Expanding the Scale of PDF Preservation
to Accommodate a State Press Association through the
Texas Digital Newspaper Program

Ana Krahmer, Digital Newspaper Program
Coordinator

ana.krahmer@unt.edu

# Overview

- What is TDNP?
- Initial PDF Newspaper Projects
- The Texas Press Association Archive
- Technology and Standards

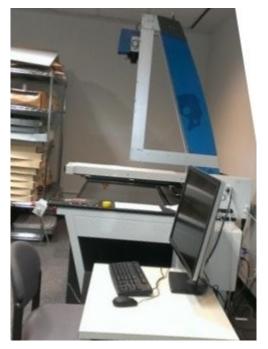# What is TDNP?

About us

ana.krahmer@unt.edu

# What is TDNP?

- Dedicated to preserving Texas newspapers, from any time or place, for any title.

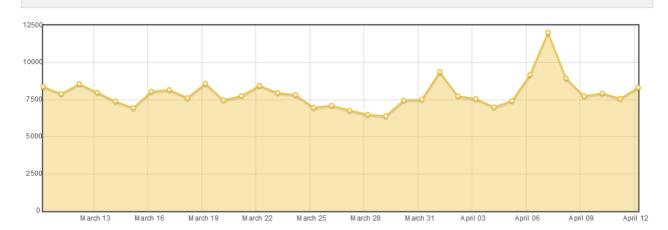- Thus far, we host nearly 3 million pages of newspapers, dating from 1829 to present.

ana.krahmer@unt.edu

About this Collection    Explore this Collection    Explore all Collections

## Statistics for Texas Digital Newspaper Program

Item Usage | Added Items | More Data

**8,681,963** Total Uses / **286,311** Total Items (2,822,790 files) / **280,026** Visible / **6,285** Hidden



### Usage by Month/Year

| Year | January | February | March | April | May | June | July | August | September | October | November | December | Total |
|------|---------|----------|-------|-------|-----|------|------|--------|-----------|---------|----------|----------|-------|
| 2015 | 220,271 | 241,844 | 238,711 | 100,847 | | | | | | | | | 801,673 |
| 2014 | 178,665 | 165,850 | 179,184 | 164,216 | 227,182 | 237,335 | 463,178 | 172,892 | 193,984 | 317,866 | 174,249 | 191,622 | 2,666,223 |
| 2013 | 236,016 | 246,394 | 260,614 | 239,745 | 214,253 | 211,501 | 230,605 | 260,006 | 230,506 | 164,797 | 174,215 | 157,953 | 2,626,605 |
| 2012 | 109,225 | 106,035 | 117,329 | 124,214 | 112,802 | 113,189 | 139,030 | 130,802 | 136,328 | 164,379 | 201,137 | 190,694 | 1,645,164 |
| 2011 | 36,983 | 41,626 | 34,802 | 36,070 | 36,739 | 35,884 | 47,272 | 47,164 | 58,449 | 67,582 | 75,419 | 86,971 | 604,961 |

About this Newspaper    Read this Newspaper    Other items in this serial (204)

The Houston Post. (Houston, Tex.), Vol. 19, No. 270, Ed. 1 Thursday, December 31, 1903

Brief Record | Full Record | Statistics



| | |
|---|---|
| **Description:** | Daily newspaper from Houston, Texas that includes local, state and national news along with extensive advertising. |
| **Creator(s):** | Unknown |
| **Location(s):** | United States - Texas - Harris County - Houston |
| **Creation Date:** | December 31, 1903 |
| **Partner(s):** | UNT Libraries<br>About \| Browse this Partner |
| **Collection(s):** | Texas Digital Newspaper Program<br>About \| Browse this Collection<br><br>Houston Daily Post<br>About \| Browse this Collection |
| **Usage:** | *Total Uses:* 22<br>*Past 30 days:* 22<br>*Yesterday:* 1 |

ana.krahmer@unt.edu

# What is TDNP?

- Dedicated to preserving Texas newspapers, from any time or place, for any title.

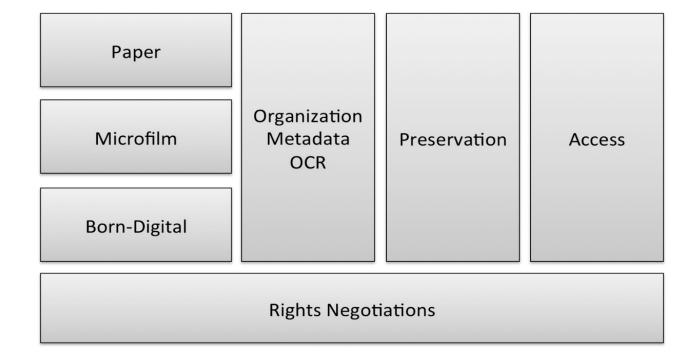- Thus far, we host nearly 3 million pages of newspapers, dating from 1829 to present.



ana.krahmer@unt.edu

# Workflow

| Paper | | | |
|---|---|---|---|
| Microfilm | Organization Metadata OCR | Preservation | Access |
| Born-Digital | | | |

Rights Negotiations

# Initial PDF Newspaper Projects

Preservation on The Portal to Texas History

ana.krahmer@unt.edu

# Initial PDF Newspaper Projects

- Began working with PDFs in 2010.
- Preserved PDF issues after receiving grants to digitize earlier, analog issues.
- Earliest PDF issue is from 18 March 1998 (University of Dallas).
- PDFs were acquired from publishers.
- Permissions granted by publishers.

# Initial PDF Newspaper Projects



ana.krahmer@unt.edu

# Initial PDF Newspaper Projects

- Flexibility of Portal to embargo gives publishers confidence.

- Example: *Cherokeean Herald*

- When publishers see one successfully-preserved title, they show interest.

# Texas Press Association Archive

Preserving Recent Texas History

ana.krahmer@unt.edu

# TPA Archive Partnership

- Collaboration with the Texas Press Association and NewzGroup out of Missouri.

- 12TB of PDF newspapers, prepared for preservation through batch processing of PDF content.

- Range from 2010-August 2014.

- Scheduled transfer of 2014-2015 newspapers from NewzGroup.

# TPA Archive Partnership: File Workflow

- QC work performed after batch processing.
- Add initial layer of metadata, pre-OCR.

# TPA Archive Partnership: File Workflow

# TPA Archive Partnership: File Workflow



| | | | |
|---|---|---|---|
| 91582_SundayJune172012_01.abbyy.xml | 3/25/2015 5:35 PM | XML File | 529 KB |
| 91582_SundayJune172012_01.jpg | 2/2/2015 10:03 AM | ACDSee 9.0 JPEG I... | 5,452 KB |
| 91582_SundayJune172012_01.txt | 3/25/2015 5:35 PM | Text Document | 8 KB |
| 91582_SundayJune172012_02.abbyy.xml | 3/25/2015 5:35 PM | XML File | 880 KB |
| 91582_SundayJune172012_02.jpg | 2/2/2015 10:03 AM | ACDSee 9.0 JPEG I... | 6,451 KB |
| 91582_SundayJune172012_02.txt | 3/25/2015 5:35 PM | Text Document | 14 KB |
| 91582_SundayJune172012_03.abbyy.xml | 3/25/2015 5:35 PM | XML File | 453 KB |
| 91582_SundayJune172012_03.jpg | 2/2/2015 10:03 AM | ACDSee 9.0 JPEG I... | 4,744 KB |
| 91582_SundayJune172012_03.txt | 3/25/2015 5:35 PM | Text Document | 7 KB |
| 91582_SundayJune172012_04.abbyy.xml | 3/25/2015 5:35 PM | XML File | 570 KB |
| 91582_SundayJune172012_04.jpg | 2/2/2015 10:03 AM | ACDSee 9.0 JPEG I... | 5,798 KB |
| 91582_SundayJune172012_04.txt | 3/25/2015 5:35 PM | Text Document | 9 KB |
| 91582_SundayJune172012_05.abbyy.xml | 3/25/2015 5:35 PM | XML File | 672 KB |
| 91582_SundayJune172012_05.jpg | 2/2/2015 10:03 AM | ACDSee 9.0 JPEG I... | 5,774 KB |
| 91582_SundayJune172012_05.txt | 3/25/2015 5:35 PM | Text Document | 11 KB |

ana.krahmer@unt.edu

# TPA Archive Partnership: File Workflow

```xml
<?xml version="1.0" encoding="UTF-8"?>
- <metadata>
      <title qualifier="officialtitle">The Monitor (Mabank, Tex.)</title>
      <title qualifier="serialtitle">The Monitor</title>
  - <creator qualifier="edt">
        <info>Managing Editor</info>
        <type>per</type>
        <name>Cantrell, Pearl</name>
    </creator>
  - <contributor>
        <info/>
        <type/>
        <name/>
    </contributor>
  - <publisher>
        <info>MediaOne, L. L. C.</info>
        <location>Mabank, Texas</location>
        <name>MediaOne, Limited Liability Company</name>
    </publisher>
      <language>eng</language>
      <description qualifier="content">Semi-weekly newspaper from Mabank, Texas that includes local Cedar Creek Lake area, state a
      <description qualifier="physical">pages : ill.</description>
      <subject qualifier="UNTL-BS">Business, Economics and Finance - Communications - Newspapers</subject>
      <subject qualifier="UNTL-BS">Business, Economics and Finance - Journalism</subject>
      <subject qualifier="UNTL-BS">Business, Economics and Finance - Advertising</subject>
      <subject qualifier="UNTL-BS">Places - United States - Texas - Kaufman County</subject>
      <subject qualifier="UNTL-BS">Places - United States - Texas - Henderson County</subject>
      <subject qualifier="LCSH">Kaufman County (Tex.) -- Newspapers.</subject>
      <subject qualifier="LCSH">Henderson County (Tex.) -- Newspapers.</subject>
      <subject qualifier="LCSH">Mabank (Tex.) -- Newspapers.</subject>
      <subject qualifier="LCSH">Gun Barrel City (Tex.) -- Newspapers.</subject>
      <subject qualifier="LCSH">Kemp (Tex.) -- Newspapers.</subject>
      <subject qualifier="KWD">Seven Points</subject>
      <subject qualifier="KWD">Tool</subject>
      <subject qualifier="LCSH">Eustace (Tex.) -- Newspapers.</subject>
      <subject qualifier="KWD">Payne Springs</subject>
      <subject qualifier="KWD">Log Cabin</subject>
      <subject qualifier="KWD">Enchanted Oaks</subject>
      <subject qualifier="LCSH">Trinidad (Tex.) -- Newspapers.</subject>
      <subject qualifier="LCSH">Malakoff (Tex.) -- Newspapers.</subject>
      <primarySource>1</primarySource>
      <coverage qualifier="placeName">United States - Texas - Kaufman County - Mabank</coverage>
      <coverage qualifier="timePeriod">mod-tim</coverage>
      <source qualifier=""/>
      <relation qualifier=""/>
      <collection>TDNP</collection>
      <collection>CCLM</collection>
      <institution>UNT</institution>
      <rights qualifier="license">copyright</rights>
      <resourceType>text_newspaper</resourceType>
      <format>text</format>
      <identifier qualifier="LCCN">sn88083774</identifier>
      <identifier qualifier="OCLC">18109674</identifier>
      <degree qualifier=""/>
      <note qualifier="display">Masthead reads, "COVERING THE ENTIRE CEDAR CREEK LAKE AREA."</note>
      <note qualifier="display">Published on Sundays and Thursdays.</note>
      <meta qualifier="hidden">False</meta>
      <meta qualifier="metadataCreator">tgieringer</meta>
  </metadata>
```

- Example of batch metadata for PDF issues
- Applied in XML file to sets divided by year and/or by content changes (managing editor, publisher, masthead, etc.)
- Newspapers uploaded with embargo are marked as "hidden"=TRUE.

ana.krahmer@unt.edu

# TPA Archive Partnership: Permissions

- Publishers are busy, hard people to catch.
- When they respond, they respond with interest, with the exception of large-city dailies.
- Embargos have ranged from the most recent 6 months to 3 years.
- The Texas Digital Newspaper Program holds a membership in the Texas Press Association.
- Krahmer attends all TPA annual conventions and summer leadership meetings, along with as many regional Press Association meetings as possible.
- Preservation is about establishing communication, trust.

ana.krahmer@unt.edu

# Collaboration with Publishers

- We have the capability to open or hide issues at publisher's request.
- We can unhide issues when the embargo period expires.

## Statistics for Cherokeean Herald

| Item Usage | More Data |
| --- | --- |

### Totals

| | |
| --- | --- |
| Total Number of Items | 4,464 |
| Total Number of Visible Items | 4,360 |
| Total Number of Hidden Items | 104 |
| Total Number of Files | 66,167 |
| Total Number of Uses | 244,400 |

Visible & Invisible Items

Visible — — Invisible

■ 4360  ■ 104

ana.krahmer@unt.edu

THE PORTAL TO TEXAS HISTORY

*Cherokeean Herald (Rusk, Tex.), Vol. 160, No. 46, Ed. 1 Wednesday, January 6, 2010*

Normal Display ▶ Sequence: 1 Permalink Print View



ana.krahmer@unt.edu

# Technology and Standards

File types, software, and metadata

ana.krahmer@unt.edu

# Filetypes

- The PDF print master is the preservation copy.
- Save this into JPG format at 400 dpi, from which derivatives are created.

# Software

- Adobe Acrobat, batch-scripting capabilities
- Batch renaming application
- Python scripts
- Microservices

ana.krahmer@unt.edu

# A file is a file is a file is a. . .

- After conversion, whether from analog to digital, or from one born-digital filetype to another, processing is very similar.

- Regardless of original format, issues are named according to yyyymmdded: 1901052301= the issue for May 23, 1901, edition 01, of a newspaper.

- Metadata, with minor differences ("physical description" & bagit information), is the same for all newspapers.

# Metadata

- Minor bag-info differences: bag-info files (BagIt) for pdfs contain the following information (red text is unique to PDF materials).

  Source-Organization: University of North Texas Libraries

  Organization-Address: P. O. Box 305190, Denton, TX 76203-5190

  Contact-Name: Mark Phillips

  Contact-Phone: 940-565-2415

  Contact-Email: mark.phillips@unt.edu

  External-Description: Newspaper issues of the "NEWSPAPER NAME HERE" published in

  [ CITY], Texas. Issues were made available from <span style="color:red">born-digital</span>

  <span style="color:red">PDF printmasters.</span> Partner institution is the [partner library here].  <span style="color:red">Master files were PDF printmasters from which derivative JPGs were created.</span>

# Questions?

Email: Ana.Krahmer@unt.edu
Call: 940-565-3367
Visit: http://texashistory.unt.edu/explore/collections/TDNP/