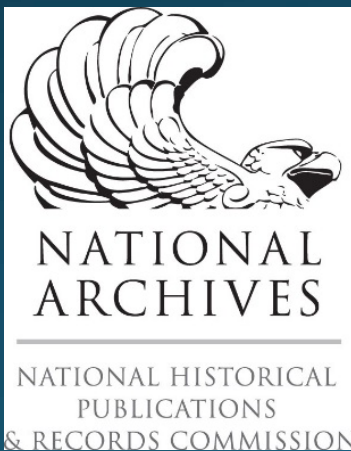


Managing and Accessing Archival Email

The TOMES Project

Camille Tyndall Watson
Digital Services Section Head
State Archives of NC

Jeremy Gibson
Systems Integration Librarian
State Archives of NC



State Archives of North Carolina
NATURAL AND CULTURAL RESOURCES

Your Program for Today

- 1) History of SANC's email program
- 2) TOMES (Capstone)
- 3) TOMES (Tool)



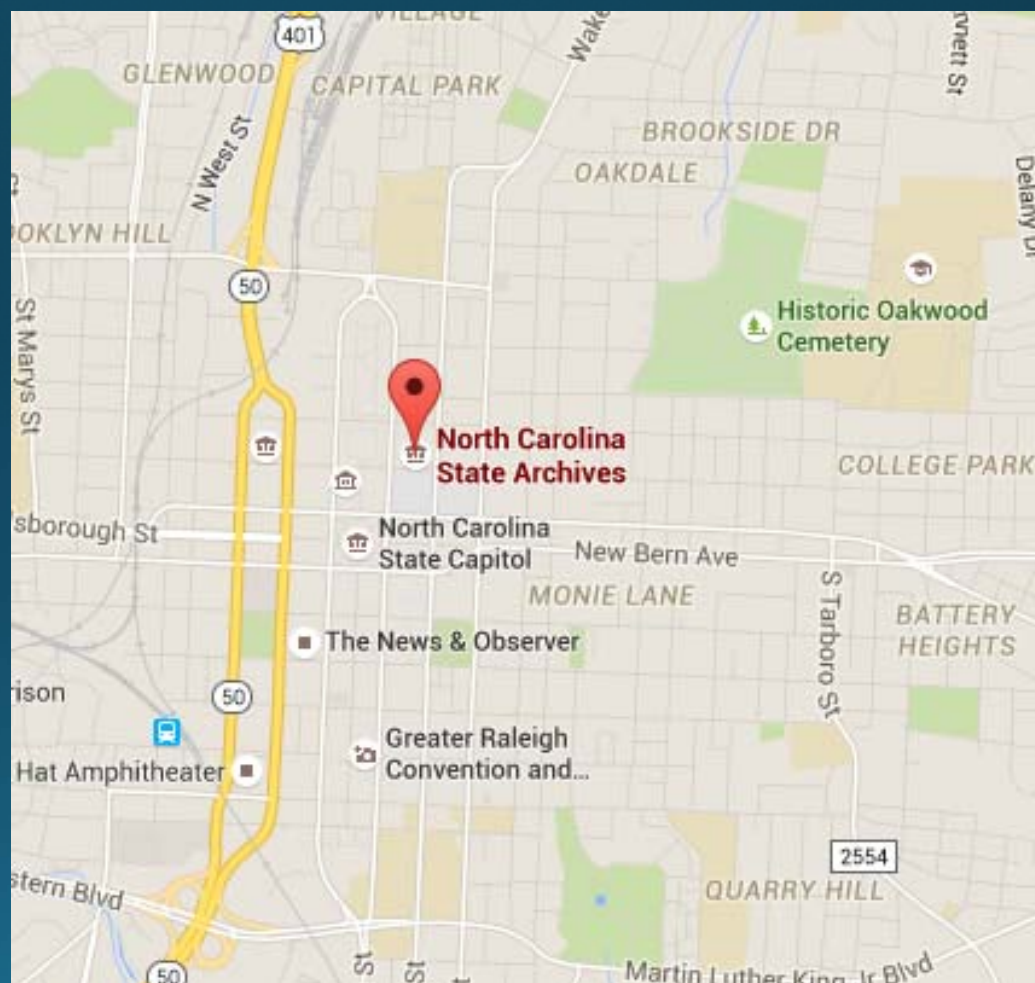
Where We've Been...

- NC processes Gov. Jim Hunt emails, 2001-2002
- EMCAP, 2007-2009
- Executive Orders 150 and 18, 2009
- Mimosa Nearpoint, 2009-2014



Where We Are...

- Email retention reduced to 5 years
- State IT has become centralized
- Email has moved to the cloud with Office 365
- Email archive migration



Transforming Online Mail with Embedded Semantics (TOMES)

- 3 year grant (2015-2018)
- Partnership between State Archives of NC, Utah State Archives, and Kansas State Historical Society
- Advisory group includes Cal Lee (UNC-Chapel Hill), Chris Prom (University of Illinois Urbana Champaign), and staff from the Library of VA



What should we keep?

Capstone Approach

Permanent - top decision-maker's e-mail

Temporary – all other staff, 7 year retention

Non-record – 1 year retention

Building Relationships

- Collaboration with Government Records Section, SANC
- Development of a series of forms to understand organizational structures of state agencies
- Outreach to State Agency CROs and CIOs for education and buy in



Collecting Data



archives.ncdcr.gov

DIVISION OF ARCHIVES AND RECORDS
GOVERNMENT RECORDS SECTION



archives.ncdcr.gov

4615 Mail Service Center

DIVISION OF ARCHIVES AND RECORDS



archives.ncdcr.gov

4615 Mail Service Center, Raleigh NC 27699-4165

919-807-7350

DIVISION OF ARCHIVES AND RECORDS
GOVERNMENT RECORDS SECTION

Agency: _____

Emails from the individuals in this category are high-level, senior-level, or otherwise significant.

For each individual identified below, please list the email accounts used by that individual.

- The head of the agency, the president, or the equivalent.** Although the one position may be held by more than one individual, list only one individual.

NAME

Predecessors (from the list above)

NAME

ASSESSMENT OF EMAIL ACCOUNTS

Agency: _____

Please list below the individuals in the suggested roles and positions of your agency. Core functions or programs that are essential to the mission statements and serve significant populations and/or groups. These individuals conduct research supporting core programs, oversee outreach programs, capture the history of core programs, or are responsible for these email accounts are not already listed under categories 1 through 4.

If the emails generated by individuals listed below are already being captured by the Part 1 form you already completed, including instances where necessary to list those individuals here. Your agency may not list individuals who are already listed on the Part 1 form.

For each individual identified below, please list the email accounts used by that individual.

- Staff assistants to heads of agencies and their deputies.** Important work is often carried out by special assistants to senior officials and/or their email account contacts. List the email accounts of senior officials they support.

NAME	POSITION TITLE/ROLE	BEACON POSITION NUMBER

ASSESSMENT OF EMAILS FOR PERMANENT RETENTION PART 3

Agency: _____

Please identify positions that create official records that document agency policies and decisions related to any of the programs or subjects listed below. These programs or subjects are not represented in the accounts of personnel already listed on Parts 1 and 2 but still may need to be retained permanently. Please note: if emails that document the programs or subjects listed below are already being captured by the email account of an agency position listed on the Parts 1 or 2 forms, including instances where the agency executive has been copied on these emails, it is not necessary to list those positions here. For each individual identified below, please list the email accounts of all predecessors in that role since January 2011.

Possible email subjects with archival value:

- Major agency policies
- Formulation of rules and monitoring standards (e.g., Administrative Code)
- Events, incidents, and situations that required a prolonged response involving multiple agencies and led or had the potential to lead to large-scale loss of life, severe damage to lands and property, or major disruption of the state's infrastructure
- Direction and planning of the core program(s) of your agency
- Cooperation with external state and/or federal agencies
- Construction and real property transactions
- Major public events, such as the State Fair, First Flight Centennial, inaugurations, etc.
- History of the state of North Carolina
- Advocacy for minorities, such as Indian tribes
- Certification, commissioning, etc.
- Management of assets held in public trust for the people of North Carolina – e.g. state parks, historic sites, artifacts, archival materials, etc.
- Evaluation of rules created by other agencies, where the agency is an established part of the rule review process

NAME	POSITION TITLE/ROLE	BEACON POSITION NUMBER	ARCHIVAL SUBJECT MATTER (insert number from list above)	EMAIL ADDRESS	BEGINNING DATE FOR EMAIL COLLECTION

Plugging In

- Department of Information Technology
 - “Tagging” accounts by function
 - Facilitating the transfer of email accounts from cloud storage
- Office of State Human Resources
 - Identifying positions by position number
 - Working with DIT to “tag” accounts



Challenges

- Capstone and Account identification
- Development of appraisal criteria
- Getting buy-in and support from agency CROs and CIOs
- Outreach and education
- Communicating needs to IT
- Figuring out who needs to be at the table



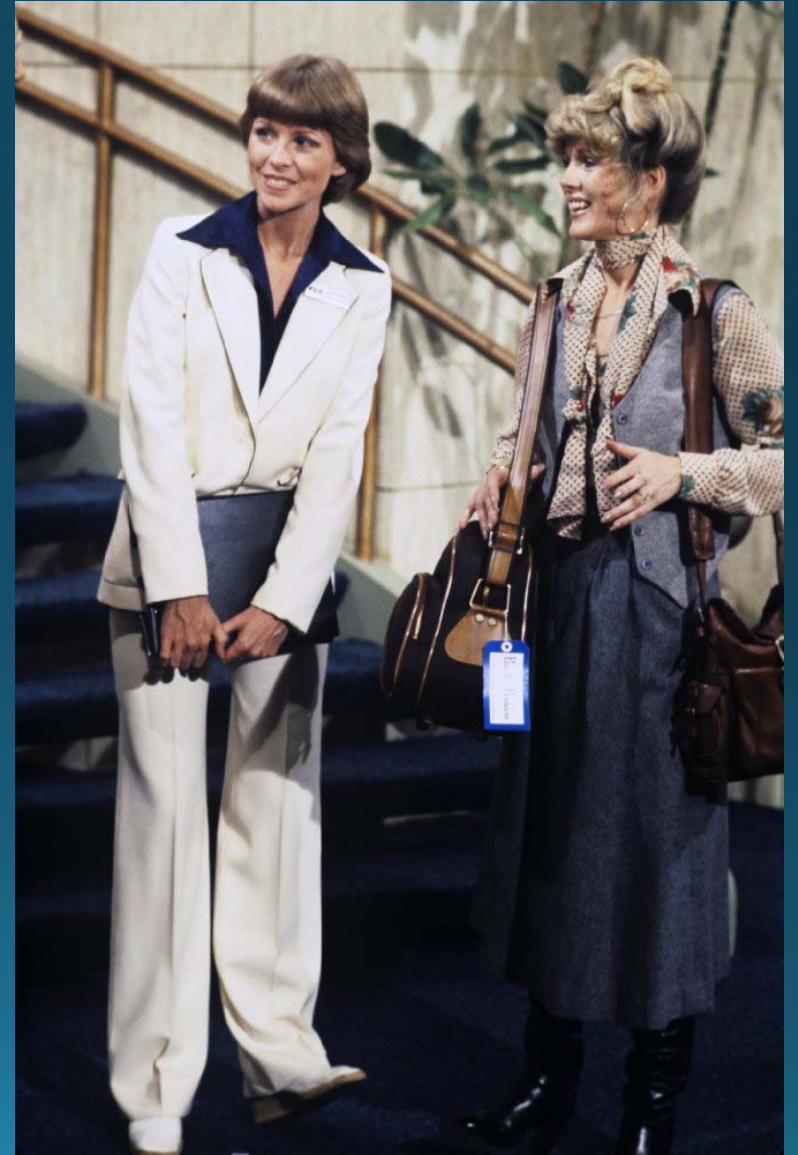
Making It Work

- Harvesting email from cloud using Office365 eDiscovery tools
- Created email processing modules written in Python as a toolbox for preserving and processing email accounts
- Developing NLP libraries for tagging



TOMES (Tool)

- 1) Architecture/Design
- 2) Original Goals
- 3) Modified Goals
- 4) V1.0



Architecture/Design

- KISS (Keep It Simple Stupid)
- Microservices

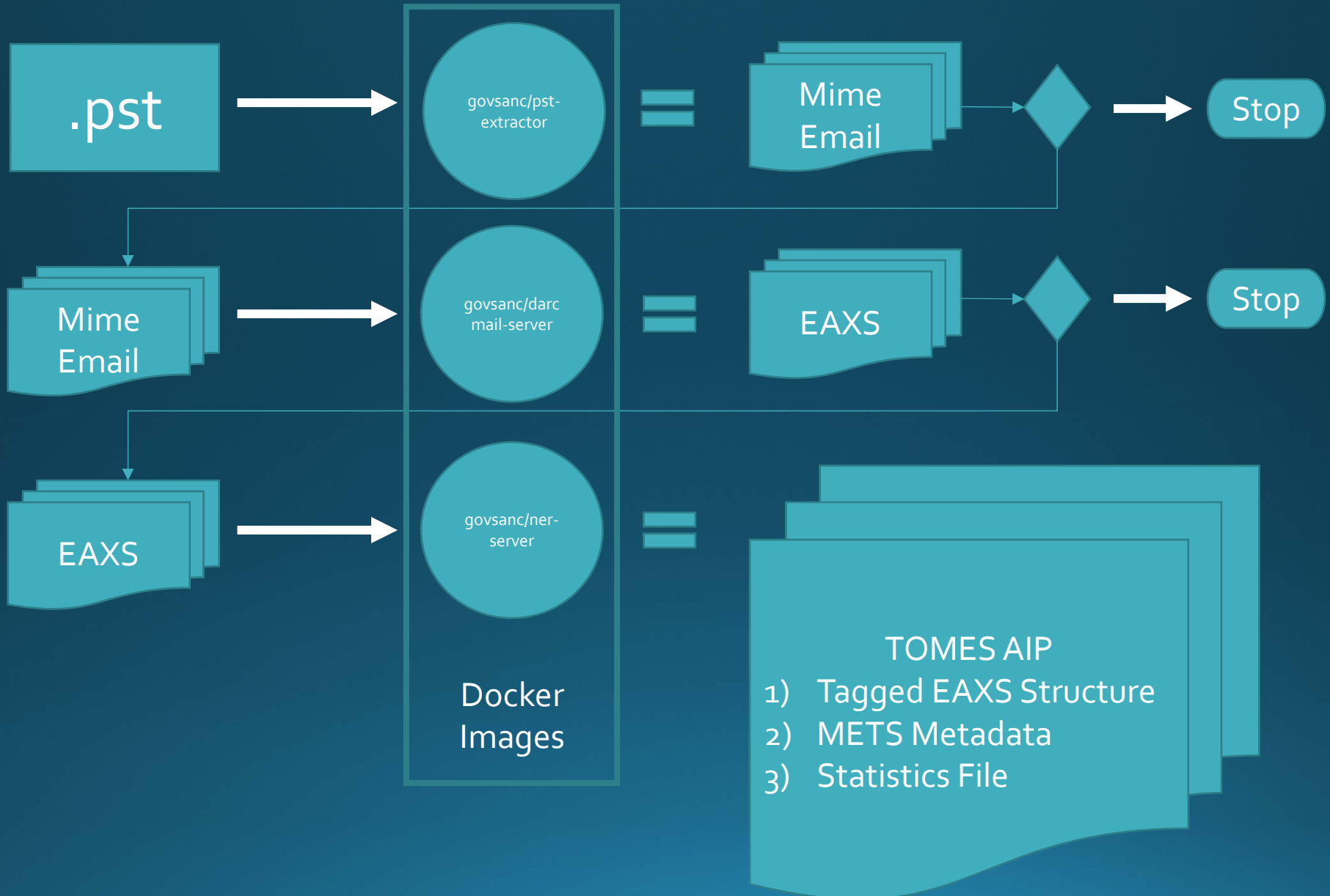


Architecture: KISS

- Use the best tool for the job
- Don't reinvent the wheel
 - Readpst
 - DarcMail
 - CoreNLP



Architecture: Microservices



Original Goals

- Move PSTs from black box to preservation format
- Add semantics to the text to aid in processing/access
 - Semantics customizable by institution
- Run on hardware/software available to under resourced institutions.



Modified Goals

- Allow for iterative processing
- Greater focus on the PII part of the semantics
- Make the package self describing and atomic



TOMES v1.0

ITS_TICKET Apps Teamwork Microsoft Office Home APPX - straycats.appx APPX - ncc433.its.stat NC Archives & Record Golf Packages Project Center Other bookmarks

Main Convert a PST Convert to EAXS Tag EAXS Documentation

TOMES (Transforming Online Mail with Embedded Semantics)

Welcome to the TOMES application. From here you can:

- Convert PSTs to MBOX or EML: We recommend EML at the moment
- Convert MBOXs or EMLs to EAXS (an XML representation of a Mailbox)
- Tag the EAXS with Named Entities
- Search the Tagged Mailbox

You are not obligated to do all of these steps you. The TOMES app will output a fully documented AIP at every step.

Technologies in TOMES

The TOMES app utilizes multiple technologies. Below are a listing of the main software components:

Docker	Docker is used to link all of the apps together and provide a simple setup and execution across platforms.
PST to Mime	The PST to Mime part of the tool utilizes an Afterlogic's MailBee .NET Objects library to convert PSTs to Mime email. Currently, those mime emails are .EMLs. MailBee is the only closed source part of this project. The project has purchased a license for the project and binaries may be distributed.
DarcMail	TOMES utilizes a fork of CmdDarcMailXML written by Carl Schaefer of the Smithsonian Institution. DarcMail converts a MBOX or an EML structure into EAXS xml.
Stanford Core NLP	Messages are tagged using Stanford's Core NLP engine. TOMES provides a method for using a tagging dictionary which allows for States and organizations to customize the tags for semantics relevant to their organizational context.

© Copyright 2017 by Tomes Project.

TOMES v1.0

Main

Convert a PST

Convert to EAXS

Tag EAXS

Documentation

Select a PST to convert to mime email.

- [-] PSTS
 - ZachAmbrose.pst
 - MikeWard.pst
 - test.pst
 - OOG_Perdue.pst

Output Mime-type folder:

MikeWard_2018_2_27

Convert

Client Connected
MikeWard is selected for conversion

TOMES v1.0

Main

Convert a PST

Convert to EAXS

Tag EAXS

Documentation

Select a directory containing Mime Emails to transform.

- Mime Sources
 - State Treasurer
 - OOG_Perdue_2018_2_14
 - MikeWard_2018_2_15
 - Dept. of Technology Services
 - Dept. of Admin Services
 - DPatton
 - Apple Valley
 - ZachAmbrose_2018_2_27
 - test_2018_2_27
 - BadEmls

Transfer Name:

State Treasurer

Chunk output? 1000

Stitch chunks?

Is source .EML?

Convert

Client connected!

TOMES v1.0

Main

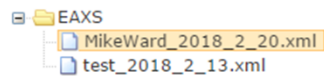
Convert a PST

Convert to EAXS

Tag EAXS

Documentation

EAXS file to transform:



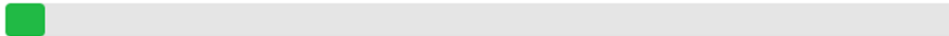
EAXS Selected

MikeWard_2018_2_20.xml

Tagged Name

MikeWard_2018_2_20_tagged.xml

Convert



```
Client Connected
Locating files...
Processing: /home/tomes/data/eaxs/MikeWard_2018_2_20/eaxs_xml/MikeWard_2018_2_20.xml
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Unzipping tokenizers/punkt.zip.
2018-02-20 20:10:18.257 - __main__ - INFO - Tagging EAXS file: /home/tomes/data/eaxs/MikeWard_2018_2_20/eaxs_xml/MikeWard_2018_2_20.xml
2018-02-20 20:10:18.257 - lib.eaxs_to_tagged - INFO - Converting '/home/tomes/data/eaxs/MikeWard_2018_2_20/eaxs_xml/MikeWard_2018_2_20.xml' EAXS file to tagged EAXS file:
/home/tomes/data/eaxs/MikeWard_2018_2_20/eaxs_xml/MikeWard_2018_2_20_tagged.xml
2018-02-20 20:10:18.260 - lib.eaxs_to_tagged - INFO - Finding number of Messages.
2018-02-20 20:10:22.756 - lib.eaxs_to_tagged - INFO - Messages Remaining: 1711
2018-02-20 20:10:22.756 - lib.eaxs_to_tagged - INFO - Processing LocalId: 0000001
2018-02-20 20:10:44.778 - lib.eaxs_to_tagged - INFO - Messages Remaining: 1710
2018-02-20 20:10:44.779 - lib.eaxs_to_tagged - INFO - Processing LocalId: 0000002
2018-02-20 20:10:44.962 - lib.eaxs_to_tagged - INFO - Messages Remaining: 1709
2018-02-20 20:10:44.962 - lib.eaxs_to_tagged - INFO - Processing LocalId: 0000003
2018-02-20 20:10:45.259 - lib.eaxs_to_tagged - INFO - Messages Remaining: 1708
2018-02-20 20:10:45.260 - lib.eaxs_to_tagged - INFO - Processing LocalId: 0000004
```

V1.0: The Final AIP

- Original Account file (.pst, .mbox, .eml)
- Untagged EAXS XML file
- Attachment XML file(s)
- Tagged EAXS XML file
- METs file

```
1 <mets:mets xmlns:mets="http://www.loc.gov/METS/" xmlns:xlink="http://www.w3.org/1999/xlink">
2   <mets:metsHdr>
3     <mets:agent ROLE="CREATOR" TYPE="OTHER" OTHERTYPE="Software Agent">
4       <SingleBody>
5         <xref href="2b29537a-fe00-4d9b-ac26-b8ca124e5d4c.xml" type="attachment"/>
6         <xref href="00010_698bcf9b-fdb2-43b5-879f-c4a72be1c67a.xml" type="attachment"/>
7         <xref href="00012_bbcf2b53-150c-4eb2-8010-cd8897fd3316.xml" type="attachment"/>
8         <xref href="00022_bbcf2b53-150c-4eb2-8010-cd8897fd3316.xml" type="attachment"/>
9       </SingleBody>
10    </mets:agent>
11  </mets:metsHdr>
12  <mets:structMap>
13    <!--This is a comment.-->
14  </mets:structMap>
15  </mets:mets>
```

Name	Date modified	Type	Size
2b29537a-fe00-4d9b-ac26-b8ca124e5d4c.xml	2/20/2018 1:23 PM	XML File	
00010_698bcf9b-fdb2-43b5-879f-c4a72be1c67a.xml	2/20/2018 1:22 PM	XML File	
00012_bbcf2b53-150c-4eb2-8010-cd8897fd3316.xml	2/20/2018 1:22 PM	XML File	
00022_bbcf2b53-150c-4eb2-8010-cd8897fd3316.xml	2/20/2018 1:22 PM	XML File	
error.log	2/20/2018 1:23 PM	Text Document	1 KB
info.log	2/20/2018 1:23 PM	Text Document	220 KB
MikeWard_2018_2_20.xml	2/20/2018 1:23 PM	XML File	77,749 KB
MikeWard_2018_2_20_tagged.xml	2/20/2018 3:18 PM	XML File	32,484 KB
00059_f6ea0834-e324-4505-b244-69bf860fd06c.xml	2/20/2018 1:22 PM	XML File	

Challenges

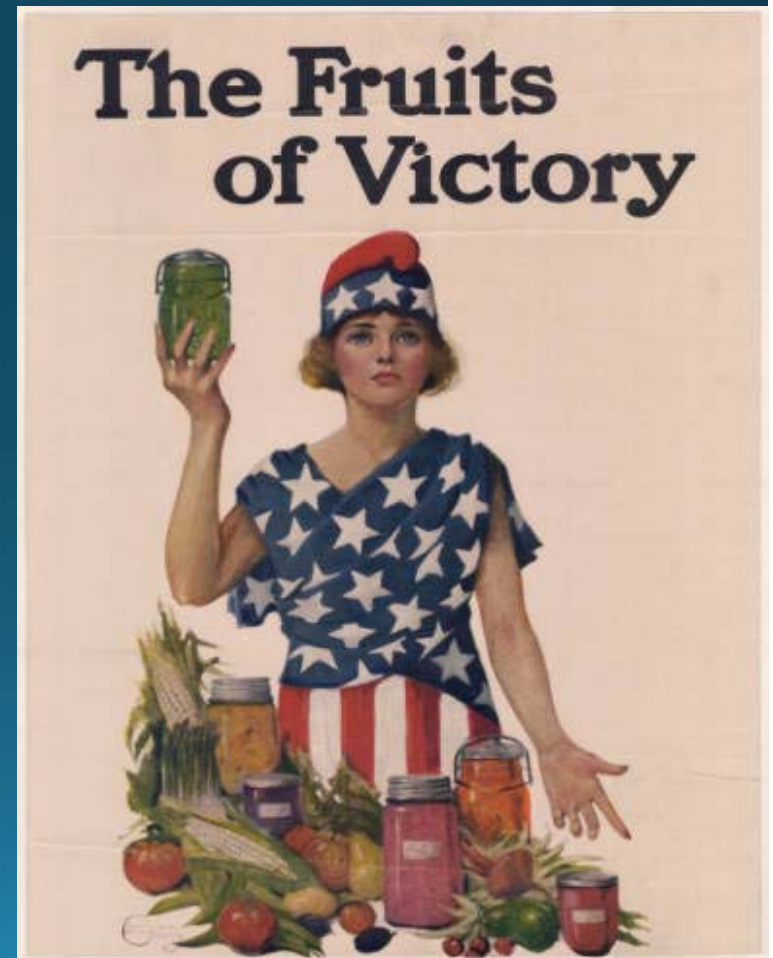
- Making imperfect NER useful.
 - Machine Learning
 - Active discovery and processing
- Making State specific libraries easier for non-technical users to develop and incorporate into the workflow.
- Handling of emails with bad encodings
 - Email is messy and comes from everywhere.



SOME PARTY

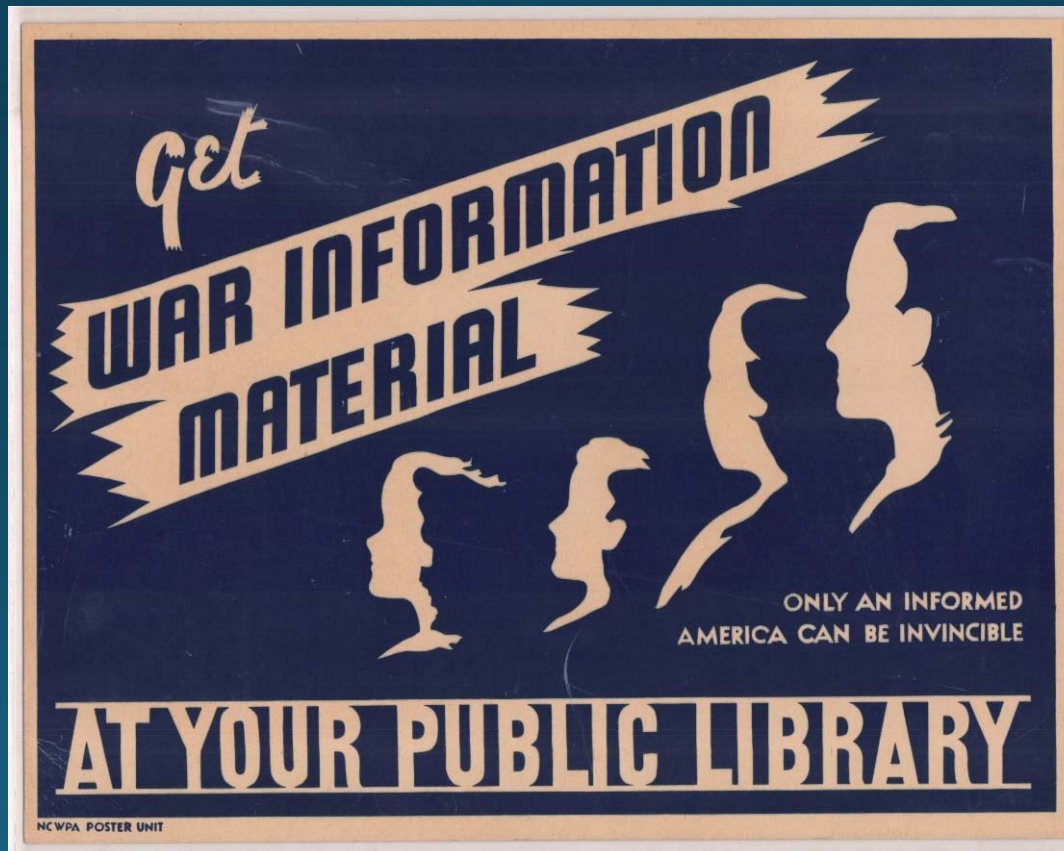
The Final Product

- Development of State Government specific NLP libraries for use in the processing of email accounts containing public records
- An MPLP approach to the arrangement and description of email
- The ability to identify materials that should be reviewed for PII before release to public
- Mediated access using iterative processing



Stay in touch!

- GitHub: <https://github.com/StateArchivesOfNorthCarolina>
- Website: <http://www.ncdcr.gov/tomes>



Questions?

Camille Tyndall Watson

Digital Services Section Head

State Archives of NC

Camille.TyndallWatson@ncdcr.gov

(919) 807 – 7359

Jeremy Gibson

Systems Integration Librarian

State Archives of NC

Jeremy.Gibson@ncdcr.gov

(919) 807-7356



State Archives of North Carolina
NATURAL AND CULTURAL RESOURCES