# End of Term Archive – 2008, 2012, 2016

Mark Phillips, Associate Dean for Digital Libraries, UNT Libraries
@vphill | mark.phillips@unt.edu

May19, 2017

# it all began a long, long, time ago, in a far away



https://flic.kr/p/4N2jHU

https://flic.kr/p/4JNkLE

National Library of Australia

nla.int-nl39859-al1-v

# original end of term web archive partners

INTERNET ARCHIVE

LIBRARY OF CONGRESS

University of California
CDL
California Digital Library

NORTH TEXAS

GPO

**for 2008/2012/2016 - all IIPC & NDIIPP/NDSA partners**

# extant gov web archiving efforts

## Capture, Preservation, & Access

- LOC: .gov, election, other
- GPO: agency sites, often ephemeral
- NARA: congressional web harvest every 2 years
- IA: global & curated crawls
- Agency-level: NIH/NLM, DOE, DOL, HHS, CMS, others, using AIT or comm tools
- UNT & Others: Topical .gov collecting

## Community Efforts

- Federal Web Archiving Group
  - most of those at left plus other feds
- Research Initiatives
  - academic
  - NGO or watchdog
- Citizen Driven
  - grassroots efforts
- End of Term
  - focused but large-scale multi-institutional project

# goals of the end of term project



United States Central Command
Sept 16, 2008

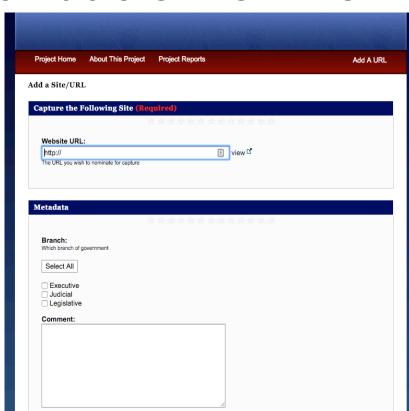U.S. Department of State Official Blog
Feb 13, 2013

Healthcare.gov Twitter
Feb 15, 2013

- □ work collaboratively to preserve public U.S. Government websites
- □ document federal agencies' presence on the web during the end of Presidential terms
- □ enhance the existing research collections of the partner institutions
- □ raise awareness about the need for preservation
- □ engage with researchers and subject experts

# eot collaborative distribution of work

- IA: crawling, preservation, access, full-text search
- LC: crawling, preservation, data transfers
- UNT: nomination tool development, crawling, nomination mgmt, preservation, access
- CDL: web portal, metadata
- GPO: URL nomination, outreach
- All: URL contributions, outreach, project management
- Others: URLs, education

Some variance of roles between 2008 & 2012 (and for 2016)
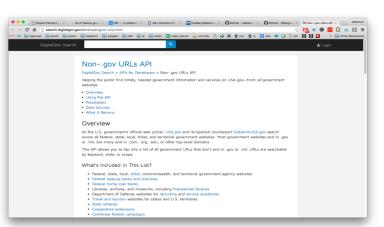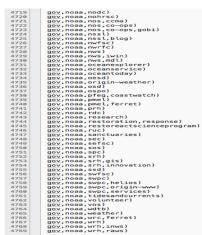
major funding brought to you by….

# no one

# defining the "government web presence"



**Stanford WebBase Project**

USA.gov

DIGITALGOV

**U.S. Digital Registry**

DATA.GOV

Non-.gov URLs API

NATIONAL ARCHIVES
2004 crawl list of URLs

unitedstates / congress-legislators

congress-legislators / legislators-social-media.yaml

# and people like you!

## End of Term Presidential Harvest 2008

**Institutions:**

Select a nominator below to see its associated URLs.

California Digital Library - 2563
Congressional Research Service, Library of Congress - 6
Cornell University ILR School - 11
Georgetown Law Library - 32
Institute of Environmental Science and Research (NZ) - 1
Internet Archive - 2184
Library of Congress - 113
Library of Congress, African and Middle Eastern Division - 17
McNeese State University - 2
New Hampshire Law Library - 38
Texas A&M University-Kingsville - 3
Thurgood Marshall Law Library, Univ. of Maryland School of Law - 1
U.S. Dept. of Labor, Office of Public Affairs-Division of Enterprise Communications - 10
UC Santa Cruz - 133
UCLA Charles E. Young Research Library - 1
University of Alabama - 4
University of Delaware - 18
University of Texas at Austin - 43
UNT Libraries - 5

## End of Term Presidential Harvest 2012

**Institutions:**

Select a nominator below to see its associated URLs.

Center for Medicare & Medicaid Innovation, Centers for Medicare & Medicaid Services, HHS - 1
Defense Commissary Agency (DeCA) - 1
Internet Archive - 3
Library of Congress - 31
National Archives and Records Administration - 1
OCC - 1
Pratt Institute - 128
Pratt Institute - SILS - 702
Pratt SILS - 683
Schoolcraft College - 1
U.S. Access Board - 2
U.S. Army Corps of Engineers, HECSA Library - 1
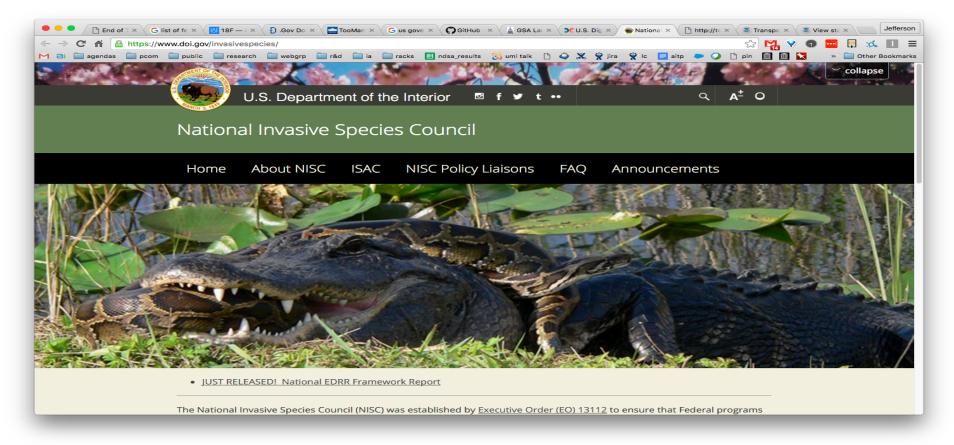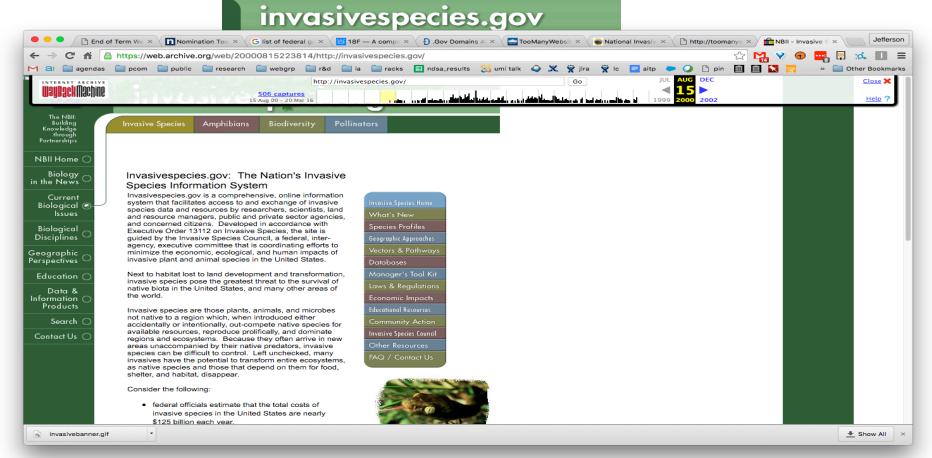UNT - 83
UNT Libraries - 4

# 2016 Was Crazy

# Some challenges of the .gov Web

# .gov websites proliferate like invasive species

# and yes, invasivespecies.gov once existed

# some are non-public or unlisted
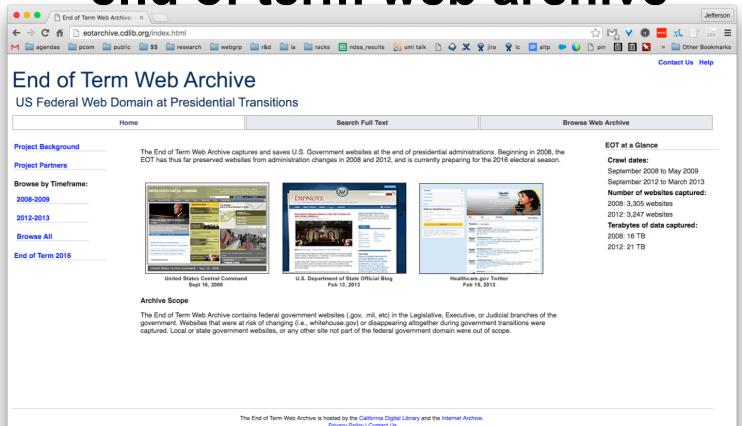
# "web waste" & preservation mentalities

# EOT Collections to Date

- 2008 – 16 TB, 160 million URLs/files

- 2012 – 35 TB, 194 million URLs/files

- 2016 – ~250 TB, >350 million URLs/files

# Currently transferring data between partners

# end of term web archive



## http://eotarchive.cdlib.org/