

Challenges facing the preservation of born-digital news applications

Katherine Boss, New York University Libraries

Meredith Broussard, New York University, Arthur L. Carter Journalism Institute

IFLA News Media Section Conference. Hamburg, Germany, April 2016



News Applications

theguardian

ZEIT  ONLINE

Chicago Tribune

politnetz.ch

What is a news application, or “news app”?

- Single data journalism project that is custom-built by the news agency
- From an archiving perspective, boundaries are hard to define

Who is creating them and why should we care to archive them?

- Creation of news apps and interactive data journalism is exploding
- Represent some of the most innovative and complex journalism



 FiveThirtyEight

The New York Times



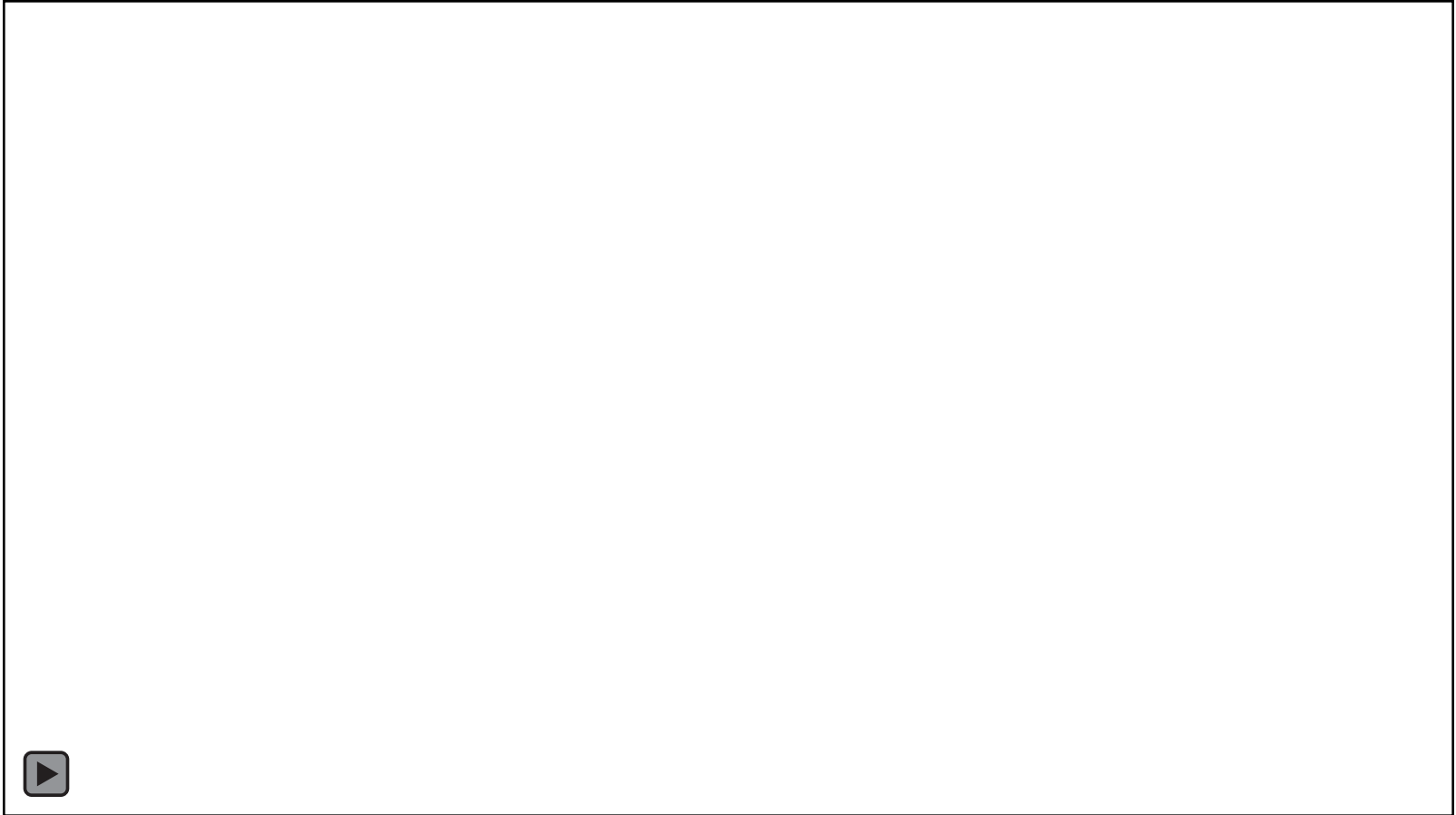


The Problem

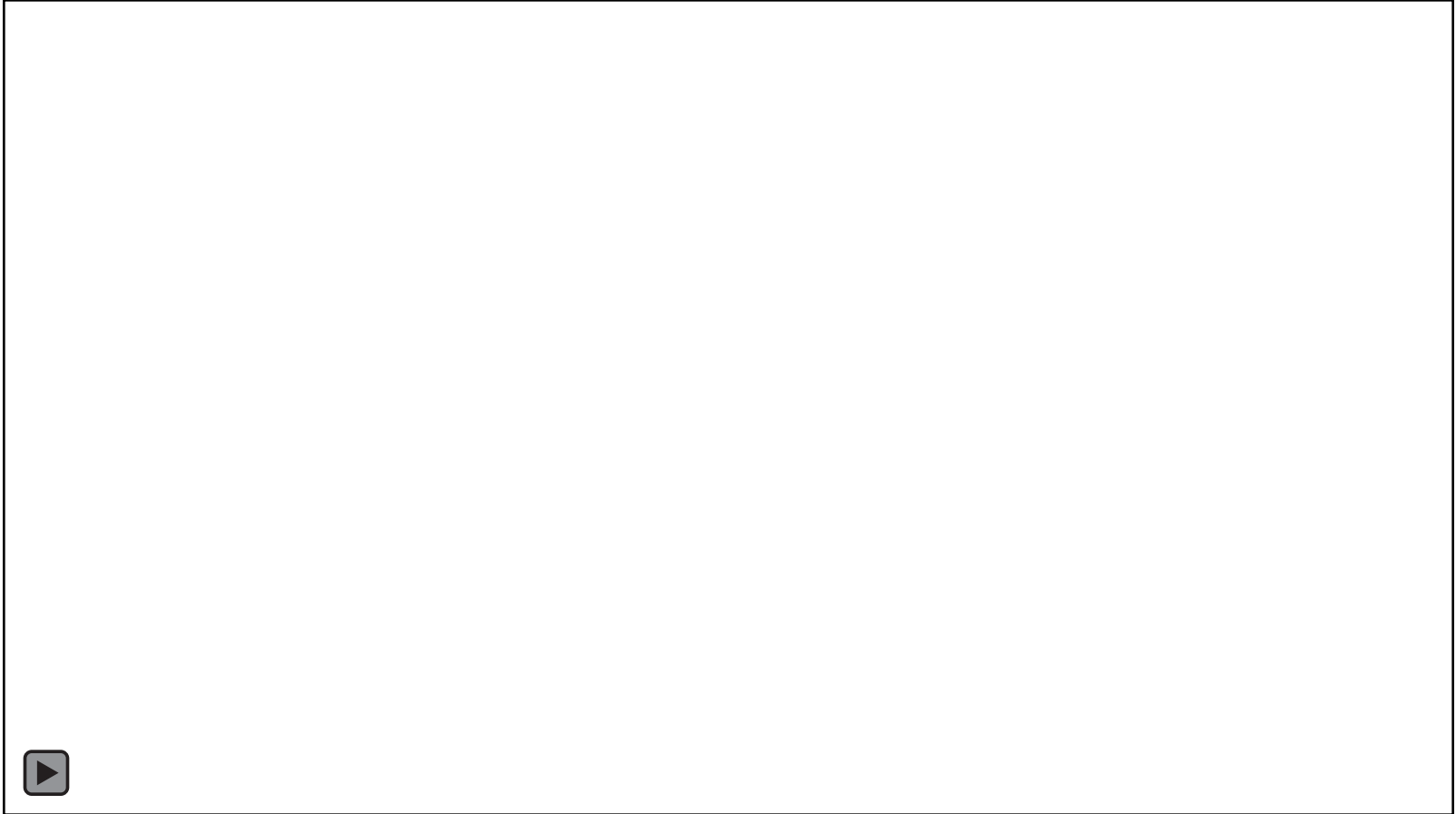
Technical Issues: News apps are dynamic, not static, digital objects

- Complex software dependencies, including the operating system, programming languages, compilers, software libraries, and so on (“dependency hell”)
- Cannot be fully captured through current web archiving
- External APIs (Google Maps etc) present extra challenges

ProPublica – screencast of live Dollars for Docs news app



Internet Archive snapshot of Dollars for Docs, Mar. 31 2016



Screencast of Live Workers' Comp Benefits news app



Internet Archive snapshot of Workers' Comp, Mar. 8 2016



http://projects.propublica.org/graphics/workers-compensation-benefits- GO

268 captures
4 Mar 15 - 17 Apr 16



Close X

Help ?



DONATE

Workers' Comp Benefits: How Much is a Limb Worth?

by Lena Groeger and Michael Grabell, ProPublica, and Cynthia Cotts, special to ProPublica, Mar. 5, 2015

If you suffer a permanent injury on the job, you're typically entitled to compensation for the damage to your body and your future lost wages. But depending on the state, benefits for the same body part can differ dramatically. [Related Story »](#)

Ever filed for workers' comp? [Help ProPublica investigate.](#)

Select a state to see the maximum it pays for different body parts.

The average maximum compensation for one **Ring Finger** in **The USA** is **\$14,660**

Oregon	Federal	Montana	Illinois	Iowa	Nevada	Pennsylvania	District of Columbia	New Hampshire	Tennessee	Vermont	Kentucky	North Carolina	Connecticut		
\$82,420	\$47,167	\$42,480	\$38,768	\$36,175	\$30,912	\$28,530	\$27,794	\$26,676	\$25,758	\$24,239	\$24,042	\$23,000	\$20,811		
New York	Delaware	Hawaii	Virginia	South Carolina	Indiana	Michigan	Missouri	Ohio	Maine	Georgia	Nebraska	South Dakota	Kansas	New Mexico	Louisiana
\$20,216	\$19,967	\$19,650	\$19,340	\$19,151	\$18,204	\$18,040	\$17,364	\$17,240	\$16,332	\$15,750	\$15,220	\$14,100	\$13,068	\$12,991	\$12,600
Arizona	Arkansas	Washington	West Virginia	Wisconsin	Idaho	Mississippi	Texas	Utah	New Jersey	Alaska	Wyoming	Oklahoma	North Dakota	Florida	California
\$11,929	\$11,328	\$10,644	\$10,638	\$10,465	\$9,474	\$9,272	\$9,030	\$8,959	\$8,892	\$8,850	\$7,548	\$7,106	\$6,400	\$6,315	\$6,090
Maryland	Alabama	Rhode Island	Massachusetts	Minnesota	Colorado										
\$5,040	\$4,840	\$4,500	\$4,131	\$3,750	\$3,047										

Need a more thorough analysis of what isn't being captured by current web archiving and why

- Insights into the technical challenges
- Test a spectrum of news apps, including projects using external APIs
- Consider if a new or updated web archiving tool could be developed
- Develop a registry of news apps or produce report on the scan

More Challenges



Conflicting incentives (legal, financial)

- Proprietary code and digital rights management issues
- News producers want to monetize their archives

Legal Strategies

- Form partnerships with non-profits and others willing to draft agreements
- Push for more open digital rights management

Landscape, Process and Resources

- Poor sense of the amount and type of news apps being created
- No automated workflows in place to capture or archive content
- Human and financial resources for digital archiving and preservation are scarce

Saving News Apps



Determine the “significant properties” of news apps

- Will affect what archiving and preservation solutions to pursue
- May vary from news app to news app
- Considerations for: content, context, rendering, structure, and behavior of the news app

- Must capture the database, the interface, data visualizations, outputs and analysis
- “Look and Feel” of the app as well as the content and functionality

The background image shows a well-organized archive or library. It features blue metal shelving units filled with numerous yellow cardboard boxes and manila-colored folders. Many of the boxes and folders have white labels with handwritten text and numbers. Some labels include names like 'DANTON', 'SLATERHOUSE', and 'MAYOR'. Numbers such as '1939', '1946', '1932', '1945', '1946', '1947', '1948', '1949', '1950', '1951', '1952', '1953', '1954', '1955', '1956', '1957', '1958', '1959', '1960', '1961', '1962', '1963', '1964', '1965', '1966', '1967', '1968', '1969', '1970', '1971', '1972', '1973', '1974', '1975', '1976', '1977', '1978', '1979', '1980', '1981', '1982', '1983', '1984', '1985', '1986', '1987', '1988', '1989', '1990', '1991', '1992', '1993', '1994', '1995', '1996', '1997', '1998', '1999', '2000' are visible. The overall scene is one of a vast, organized collection of historical or archival documents.

Determine a technique for capture and archiving

- Will depend on the “significant properties,” and what level of human and financial resources can be devoted to the process
- Possible techniques include: technical preservation (computer museums, not very tenable), migration (better for static, not dynamic, digital objects), and emulation
- Digital preservation community currently favors emulation – need consensus on which emulator to use long-term, can migrate the emulator



Adopt a framework for archiving and preservation

- Migration or web archiving framework
- Emulation framework: Performance Model Framework for the Preservation of a Software System (Matthews et al, 2010)

Adopt a metadata description schema for best interoperability, discovery and access to news apps

- If web archiving can be a solution, we can use and update the current XML metadata schema
- If emulation must be used, there are several options, including PREMIS from the Library of Congress

Mapping the framework to a metadata schema

Framework Category (Matthews, et al., 2010)	Framework Description (Matthews, et al., 2010)	Examples (Matthews, et al., 2010)	Equivalent OAIS Terms	Equivalent PREMIS v.3 categories (select examples)	News Apps Registry Fields
Functionality	-Description of the typical characteristics of software. -Useful for efficient discovery and accessibility of the software in future	-Description of inputs and outputs -Description of operation and algorithms -Description of the domain addressed	-Descriptive Information	-Intellectual Entities or Object Entities (significantProperties, objectIdentifier, etc)	-Descriptive information on the application for discoverability: the name of the application and a description of its purpose -URL of the news app
Software Composition	-Description of the components that constitute software -Useful for rebuilding and reusing the software in future -Detailed history of version changes and other significant changes that a software product has undergone facilitates verification of its authenticity	-A typical record: binary files, source code, user manuals and tutorials. -A more complete record: requirements and design documentation, test cases and harnesses, prototypes, formal proofs.	-Representation Information -Preservation Description Information (PDI)	-Intellectual Entities or Object Entities (environmentFunction, environmentDesignation, etc)	-Description of the components that constitute the software -Link to the open source software repository, if available (Github ¹ , SourceForge ² , Launchpad ³ , etc)
Provenance and Ownership	-Different software components have different and complex licensing conditions. -Needs to be included in the preservation planning	-Software owner and licence information, e.g. Microsoft for MS Word®	-Provenance Information category of Preservation Description Information (PDI)	-Rights Entities (rightsStatement, licenseInformation, etc.) -Agent Entities (agentIdentifier, etc.)	-Rights statements, software owner and license information, the news application development team member names and roles
User Interaction	-Description of expected mode of interaction between user and software -The 'Look and Feel' and the model of user interaction can play a	-The inputs which a user enters through a keyboard, pointing device or other input devices, such as	-Not comprehensively addressed in the OAIS – may be categorized as the Significant	-Event Entities (eventDetail, eventOutcome, etc)	-User inputs and application outputs

- Framework for Preservation of a Software System => PREMIS v. 3
- Includes categories for: Functionality, Software composition, Provenance and ownership, User interaction, Software environment, Software architecture, and Operating performance

¹ <https://github.com/>

² <https://sourceforge.net/>

³ <https://launchpad.net/>

Needed Work

Saving News Apps is a group effort. We need:

- Conversations with media organizations about preserving their digital history and legacy
- More digital archivists and web developers
- More grant funding and institutional resources
- Case studies to test and update these recommendations
- Advocacy for more open digital rights management
- Establishment of best practices, workflows

Thank you! Questions?

Katherine.Boss@nyu.edu

MerBroussard@nyu.edu