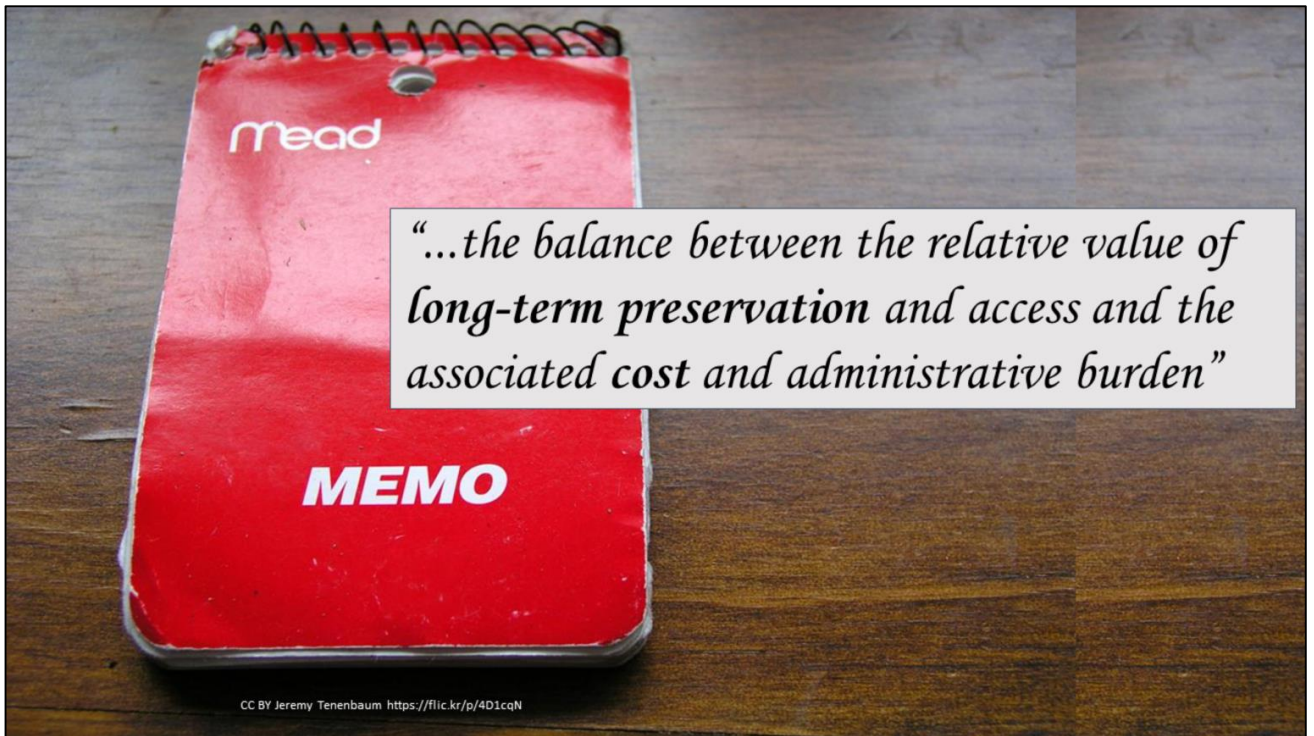


Many academic libraries were ready to act after the release of the US Federal Government's OSTP memo on access to federally funded research data. There was an expectation that in the plans generated by federal agencies in responses to the memo that libraries would play an important role. Though largely the plans did not end up mentioning institutional repositories for storage and access to research data, libraries and librarian can - regardless - take a leadership position in setting norms and expectations for preservation.

Photo based on <https://flic.kr/p/aQYUVM>



The OSTP memo listed specific things that the agency's plans needed to address. For preservation, it said, to find "the balance between the relative value of long-term preservation and access and the associated cost and administrative burden".

This directive explains the language seen in the responses. As a part of trying to craft policy for our future data repository, I examined all of the responses released so far, hoping to ensure our policies and system design are in line with requirements. Of the 17 agencies that have put out a plan, 14 mention preservation. By and large they use the phrasing directly from the memo "long-term" and a few talk about the cost-benefit analysis. Only one, USGS, gets more specific that you should follow the same rules as you would for USGS record retention.

[refs:https://www.whitehouse.gov/sites/default/files/microsites/ostp/ostp_public_access_memo_2013.pdf, this document links to all of the individual responses: <https://docs.google.com/spreadsheets/d/1PYOhBh6bglh6BkQFlpvNLOWlpzvQyguWAG8AkQMtU0s/edit#gid=0>]



Forever?

So, what does long-term mean in the eyes of the feds, does the absence of a definition imply forever? More importantly are researchers assuming that we will keep their stuff forever?

It is important to note that these plans are all drafts at this point. Things can, probably and hopefully will change.



Interested in clues as to the direction the drafts might head in, I looked at data sharing policies that predated the OSTP memo. 13 of the NSF directorates had a data sharing policy. Only 3 of these mentioned preservation in any capacity. And one had a time frame, Engineering, which said that the data needs to be publicly available for 3 years. So, that's one idea. IES, part of the dept of ed, also had a pre OSTP data sharing policy and it said the data needed to be made available for at least 10 years. Quickly, other countries: Canada's TriCouncil draft plan says to follow community norms and the UK's Concordant draft says 10 years, unless directed otherwise by a funder.

[refs: https://ies.ed.gov/funding/datasharing_implementation.asp,
<http://www.science.gc.ca/default.asp?lang=En&n=83F7624E-1>, <http://www.rcuk.ac.uk/RCUK-prod/assets/documents/documents/ConcordatOpenResearchData.pdf>]

010001010010000111010010101010000101010001010101101000101001010100101010010101000000101000001000101001000011101
0010101010000101010001010111010001010100101010010101001010100000101000010001010010000111010010101000010100
010101011010001010010101001010100101010000001010000100010100100001110100101010000101000101011010001010
0101010010101001010101000000101000010001010010001110100101010100001010100010101000101010001010100101
010100C 0010
001010C 1001
010101C 0010
101011C 0010
101001C 1010
100000C 0001
010010C 1010
101000C 0101
011010C 0101
0010101 0100
0000101 1010
0100001 0101
0000101 1011
010001010001010010101001010100000010100001000101001000111010010101000010101000101010001010100010101001
01010100101010100000010100000100010100100011101001010101000010101000101010100010101000101010001010100000
01010000010001010010000111010010101000010100010101011010001010010101001010100101010010000001010000010001010010
000111010010101000010101000101011010001010010101000101010001010100010101000000101000001000100010001000
0101010001010101010001010100101010010101000000101000001000101001000011101001010101000010101000101011010
00101010010101001010100101010100000010100001000101001000111010010101000010101000101010100010101001010
101001010101000000101000010001010010001110100101010000101000010101010001010100010101001010100000010
1000001000101001000011101001010100001010001010101101000101001010010101001010101000000101000001000101001000
111010010101010000101010001010101101000101001010100101010000001010000v010001010010000111010010101000010
101000101010100010101001010100101010010101000000101000010001010010000111010010101000010101000101011010001
01010010100101010010101010000001010000100010100100011101001010100001010101000010101010000101001010010101
00101010100000010100001000101001000111010010101000010100010101101000101001010100101010000001010

“Data preservation, or more specifically, *digital* data preservation, refers to the series of managed activities necessary to ensure continued access to digital materials for as long as necessary. This broad definition of data preservation refers to all of the actions required to maintain access to digital materials beyond the limits of media failure or technological change.”

From: International Federation of Data Organizations for Social Science, http://ifdo.org/wordpress/?page_id=18 (editorial emphasis)

The situation that Mahria has laid out is obviously less than ideal, but with little guidance from the Government, how long is long enough? Their lack of specificity does not justify ours; Libraries are currently moving to fill the perceived need and building out capacity to collect digital data. We have to make our own decisions, otherwise they will likely, unwittingly, be made for us, by our technology choices and the vagaries of funding, or lack thereof.

But let’s take a quick step back and talk about data preservation—what it is, why it’s important, and how the work gets done.

Like all good librarians, I’ll start with a definition. Already, I hope many of you know the challenges of digital preservation—and many come in that last clause: “beyond the limits of media failure or technological change” –and in a minute I’ll get into why that is especially relevant for data.

But first, why preserve data? Why worry about ensuring persistent access to it, over a large span of time?



"L'evoluzione della specie" by Aldo Cavini Benedetti on Flickr

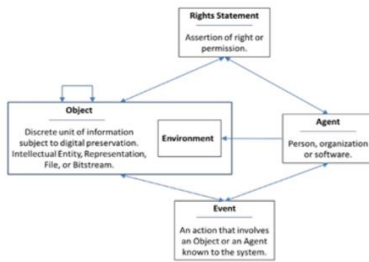
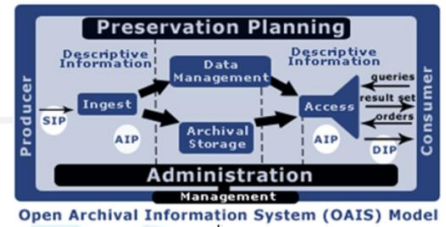


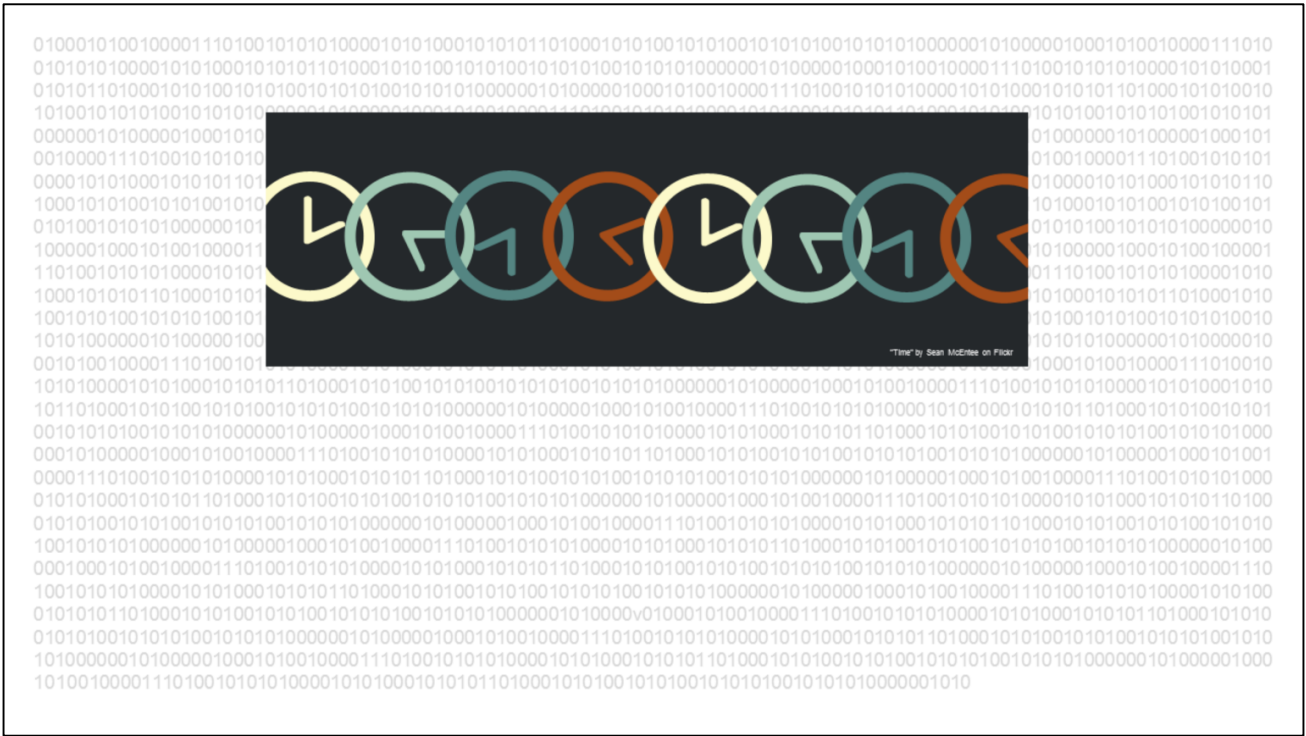
Figure 1: The PREMIS Data Model

Preserving research data is important. How do we do it? Obviously, my time is limited, but I can say you can use the same bag of tricks as other digital materials, geographically dispersed replication, checksums and monitored storage, harvesting and storing technical metadata, recording preservation metadata, migration, normalization—emulation, even. OAIS-compliant repositories, PREMIS metadata, these standards are out there. But there are two things about data that makes it a little different from other digital materials.

Premis data model: <http://www.loc.gov/standards/premis/v3/index.html>

Team digital preservation: <https://youtu.be/pbBa6Oam7-w>

OAIS: <http://public.ccsds.org/publications/archive/650x0m2.pdf>



But this brings us right back—what is the “right” amount of time to preserve data? How long is long enough? My answer is : “it depends”, some might be worth more effort than others because of the factors I have already mentioned. That is why planning for selection in data repositories is my last rant of the day. We all can attest that digital storage is NOT free, and that as computing power increases, so does the size and complexity of data. But then again, who could have predicted the longevity of Bumpus’ sparrows? Is it cheating to promise or imply “long term preservation” without the people, resources, and systems in place to assure it? Accept all materials university researchers care to self-deposit and only *recommend* appropriate file formats? Say “we’ll reevaluate when we run out of storage” or “we’ll reevaluate in ten years” trusting a monitored storage system to provide bit-level preservation and increasing technological savvy to take care of the obsolete hardware and software? Questions abound, but not asking them is asking for trouble.

