that about half of the resources that are associated with broken links and needed to reproduce published results cannot be found online at all, either because they have been supplanted by later versions or because they have vanished altogether.

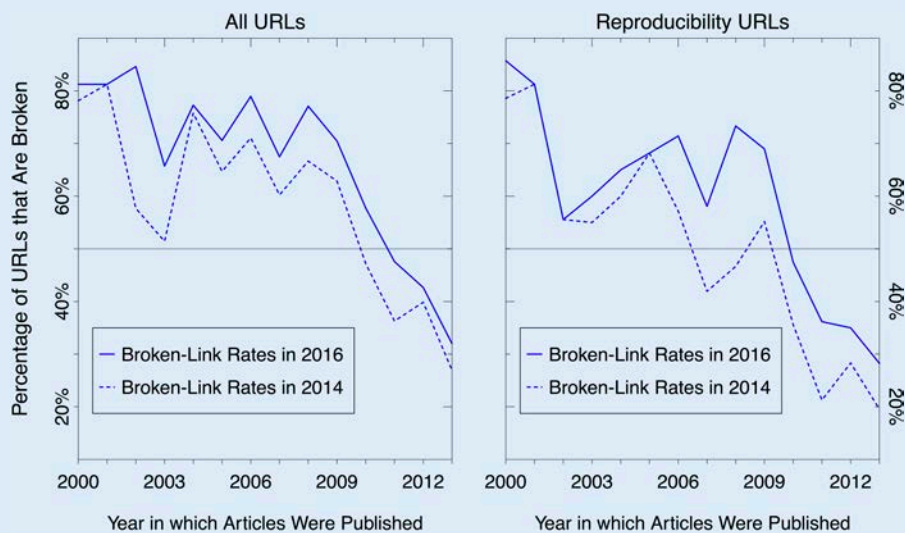
RECOMMENDATIONS

We close with three recommendations. To the extent feasible, scholars should archive their original data and related materials in trustworthy digital repositories, use links that incorporate persistent digital identifiers, and create long-lasting, archival versions of the webpages to which they link. Each of these practices is related to the other, and each has been recommended as a way to promote research transparency in a general sense. But as reference rot is our focus, we draw out the particular ways in which each practice can mitigate reference rot.

Host Data in Trusted Digital Repositories

One reason for reference rot is the disappearance from the Internet of resources—data, code, and even entire websites—that it once contained. Another reason is that the structure of

Figure 2
Broken Links in the American Political Science Review, 2000–2013: Broken-link Rates Measured in 2014 and 2016



websites changes over time, rendering old URLs useless. Depositing data and related materials into trustworthy digital repositories will mitigate both problems. These repositories exist to preserve replication materials and to ensure that they remain accessible to a wide audience. They are typically easy to access and free to use. And inasmuch as they are backed by institutions of long standing, they are unlikely to disappear. For example, the

your old number even after you move from one state to another, persistent identifiers permit readers to find relevant materials even after those materials move from one location to another. From the author's perspective, the point is simple: links created today are more likely to work in the future if they incorporate persistent identifiers. (For general discussions, see Askitas 2010 and Altman and King 2007, Section 3.)

Persistent identifiers are thus roughly analogous to mobile telephone numbers: much as your friends can continue to call you via your old number even after you move from one state to another, persistent identifiers permit readers to find relevant materials even after those materials move from one location to another.

repository of the Interuniversity Consortium for Political and Social Research (ICPSR) has been run out of the University of Michigan since 1962 (Center for Research Libraries 2006, 7); another prominent repository, the largest Dataverse network, is backed by Harvard's Institute for Quantitative Social Science (Harvard Dataverse 2015; King 2007).

Moreover, many digital repositories are party to "syndicated storage" agreements that oblige them to assume the task of publishing data from other repositories if those other repositories disappear. The point of these agreements is to protect data from physical threats (e.g., earthquakes), human threats (e.g., electronic attacks), and institutional threats (e.g., economic failure). Most notably, the repositories of the Data Preservation Alliance for the Social Sciences—including the ICPSR, Harvard's IQSS, and the University of North Carolina's Odum Institute—are all bound by a joint syndicated storage agreement (Data-PASS 2014).

When Possible, Use Links that Incorporate Persistent Digital Identifiers

Many online resources—especially datasets and journal articles—can be found through multiple URLs. For example, the 2012 American Identity and Representation Survey can be found at <http://www.icpsr.umich.edu/icpsrweb/RCMD/studies/36410/version/1> or at <http://doi.org/10.3886/ICPSR36410.v1>. But only one of these links—the latter—incorporates a *persistent digital identifier*. When scholars can choose which type of link to use, they should choose the type that incorporates a persistent identifier.

Persistent identifiers exist in multiple formats, but all are strings of characters that uniquely identify a digital resource independent of its location in the Internet. They can be embedded in URLs. They cannot be created by individual scholars, but they can be created by publishers (who assign them to journal articles) and by digital repositories (which can assign them to individual files). And as figure 3 shows, they combat reference rot by inserting a layer of abstraction between URLs and the resources to which they refer.

When persistent identifiers are used, the many links to a resource (say, a dataset) need not change when the resource's location changes. Only the connection between the persistent identifier and the resource—which is defined by the organization that issues the identifier—will need to be updated. Persistent identifiers are thus roughly analogous to mobile telephone numbers: much as your friends can continue to call you via

An example will further clarify the logic of persistent digital identifiers. When Deborah Schildkraut deposited the data from the American Identity and Representation Survey into the ICPSR repository, the ICPSR assigned a unique identifier to the dataset: 10.3886/ICPSR36410.v1. Like other digital identifiers, this one can be embedded within a URL: <http://doi.org/10.3886/ICPSR36410.v1> is a link to the dataset.⁵ Critically, the identifier will remain the same even if the location of the dataset changes—for example, even if the file structure of the ICPSR repository changes. As a member of the Data-PASS network, the ICPSR has committed itself to maintaining the identifier so that it always remains current. The link that contains the identifier will therefore always work as well.

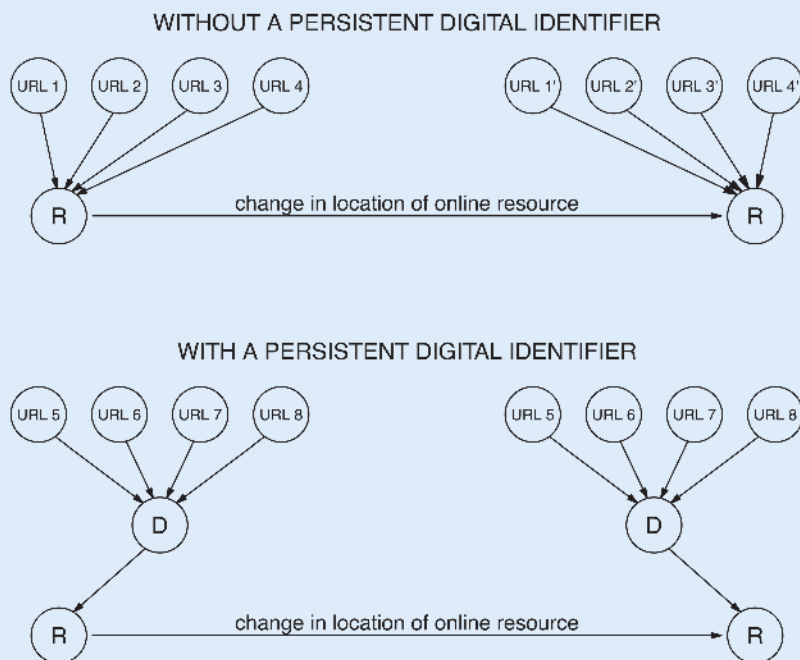
One may be skeptical of institutional promises to maintain persistent digital identifiers. But our data suggest that those promises are being kept. Only thirteen links in our dataset incorporate persistent digital identifiers, but of those thirteen, twelve are still working as intended. This working-link rate, 92%, far exceeds the overall working-link rate in any of the years that we analyzed. Moreover, many organizations have committed themselves to maintaining identifiers that were originally issued by other organizations, should those other organizations go out of business (International DOI Foundation 2015). This commitment, too, should instill confidence in the persistence of digital identifiers.

Archive Webpages

As figure A1 of the online appendix shows, nearly half of all APSR URLs are "bibliographic URLs" that are used to support inherently irreproducible claims. Most links to speeches, press releases, government documents, and characterizations of historical or current events (such as one finds in newspapers) fall within this category.

Documents of this sort typically do not have persistent digital identifiers. In principle, they can be archived in trustworthy digital repositories, just as datasets are. But in practice, it is often difficult for authors to archive webpages in this way, and it is often difficult for conventional repositories to store them and display them. This is especially true because many modern webpages are complex amalgams of dynamic content from many different sources. Attempting to save such pages without specialized software will typically fail to yield a true copy of their content (Van de Sompel and Treloar 2014, Section 4.2; see also Ainsworth, Nelson, and Van de Sompel 2015).

Figure 3
Using Persistent Digital Identifiers to Combat Reference Rot



In the top panel, all URLs point directly to an online resource (“R”). When the location of the resource changes, every URL must change or it will fail to work. The bottom panel shows how the use of persistent digital identifiers remedies the problem. In this scenario, URLs point not to the resource itself but to a digital identifier (“D”). When the location of a resource changes, the URLs do not need to change as well. Only the digital identifier—typically maintained by a publisher or a repository—needs to be updated.

Several solutions to the problem have long existed, the best-known of which is the Internet Archive (Lepore 2015). But the Internet Archive and related tools have two large drawbacks. First, the Archive’s contents are entirely at the discretion of website owners, who can cause any content that they own to be deleted; like its peers, the Internet Archive is “not interested in offering access to websites or other Internet documents whose authors do not want their materials in the collection” (Internet Archive 2015). Second, these organizations are supported largely by volunteers, which gives rise to concerns about their existence in the long term (Lepore 2015; Zittrain, Albert, and Lessig 2014). A system of distributed storage, buttressed by succession-planning agreements between partner organizations that take effect should any one organization fail, is preferable.

Such an organization now exists. Perma.cc, a service developed by the Harvard Library Innovation Lab and supported by a large consortium of academic libraries, was designed to combat the problem of reference rot in bibliographic URLs (Zittrain, Albert, and Lessig 2014). Perma (<https://perma.cc>) makes it trivial to archive webpages and to create persistent links to those archival versions. It resembles the URL shorteners with which many readers are already familiar, e.g., <http://bit.ly>, <http://goo.gl>: one enters a URL—perhaps a very long one—and Perma returns a new, condensed URL that is suitable for print. But unlike conventional URL shorteners, Perma archives the material at the original URL, ensuring both that the material will continue to exist and that the new URL will link to it even if the original link expires. Standards for the citation of archived pages have yet to be developed, but

emerging practice entails listing both the original and the archival URL if either is listed (e.g., Zittrain, Albert, and Lessig 2014; see also the references section of this article).

There are two qualifications to the recommendation that political scientists use Perma to archive webpages. The first is that they may not be able to use it. It was developed by law libraries, and although it is now expanding to other libraries, it remains predominantly a law-library service. Unless one is affiliated with a registered library,⁶ the number of permanent archives that one can create with Perma is limited. (As of November 2016, the limit is 10 records per month.) If this limit proves an important constraint, we recommend archiving with the Internet Archive or another service. Notwithstanding the limitations of these alternatives, archiving webpages with them is preferable to not archiving at all.

The second qualification is that not every link that appears in an article should be a persistent link to an archived webpage. In some cases, authors may be confident that the institution currently serving a webpage will preserve it into perpetuity, and archival links may be superfluous in these cases. In other cases, authors may simply want to alert their readers to the existence of certain websites—www.census.gov, www.voteview.com, and so on—without calling their attention to specific information on those sites.

Here, too, there may be no reason to link to archival versions of the sites. But when authors are linking to pages that contain particular information that they have used to make their arguments, archival links such as those produced by Perma are usually in order.

CONCLUSION

Political scientists are extremely reliant on URLs to promote research transparency, but those URLs are often broken, rendering research anything but transparent. We have shown that most of the URLs published in the *American Political Science Review* between 2000 and 2013 no longer work as intended. Nearly three-fourths of the URLs published in the *APSR* in 2009 are broken, and the situation is worse for most earlier years. Even as late as 2012, more than 40% of the links published in the *APSR* are broken. And these percentages are certain to increase as time passes.

To combat the problem, we recommend three practices: authors should host data in digital repositories, use links that contain persistent identifiers, and create long-lasting, archival versions of many of the pages to which they link. We are not alone in urging some of these practices. In particular, the editors of the *APSR* and many other journals recently signed the Data Access and Research Transparency Joint Statement, thereby pledging to require authors to use trustworthy digital repositories and persistent digital identifiers (Data Access and Research Transparency 2014).⁷

These aspects of the DA-RT statement should mitigate reference rot—which seems, by our analysis of the *APSR*, to be a large impediment to research transparency in political science. Only by embracing new practices—including but not limited to those

prescribed by the DA-RT statement—can we ensure that future readers will be able to examine the sources on which we rely when we make our arguments.

SUPPLEMENTARY MATERIAL

To view supplementary material for this article, please visit <https://doi.org/10.1017/S1049096516002353>

ACKNOWLEDGMENTS

The authors thank Alex Branham, Allan Dafoe, Stephen Jessee, Thomas Leeper, and Chris Wlezien for comments on an earlier draft, Kim Dulin and Adam Ziegler for a helpful discussion, and Alex Branham again for excellent research assistance. ■

NOTES

1. We did not record links to personal websites if they appeared only in notes that provided authors' contact information.
2. Our analysis incorporated a number of finer distinctions—for example, between different types of broken links, and between “bibliographic,” “database,” and “reproducibility” links—that we lack the space to discuss here. See the online appendix for details.
3. “Personal-site” URLs include links to personal sites that are associated with .edu domains. See the online appendix for details.
4. During these 18 months, reference rot increased for every one of the 14 *APSR* volumes, save for the 2001 volume. Sixty-four links were published in the 2001 volume of the *APSR*, and none of those links broke during the 18 months of our study. This exception arises partly because almost all of those links were already broken by November 2014, when we began our investigation.
5. This particular identifier has the DOI format, defined by the International DOI Foundation (<http://doi.org>). Hence the “doi.org” prefix in the URL.
6. It is extremely simple for unregistered libraries to register with Perma. See <https://perma.cc/libraries>.
7. One qualification: the DA-RT statement mandates the use of persistent identifiers for datasets but not for other resources.

REFERENCES

Ainsworth, Scott G., Michael L. Nelson, and Herbert Van de Sompel. 2015. “Only One Out of Five Archived Web Pages Existed as Presented.” Presented at the Annual Meeting of the ACM Conference on Hypertext and Social Media, Cyprus. <http://dx.doi.org/10.1145/2700171.2791044> (accessed April 30, 2016).

Altman, Micah and Gary King. 2007. “A Proposed Standard for the Scholarly Citation of Quantitative Data.” *D-Lib Magazine* 13 (March-April).

Askatas, Nikos. 2010 June. “What Makes Persistent Identifiers Persistent?” German Data Forum. Manuscript. <http://ssrn.com/abstract=1639996> (accessed July 26, 2015). Archived at <http://perma.cc/S3FJ-NKZG>.

Center for Research Libraries. 2006. “ICPSR Audit Report.” http://www.crl.edu/sites/default/files/d6/attachments/pages/ICPSR_final.pdf (accessed August 27, 2015). Archived at <http://perma.cc/H47N-MJUA>.

Dafoe, Allan. 2014. “Science Deserves Better: The Imperative to Share Complete Replication Files.” *PS: Political Science & Politics* 47 (1): 60–66.

Data Access and Research Transparency. 2014. “Data Access and Research Transparency (DA-RT): A Joint Statement by Political Science Journal Editors.” http://media.wix.com/ugd/fa8393_dao17d3fed824cf587932534c860ea25.pdf (accessed August 27, 2015). Archived at <http://perma.cc/XTZ7-KGNM>.

Data Preservation Alliance for the Social Sciences. 2014. “Memorandum of Understanding.” http://data-pass.org/sites/default/files/Data-PASS_MoU_201504.pdf (accessed August 27, 2015). Archived at <http://perma.cc/NS5H-KS5H>.

Dimitrova, Daniela V. and Michael Bugeja. 2007. “The Half-Life of Internet References Cited in Communication Journals.” *New Media & Society* 9 (5): 811–26.

Gerber, Alan and Neil Malhotra. 2008. “Do Statistical Reporting Standards Affect What Is Published?” *Quarterly Journal of Political Science* 3 (October): 313–26.

Harvard Dataverse. 2015. “Harvard Dataverse Preservation Policy.” <http://best-practices.dataverse.org/harvard-policies/harvard-preservation-policy.html> (accessed July 27, 2015). Archived at <http://perma.cc/P5DL-KG2Y>.

International DOI Foundation. 2015. “DOI System and the Handle System.” <http://www.doi.org/factsheets/DOIHandle.html> (accessed August 27, 2015). Archived at <http://perma.cc/7A6B-ULML>.

Internet Archive. 2015. “Removing Documents from the Wayback Machine.” <https://archive.org/about/exclude.php> (accessed July 26, 2015). Archived at <http://perma.cc/ZY3P-RX45>.

King, Gary. 2007. “An Introduction to the Dataverse Network as an Infrastructure for Data Sharing.” *Sociological Methods & Research* 36 (2): 173–99.

Lepore, Jill. 2015. “The Cobweb.” *The New Yorker*, January 26, 33–41.

Lupia, Arthur and Colin Elman. 2014. “Openness in Political Science: Data Access and Research Transparency.” *PS: Political Science & Politics* 47 (1): 19–24.

Monogan, James E III. 2015. “Research Preregistration in Political Science: The Case, Counterarguments, and a Response to Critiques.” *PS: Political Science & Politics* 48 (3): 425–29.

Van de Sompel, Herbert and Andrew Treloar. 2014. “A Perspective on Archiving the Scholarly Web.” Proceedings of the 11th International Conference on Digital Preservation. http://public.lanl.gov/herbertv/papers/Papers/2014/iPres2014_Sompel_Treloar.pdf (accessed August 27, 2015). Archived at <http://perma.cc/F6QX-KJEU>.

Wagner, Cassie, Meseret D. Gebremichael, Mary K. Taylor, and Michael J. Soltys. 2009. “Disappearing Act: Decay of Uniform Resource Locators in Health Care Management Journals.” *Journal of the Medical Library Association* 97 (2): 122–30.

Zittrain, Jonathan, Kendra Albert, and Lawrence Lessig. 2014. “Perma: Scoping and Addressing the Problem of Link and Reference Rot in Legal Citations.” *Harvard Law Review* 127 (February): 176–99.